# Applied Machine Learning
# Part II

BMI/CS 776

www.biostat.wisc.edu/bmi776/

Spring 2021
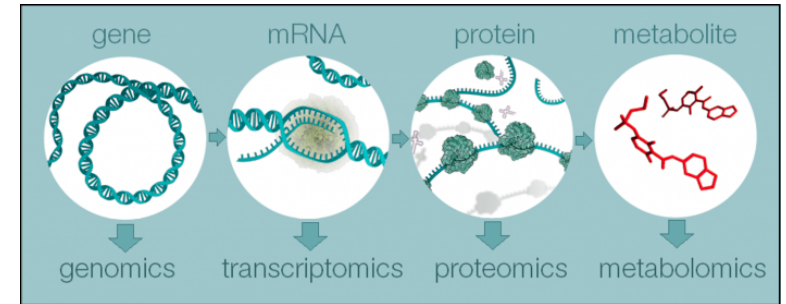
Daifeng Wang

daifeng.wang@wisc.edu
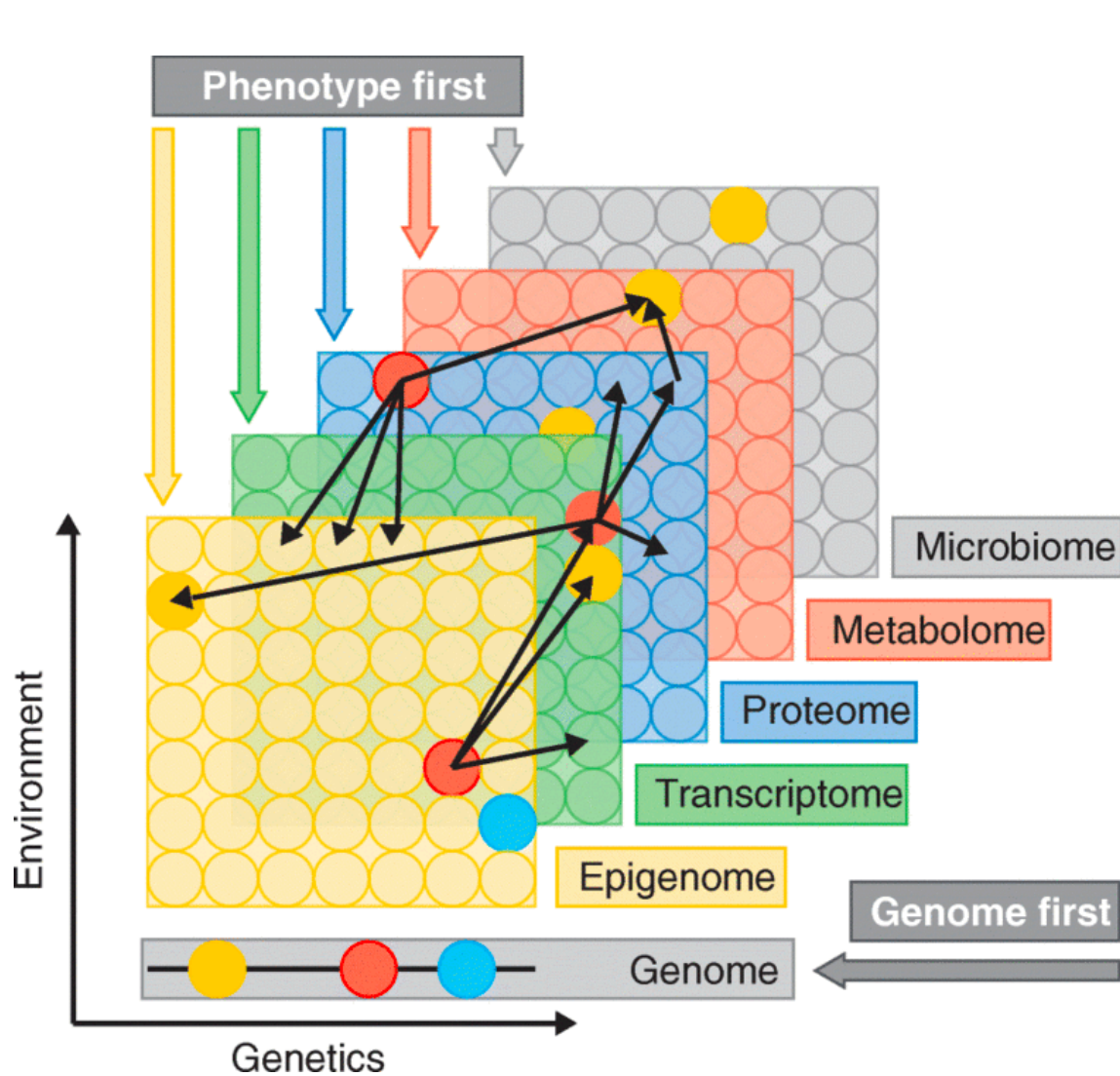
1

# Goals for lecture

- Multi-omics data

- Machine learning modeling
  - Empirical risk minimization (ERM)

- Multi-layer network clustering

- Dimensionality reduction & Spectral methods

- Decision tree

- Neural network

# Goals for lecture

- **Multi-omics data**
- Machine learning modeling
  - Empirical risk minimization (ERM)
- Multi-layer network clustering
- Dimensionality reduction & Spectral methods
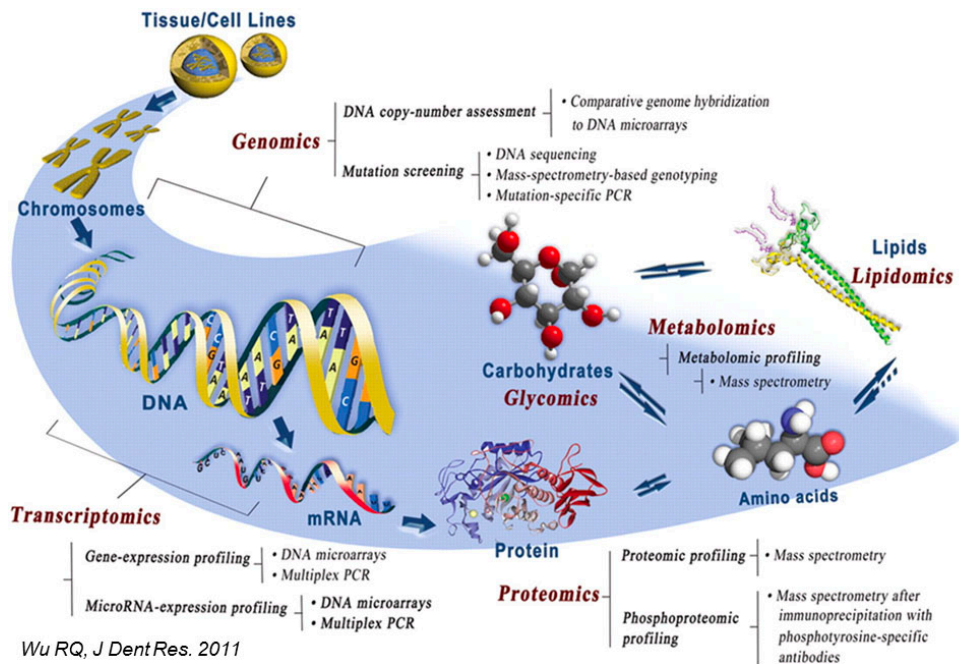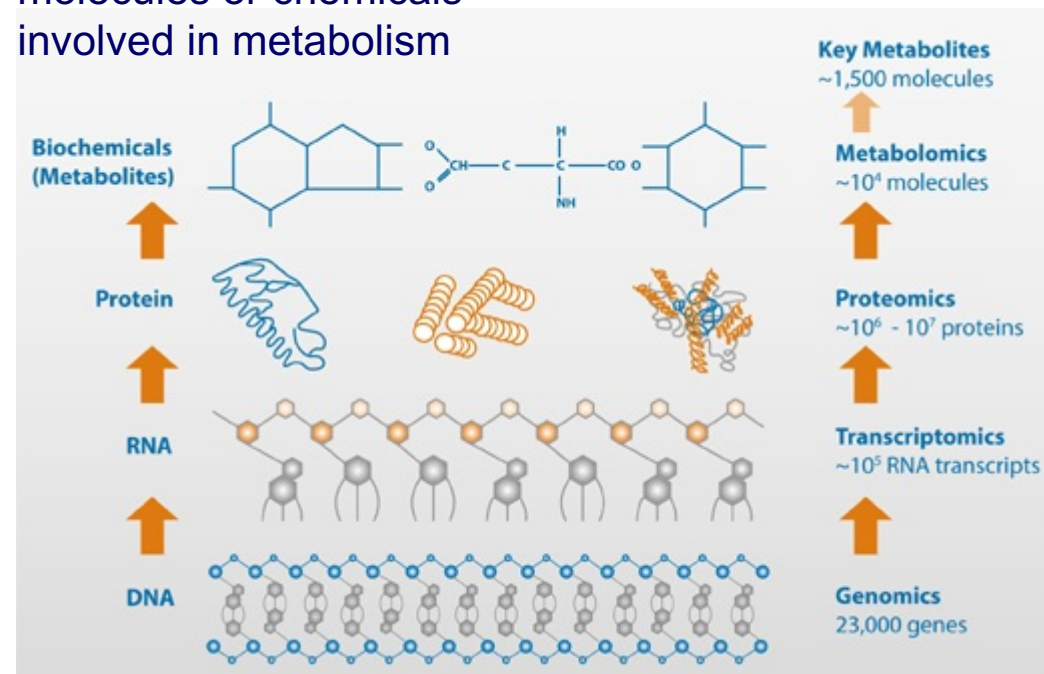- Decision tree
- Neural network

# Multi-omics

# Metabolites and Metabolomics



Wu RQ, J Dent Res. 2011

Metabolites are small molecules or chemicals involved in metabolism

# Multi-scale mechanisms



Nature Reviews | Genetics

# Functional genomics to understand mechanisms



Disease-associated genomic variants

How do variants function?

Gandal et al., Nature Neuroscience, 2016

7

# Example

# Hierarchical understanding from genotype to phenotype

**Elements**
- variants
- genes
- regulatory regions

→

**Interactions**
- gene regulation
- chromatin interaction
- TF binding

→

**Mechanisms**
- pathways
- circuits
- functions

→

**Prediction & Prioritization**
- disease variants & genes
- networks
- cell types

# Multi-omics for understanding functional genomics and gene regulation



Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

# Some multi-omics datasets

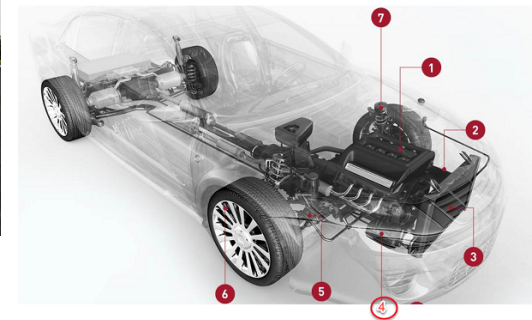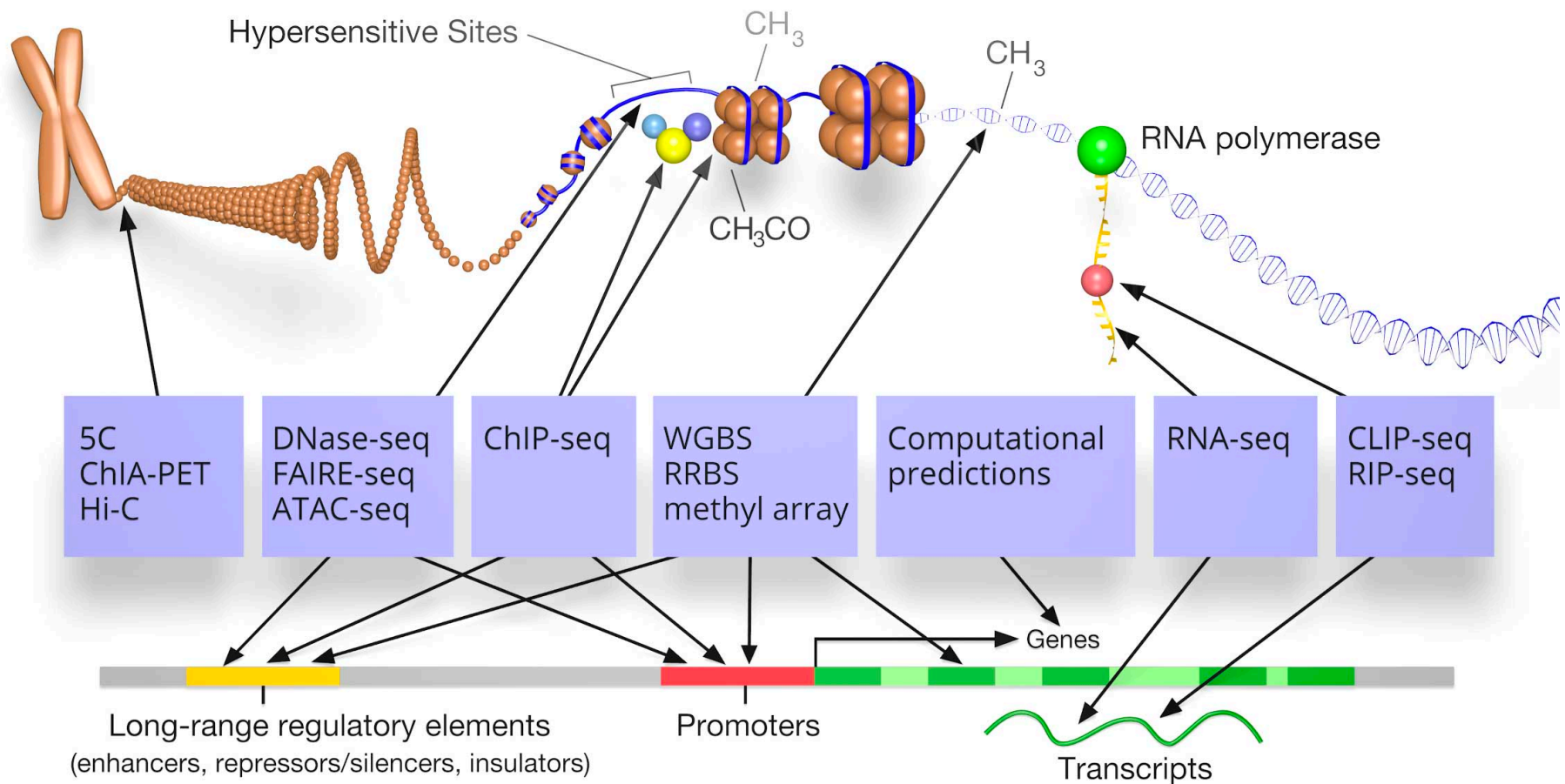| Human | 20,000 genes (2% genome) | Other genomic elements: non-coding RNAs, gene regulatory regions, repeats, and so on… (98% genome) |
|---|---|---|
| Cell lines | ENCODE (Encyclopedia of DNA Elements) Consortium (> 300 cell types) | |
| Tissues | Genotype-Tissue Expression (GTEx) (> 40 tissues) | |
| Cancers | The Cancer Genome Atlas (TCGA) (> 40 cancer types) | |
| Development | (13 developmental stages, 16 brain regions) | |
| Psychiatric disorders | PsychENCODE Consortium (~2,000 tissues incl. health, Schizophrenia, Autism, Bipolar) | |
| Neurodegenerative diseases | Religious Orders Study and Memory and Aging Project (ROSMAP) | International Parkinson's Disease Genomics Consortium (IPDGC) |

# Goals for lecture

- Multi-omics data

- Machine learning modeling
  - Empirical risk minimization (ERM)

- Multi-layer network clustering

- Dimensionality reduction & Spectral methods

- Decision tree

- Neural network

# Multi-omics data integration



**a** Concatenation-based integration   **b** Transformation-based integration   **c** Model-based integration

SNP matrix — Gene expression matrix — miRNA matrix

Patients: Patient 1, Patient 2, Patient 3, ..., Patient n

SNP variables (SNP 1, SNP 2, SNP 3, ..., SNP i)

Gene expression variables (Gene 1, Gene 2, Gene 3, ..., Gene j)

miRNA variables (miRNA 1, miRNA 2, miRNA 3, ..., miRNA k)

Phenotype 1
Phenotype 2

http://www.nature.com/nrg/journal/v16/n2/full/nrg3868.html

# Multi-omics data modeling



Phenotype

Higher-order groupings (e.g., pathways, circuits)

**Biological Interpretation**

Cell types

Regulatory elements

Modules

Genes

Genotype

**Abstract** — **Specified**

**Relationship** — Co-expression
Gene 1
Gene 2
expression
sample

**Topology** — Association network

**Influence** — Causal network

**Logic** — Circuit

**Dynamics** — Process
$X_{t+1} = AX_t + BU_t$

**Interactions between elements** (e.g., co-expression, regulation)

**Integrative & Predictive model** (e.g., deep neural network)

**Genomic elements** (e.g., GWAS, differentially expressed genes)

Raw data — Pre-processing — Clean data

```
A C G T C
G C G T A
G T C C G
T T A G T
C G T A G
G A G A A
```

Feature extraction — Features — Training — Model — Evaluation — Results

\* Spectrum inspired from Ideker & Lauffenburger, Trends in Biotechnology, 2003
\* Christof Angermueller et al. Mol Syst Biol 2016;12:878

# Multiview learning for understanding functional multi-omics
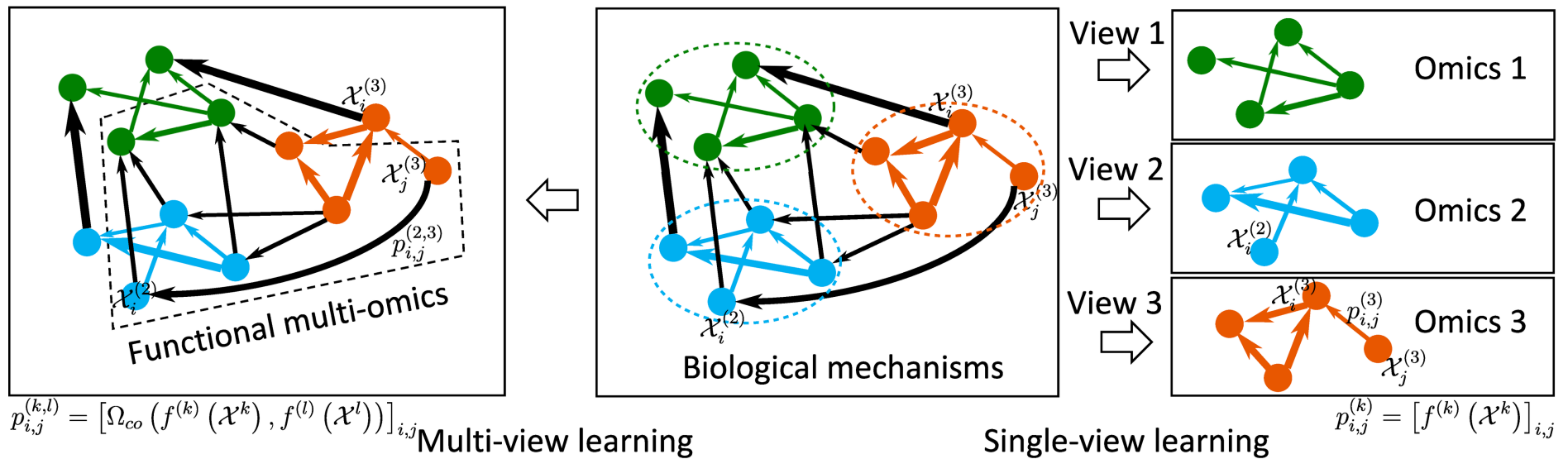


$$p_{i,j}^{(k,l)} = \left[ \Omega_{co} \left( f^{(k)} \left( \mathcal{X}^k \right), f^{(l)} \left( \mathcal{X}^l \right) \right) \right]_{i,j}$$ Multi-view learning

$$p_{i,j}^{(k)} = \left[ f^{(k)} \left( \mathcal{X}^k \right) \right]_{i,j}$$

Single-view learning
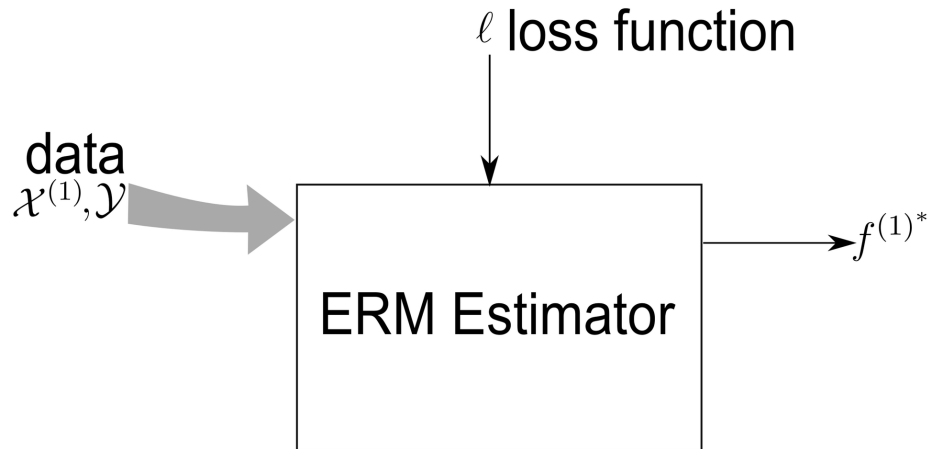
- For example, gene regulation can relate to

  1. Genomics; e.g., SNPs

  2. Transcriptomics; e.g., genes

  3. Proteomics; e.g., transcription factors (TFs)

Cross-omics interactions

$\Omega_{co}(f^{(1)}, f^{(3)})$: SNPs break TF binding sites

$\Omega_{co}(f^{(2)}, f^{(3)})$: TFs control gene expression

$\Omega_{co}(f^{(1)}, f^{(2)})$: SNPs associate with gene expression (e.g, eQTLs)

15

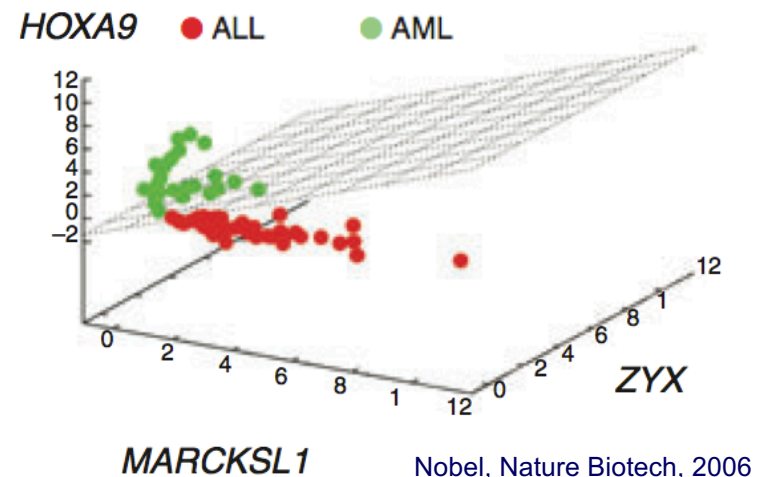# Empirical risk minimization (ERM) for machine learning modeling

$\ell$ loss function

data
$\mathcal{X}^{(1)}, \mathcal{Y}$

ERM Estimator $\longrightarrow f^{(1)*}$

hypothesis space $\mathcal{F}^{(1)}$

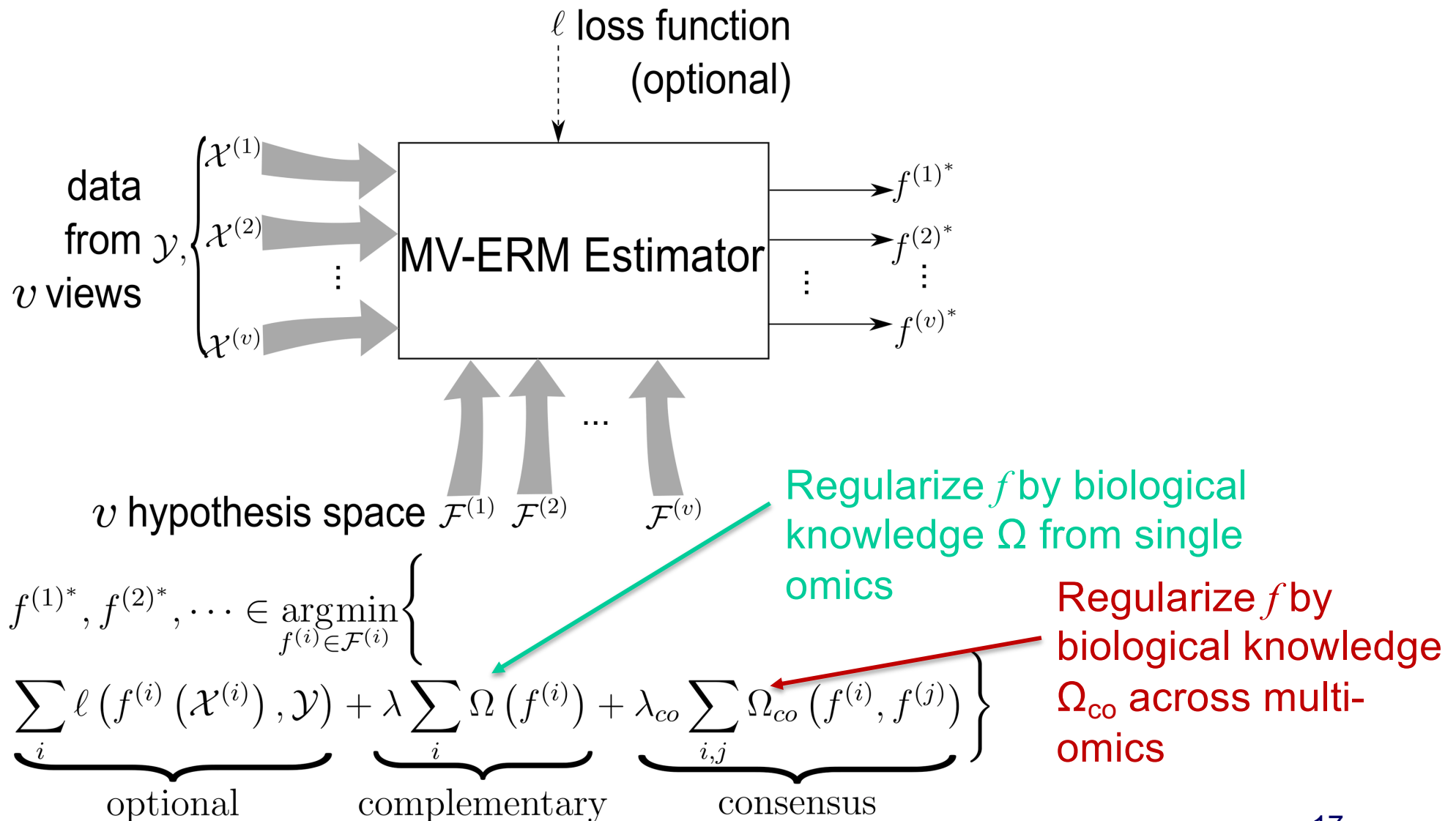$$R(f) = \frac{1}{|S|} \sum_{(x_i, y_i) \in S} \ell(f(x_i), y_i)$$

$$f^* \in \underset{f}{\text{argmin}}\{R(f) + \lambda\Omega(f)\}$$

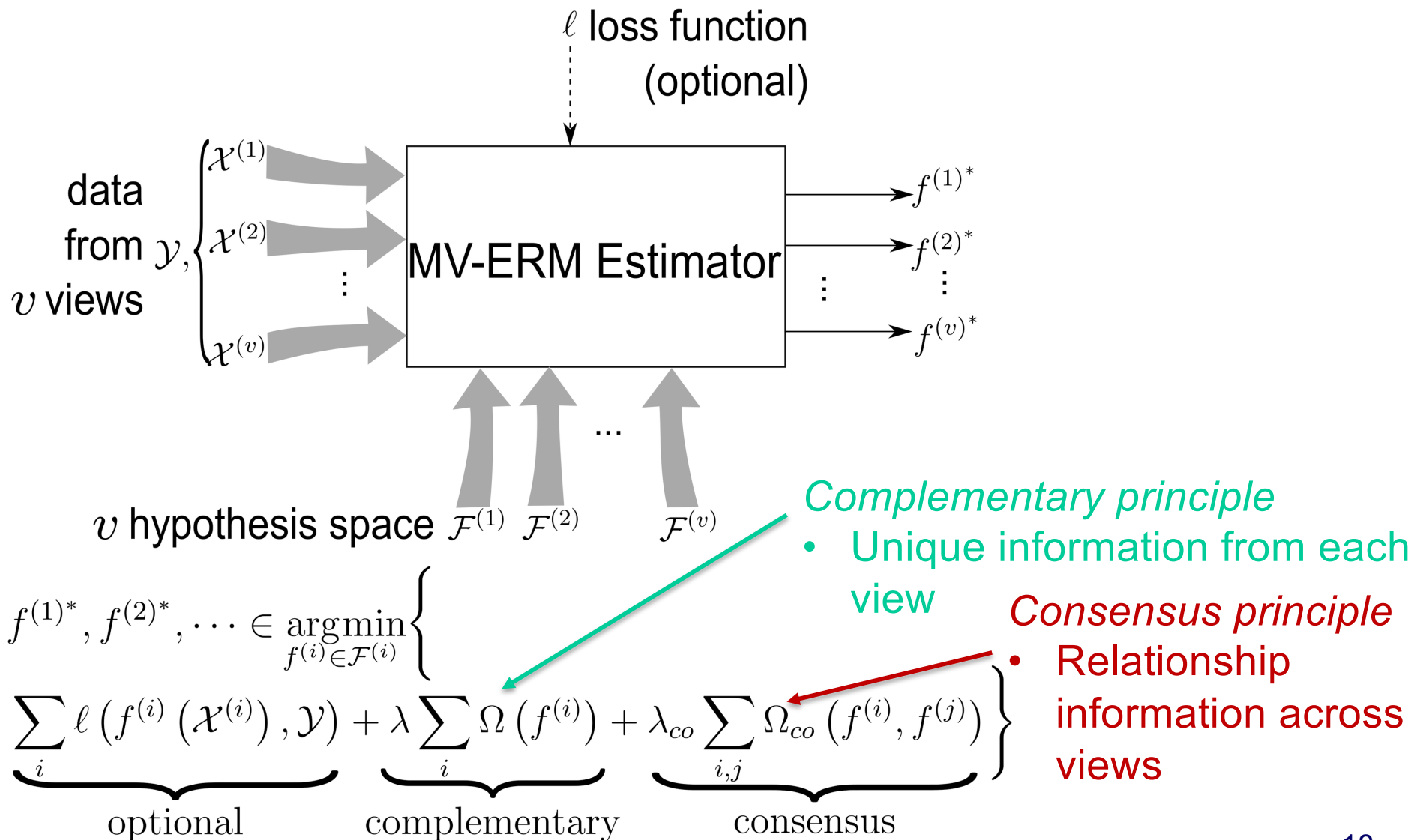Regularize $f$ by biological knowledge $\Omega$

- e.g., Leukemia patient classification
  - $y_i$: Acute lymphoblastic leukemia (ALL) vs. Acute myeloid leukemia (AML)
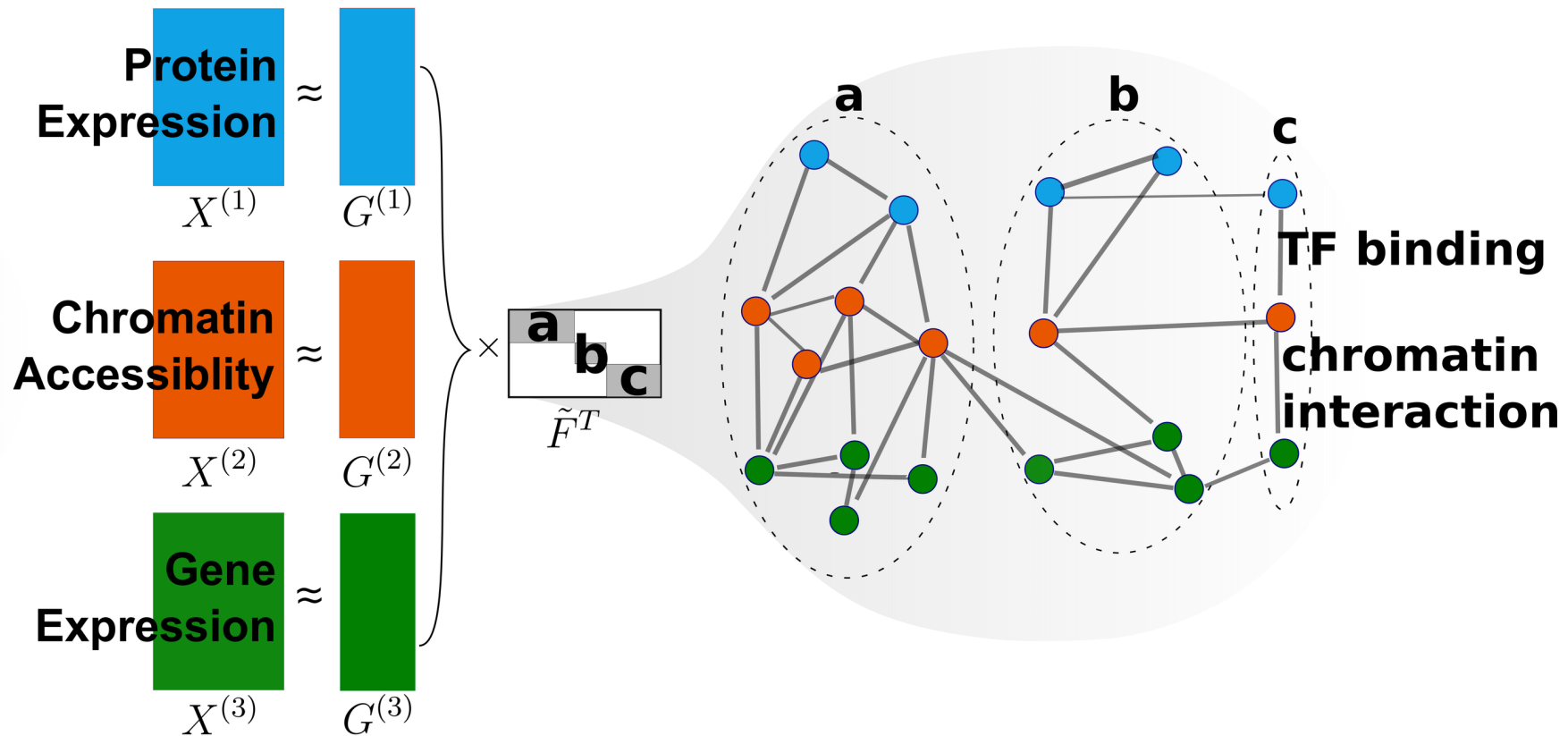  - $x_i$: gene expression
  - $f$: SVM



HOXA9   ● ALL   ● AML

MARCKSL1    ZYX

Nobel, Nature Biotech, 2006

16

# Empirical risk minimization for multi-view learning (MV-ERM)



$\ell$ loss function (optional)

data from $\mathcal{Y}$, $v$ views
$\left\{\begin{array}{l}\mathcal{X}^{(1)}\\ \mathcal{X}^{(2)}\\ \vdots\\ \mathcal{X}^{(v)}\end{array}\right.$

MV-ERM Estimator

$f^{(1)*}$
$f^{(2)*}$
$\vdots$
$f^{(v)*}$

$v$ hypothesis space $\mathcal{F}^{(1)}$ $\mathcal{F}^{(2)}$ $\mathcal{F}^{(v)}$

Regularize $f$ by biological knowledge $\Omega$ from single omics

Regularize $f$ by biological knowledge $\Omega_{co}$ across multi-omics

$$f^{(1)*}, f^{(2)*}, \cdots \in \underset{f^{(i)} \in \mathcal{F}^{(i)}}{\operatorname{argmin}} \left\{ \underbrace{\sum_i \ell\left(f^{(i)}\left(\mathcal{X}^{(i)}\right), \mathcal{Y}\right)}_{\text{optional}} + \lambda \underbrace{\sum_i \Omega\left(f^{(i)}\right)}_{\text{complementary}} + \lambda_{co} \underbrace{\sum_{i,j} \Omega_{co}\left(f^{(i)}, f^{(j)}\right)}_{\text{consensus}} \right\}$$

# Consensus and complementary principles



$\ell$ loss function (optional)

data from $\mathcal{Y}$, $v$ views

$\begin{cases} \mathcal{X}^{(1)} \\ \mathcal{X}^{(2)} \\ \vdots \\ \mathcal{X}^{(v)} \end{cases}$

MV-ERM Estimator

$f^{(1)*}$
$f^{(2)*}$
$\vdots$
$f^{(v)*}$

$v$ hypothesis space $\mathcal{F}^{(1)}$ $\mathcal{F}^{(2)}$ $\mathcal{F}^{(v)}$

*Complementary principle*
- Unique information from each view

*Consensus principle*
- Relationship information across views

$$f^{(1)*}, f^{(2)*}, \cdots \in \underset{f^{(i)} \in \mathcal{F}^{(i)}}{\arg\min} \left\{ \underbrace{\sum_i \ell\left(f^{(i)}\left(\mathcal{X}^{(i)}\right), \mathcal{Y}\right)}_{\text{optional}} + \lambda \underbrace{\sum_i \Omega\left(f^{(i)}\right)}_{\text{complementary}} + \lambda_{co} \underbrace{\sum_{i,j} \Omega_{co}\left(f^{(i)}, f^{(j)}\right)}_{\text{consensus}} \right\}$$

18

# Factorization-based MV-ERM framework



$$\left(G^{(1)*}, G^{(2)*}, G^{(3)*}, \tilde{F}^*\right) \in \underset{G^{(i)}, F^{(i)}, \tilde{F} \geq 0}{\arg\min} \sum_{i=1}^{3} \left\{ \left\| X^{(i)} - G^{(i)} F^{(i)T} \right\|_F^2 + \lambda \left\| F^{(i)} - \tilde{F} \right\|_F^2 \right\}$$

*Complementary* $G^{(i)}$     *Consensus*

- e.g., solved by Multi-view NMF (Liu et al., SIAM ICDM, 2013)   $\tilde{F}$
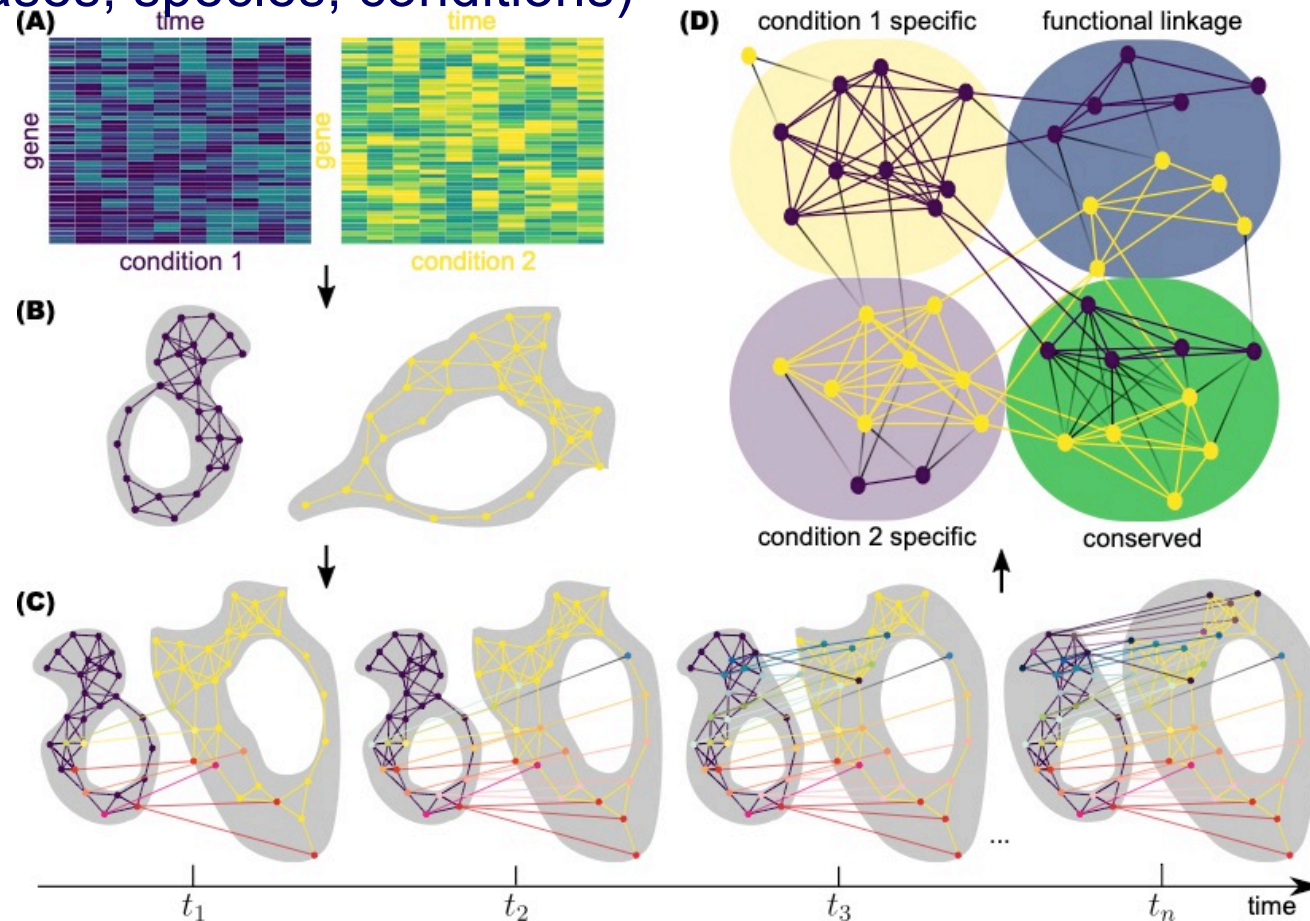
# Alignment-based MV-ERM



- For instance, Canonical correlation analysis (CCA)
  - Consensus only
  - $\Omega_{co}(.) = -tr(F_1^T X_1 X_2^T F_2)$ for two views $X_1$ and $X_2$ with linear projections $F_1$ and $F_2$

# ManiNetCluster: manifold alignment to reveal the functional links between gene networks

Multi-view datasets (e.g., diseases, species, conditions)
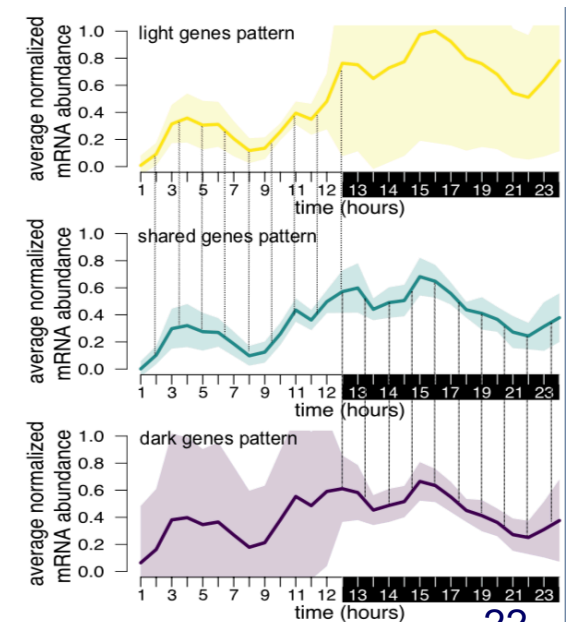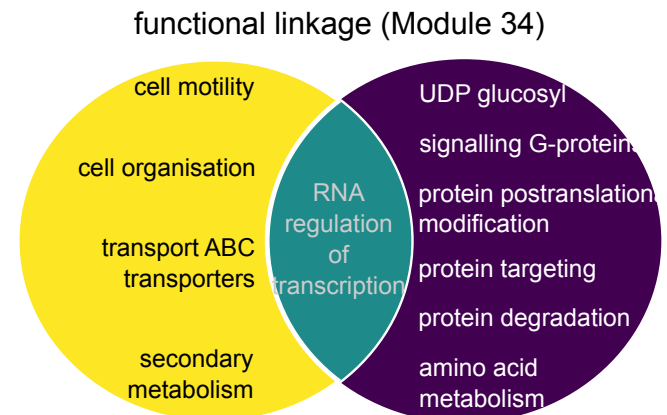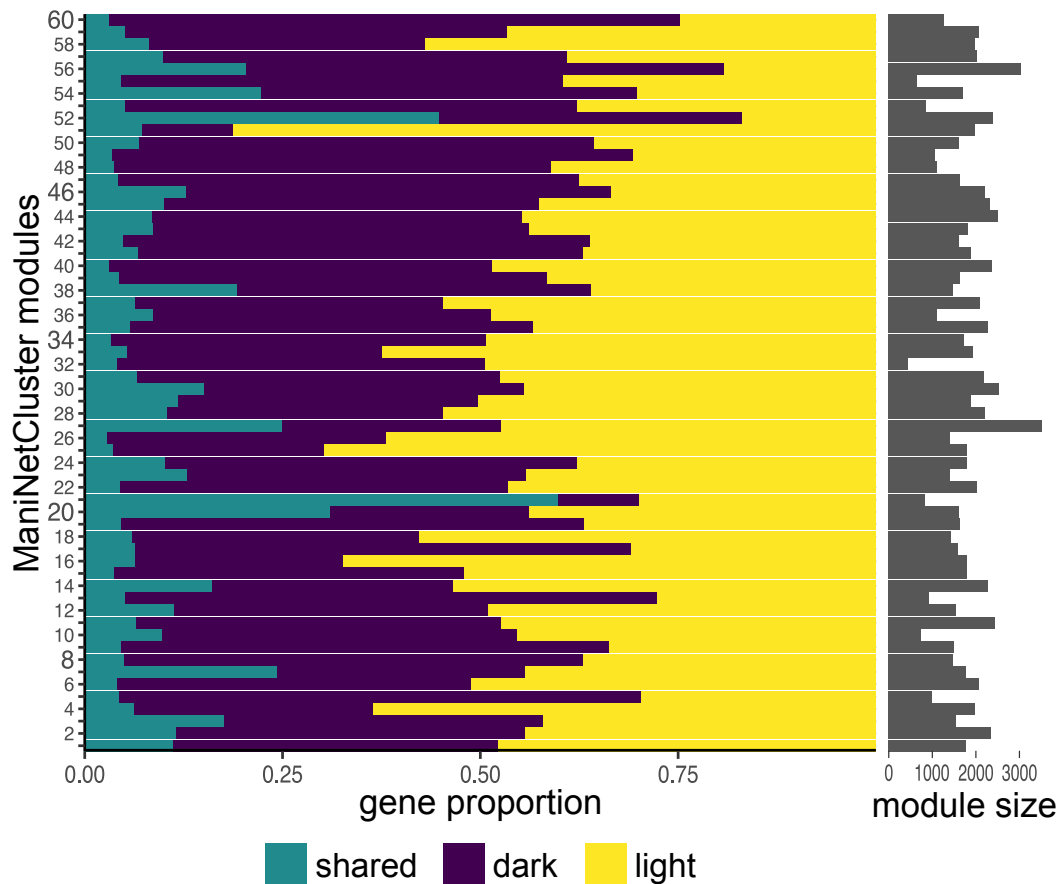
Functional linkages across dimensions



$$argmin_{f_x,f_y} \; \lambda \sum_{i,j} \left\| f_x(X_i) - f_x(X_j) \right\|^2 S_x(i,j) + \lambda \sum_{i,j} \left\| f_y(Y_i) - f_y(Y_j) \right\|^2 S_y(i,j) + (1-\lambda) \sum_{i,j} \left\| f_x(X_i) - f_y(Y_j) \right\|^2 W(i,j)$$

21

# ManiNetCluster: manifold alignment to reveal the functional links between gene networks

Application: genomic functional linkages between light and dark periods of green alga
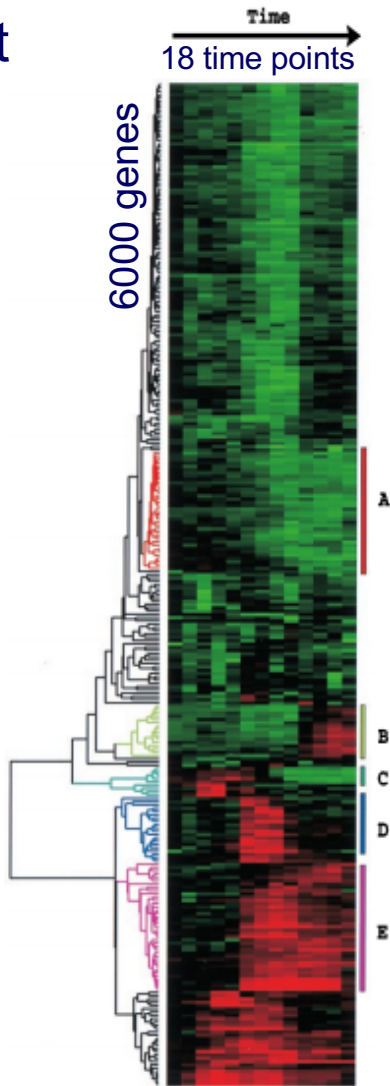


functional linkage (Module 34)

cell motility

UDP glucosyl

light genes pattern

shared genes pattern

dark genes pattern

shared   dark   light

gene proportion

module size

ManiNetCluster modules

Nguyen, Blaby, Wang, BMC Genomics, 2019

# Goals for lecture

- Multi-omics data
- Machine learning modeling
    - Empirical risk minimization (ERM)
- **Multi-layer network clustering**
- Dimensionality reduction & Spectral methods
- Decision tree
- Neural network

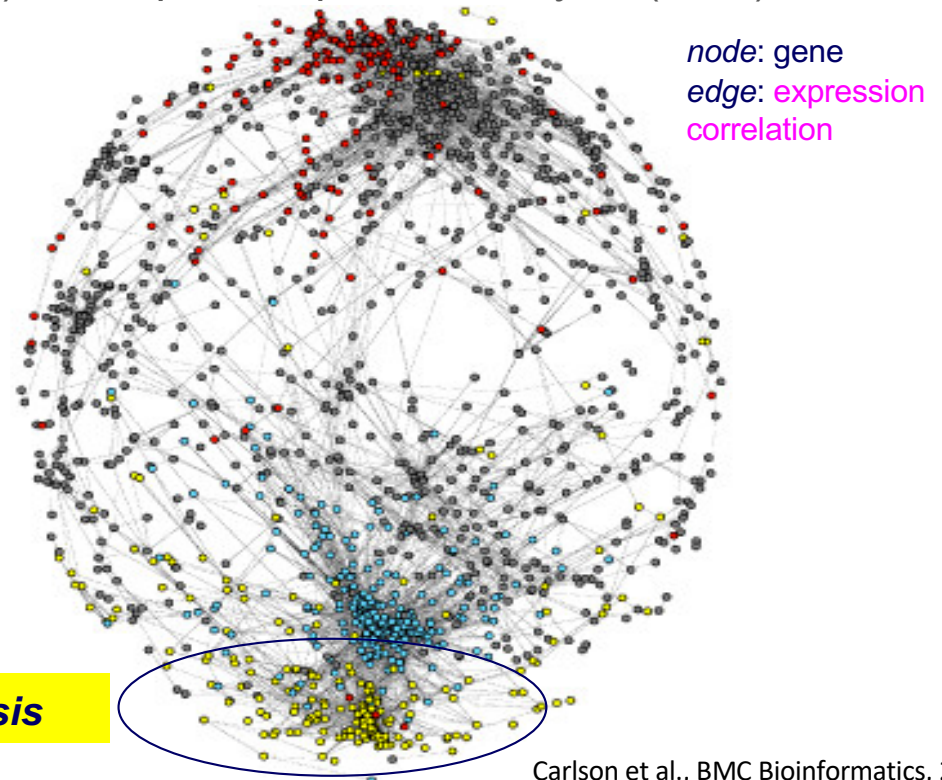# Co-expressed genes have similar functions in single species

Yeast cell cycle

**Time**
18 time points

6000 genes



A gene co-expression network (relationship) can reveal *functional groupings*

- Hierarchal clustering, K-means, Gaussian mixture model (GMM), Principal component analysis (PCA), …

*node*: gene
*edge*: expression correlation

**Protein synthesis**



Eisen et al., PNAS, 1998.

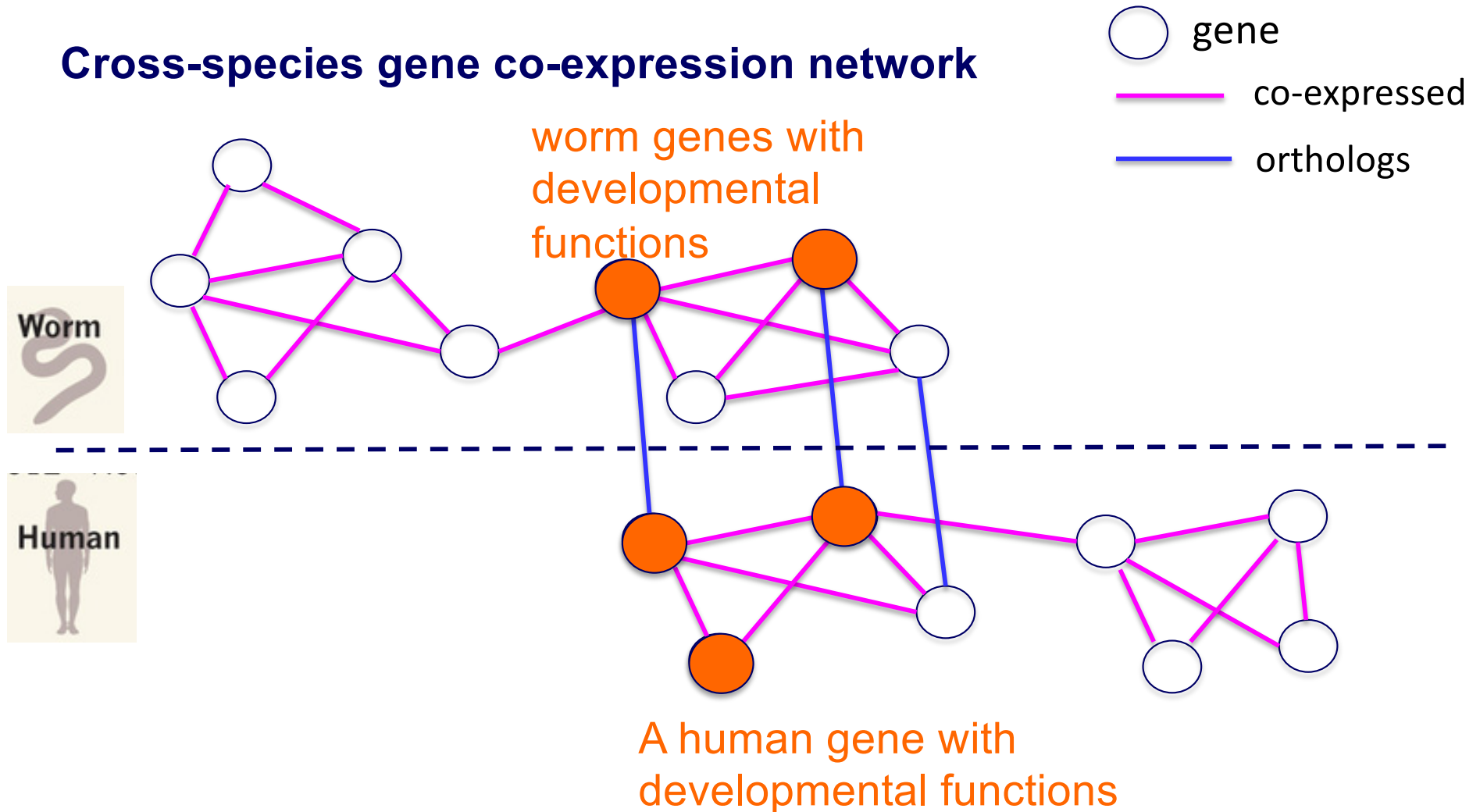Carlson et al., BMC Bioinformatics, 2006.

24

# Limited knowledge in single species

**~ half human genes, 1% human genome plus other 98% genomic elements (non-coding regions) with unknown functions**
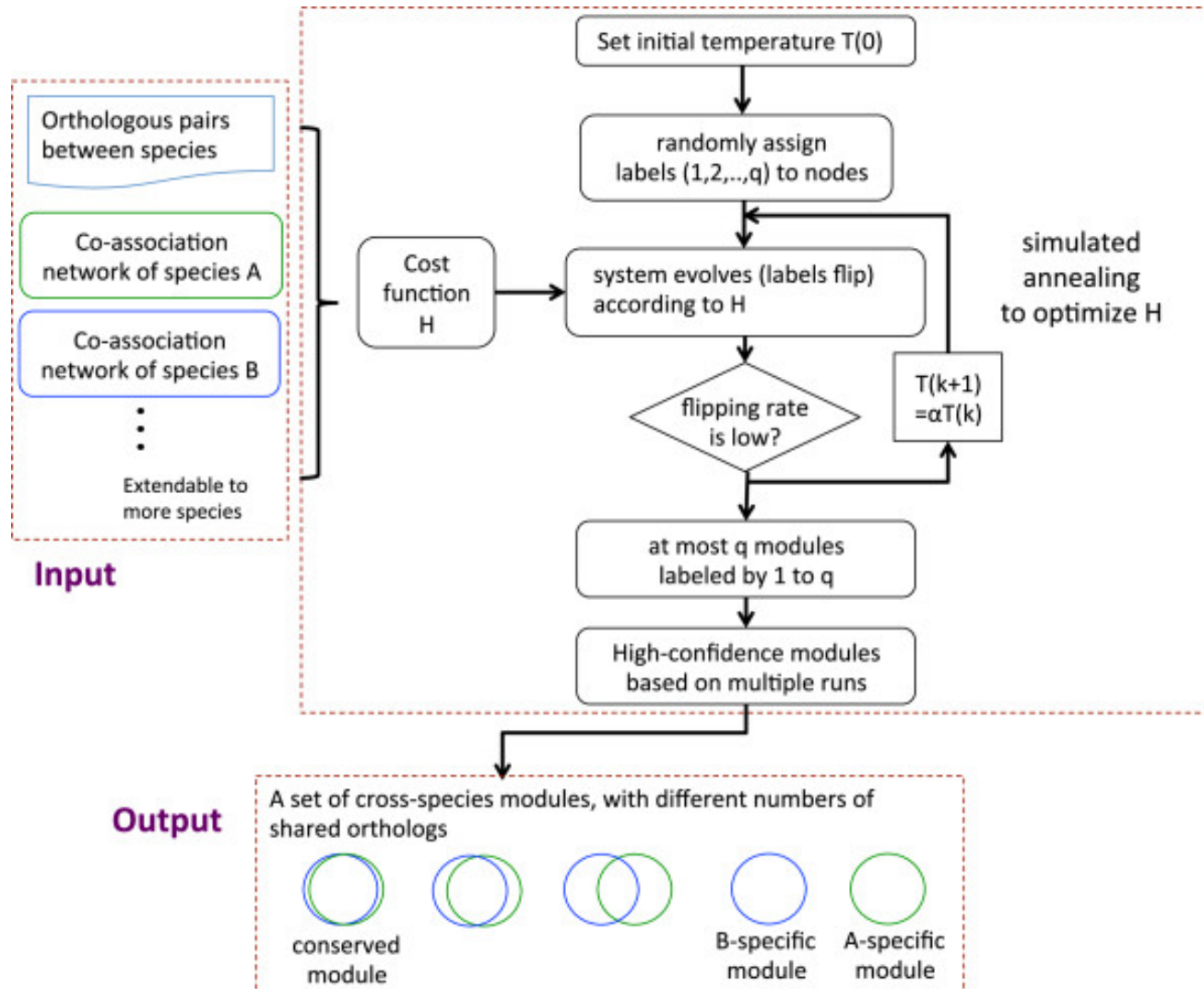


Cell adhesion (577, 1.9%)
Miscellaneous (1318, 4.3%)
Viral protein (100, 0.3%)
Transfer/carrier protein (203, 0.7%)
Transcription factor (1850, 6.0%)
Nucleic acid enzyme (2308, 7.5%)
Signaling molecule (376, 1.2%)
Receptor (1543, 5.0%)
Kinase (868, 2.8%)
Select regulatory molecule (988, 3.2%)
Transferase (610, 2.0%)
Synthase and synthetase (313, 1.0%)
Oxidoreductase (656, 2.1%)
Lyase (117, 0.4%)
Ligase (56, 0.2%)
Isomerase (163, 0.5%)
Hydrolase (1227, 4.0%)

Chaperone (159, 0.5%)
Cytoskeletal structural protein (876, 2.8%)
Extracellular matrix (437, 1.4%)
Immunoglobulin (264, 0.9%)
Ion channel (406, 1.3%)
Motor (376, 1.2%)
Structural protein of muscle (296, 1.0%)
Proto-oncogene (902, 2.9%)
Select calcium-binding protein (34, 0.1%)
Intracellular transporter (350, 1.1%)
Transporter (533, 1.7%)

nucleic acid binding
signal transduction
enzyme
none

Molecular function unknown (12809, 41.7%)

How do we know human gene functions during embryonic or brain development?

25

# Integration of co-expressed and orthologous genes across species to transfer function information
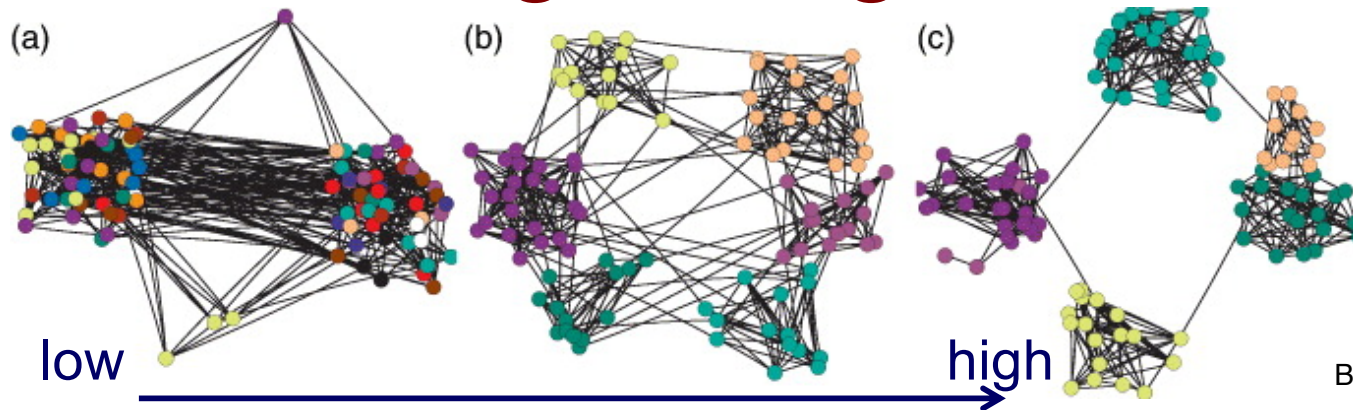
**Cross-species gene co-expression network**

gene

co-expressed

orthologs

worm genes with developmental functions

A human gene with developmental functions

Worm

Human



26

# OrthoClust: an orthology-based method for clustering cross-species networks (e.g., co-expression networks)

# Maximize "modularity" for clustering a single network



(a) low → (b) → (c) high

Brede, Europhysics Letters, 2010.

**Modularity** $Q$: measurement on strength of network division

$$Q = \frac{1}{2m} \sum_{i,j} \left( W_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\sigma_i \sigma_j}$$

sum over nodes within a group (module)

normalization
$m$: total number of edges

edge weight between nodes i and j

$\dfrac{k_i k_j}{2m} = p_{ij}$ =expected edge weight that would go between i and j

**Clustering goal:** assign each node a module to maximize "modularity" as an objective function (module is a group of highly connected nodes)

Newman, PNAS, 2006.

28

# OrthoClust: an orthology-based method for clustering cross-species networks

Every node i is assigned with a module number $\sigma_i$.



Objective function

$$H = \boxed{Q_A + Q_B} + \boxed{\kappa \sum_{(i,j') \in Ortho} \delta_{\sigma_i \sigma_{j'}}}$$

reward an orthologous pair in the same module

"**Modularity**" in species A  +  "**Modularity**" in species B  + **consistency** between A & B

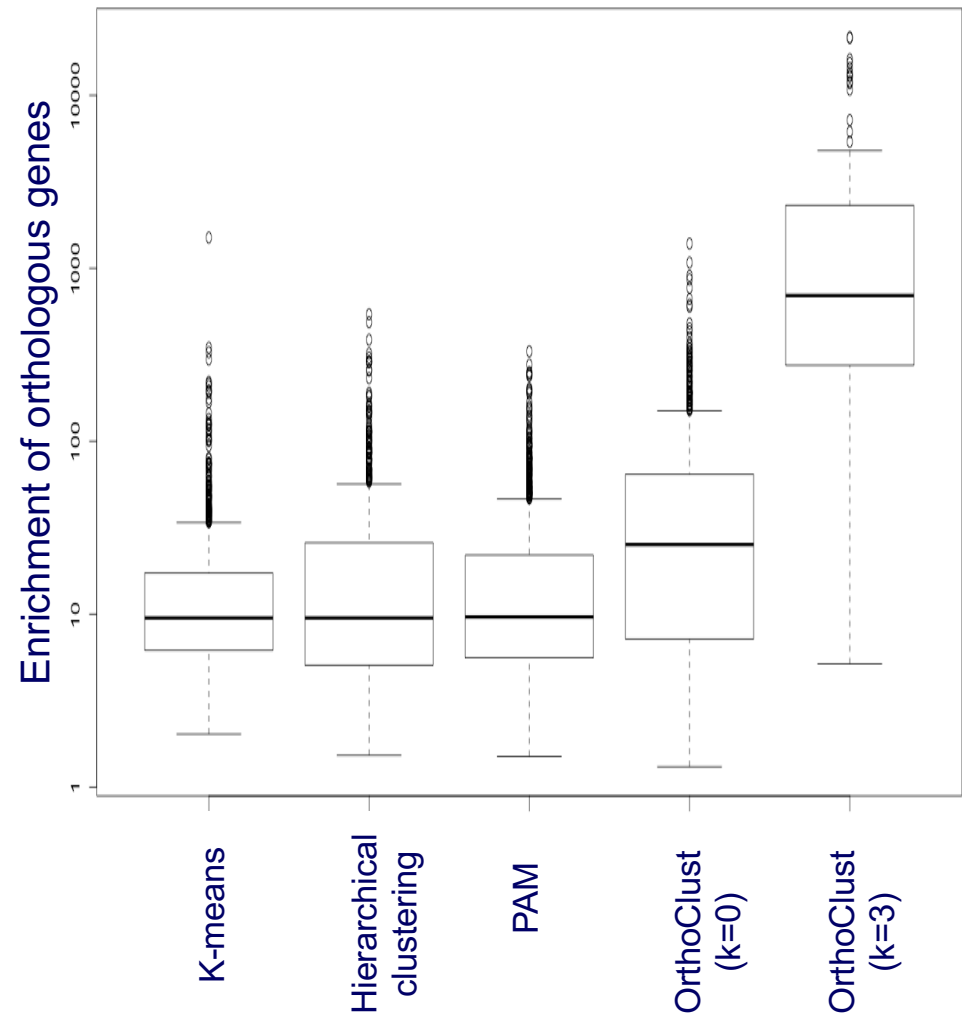# Conserved gene co-expression modules discovered human genes having developmental functions



human-worm-fly conserved

worm-fly conserved

worm-specific

Gene-gene co-association

1

0

mod ENCODE
The National Human Genome Research Institute
model organism ENCyclopedia Of DNA Elements

Worm
datasets: 219
factors: 93

Fly
datasets: 93
factors: 52

Human
datasets: 707
factors: 165

20,377 gene co-expression network across 33 developmental stages

13,623 gene co-expression network across 30 developmental stages

19,901 gene co-expression network across 19 cell lines

# of Genes    1000       500       0

Human

Worm

Fly

16 human-worm-fly conserved modules

Signal transduction, cytoskeletal
Morphogenesis, epidermal GF
Histone mRNA proc., nuc. export
Topoisomerase, RNA POL II
Cell cyc. ctrl, signal transduction
Ribosome
Translocase, folding, G1S cell cyc.
La autoantigen
Signal transduction, Integrins
Spliceosome

**19,901**    **20,377**    **13,623**    genes (~55,000)

30

Gerstein*,…, **Wang**\*, et al., *Nature*, 2014

# OrthoClust reveals better genomic functional groups

**OrthoClust's modular genes have similar functions**

**OrthoClust clusters more orthologs than other clustering methods**



Left figure axis: Gene ontology (GO) similarity of gene pairs; annotation * P=3x10$^{-83}$; categories: Pairs between OrthoClust modules, Pairs in OrthoClust modules

Right figure axis: Enrichment of orthologous genes; categories: K-means, Hierarchical clustering, PAM, OrthoClust (k=0), OrthoClust (k=3)

# Developmental hourglass behavior across conserved modules in a species

**Inter-organism**
Temporal differences among ortholog expression levels are minimized at phylotypic stage **across different species**.

Kalinka et al. Nature, 2010

Hourglass

Expression

**Intra-organism**
Temporal differences among ortholog expression levels are minimized at phylotypic stage **across conserved modules in a species (fly)**.

Gerstein*,…, **Wang***, et al., Nature, 2014

# Human and Rhesus brain developmental "hourglass"



Li, …, **Wang**, ..., Sestan, *Science*, 2018

Ying, ..., Sestan, *Science*, 2018

# Goals for lecture

- Multi-omics data

- Machine learning modeling
  - Empirical risk minimization (ERM)

- Multi-layer network clustering

- Dimensionality reduction & Spectral methods

- Decision tree

- Neural network

# Reading list for spectral methods

- O Alter et al. (2000). "Singular value decomposition for genome-wide expression data processing and modeling." PNAS 97: 10101

- Langfelder P, Horvath S (2007) Eigengene networks for studying the relationships between co-expression modules. BMC Systems Biology 2007, 1:54

- Z Zhang et al. (2007) "Statistical analysis of the genomic distribution and correlation of regulatory elements in the ENCODE regions." Genome Res 17: 787

- TA Gianoulis et al. (2009) "Quantifying environmental adaptation of metabolic pathways in metagenomics." PNAS 106: 1374.

# What is Principal component analysis (PCA) ?

- A technique used to reduce the dimensionality of a data set by finding directions of maximum variability
- Projection (typically a rotation) into new axes
- But still retains the dataset's variation

Adapted from http://www.astro.princeton.edu/~gk/A542/PCA.ppt

# PCA Matrix

1. Start with dataset of k variables $X = x_1, x_2 ... x_k$ and n observations.
2. Construct **covariance or correlation matrix** for variables.
3. **The Eigenvalue Problem** or Eigenanalysis: matrix diagonalization and solve for eigenvalues and eigenvectors

E.g. Start with a bunch of coordinates



| Observations | X1 | X2 |
|---|---|---|
| 1 | 2 | 5 |
| 2 | 5 | 6 |
| 3 | 4 | 2 |
| 4 | 3 | 7 |
| 5 | 9 | -5 |
| ... | | |
| n | -5 | -8 |

37

# Interpretation:
# Eigenvalues & Eigenvectors

Adapted from http://myweb.dal.ca/~hwhitehe/BIOL4062/pca.ppt

# Quick Refresher on Matrices

$$\begin{pmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x_3 & y_3 & z_3 \end{pmatrix} * \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} ax_1 + by_1 + cz_1 \\ ax_2 + by_2 + cz_2 \\ ax_3 + by_3 + cz_3 \end{pmatrix}$$

Matrix **A** is 3x4

$$\begin{bmatrix} 8 & 3 & 0 & 1 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

Matrix **B** is 4x4

$$\begin{bmatrix} 5 & \cdot & \cdot & \cdot \\ 4 & \cdot & \cdot & \cdot \\ 3 & \cdot & \cdot & \cdot \\ 1 & \cdot & \cdot & \cdot \end{bmatrix}$$

Matrix **C** is 3x4

$$= \begin{bmatrix} 53 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

$$\text{because } c_{11} = \sum_{k=1}^{4} a_{1k} b_{k1} = 8 \cdot 5 + 3 \cdot 4 + 0 \cdot 3 + 1 \cdot 1 = 53$$



http://eli.thegreenplace.net/2015/visualizing-matrix-multiplication-as-a-linear-combination/

39

http://www.catonmat.net/blog/mit-linear-algebra-part-three/

# SVD for gene expression data
# (Alter et al, PNAS 2000)

$$X = USV^{\mathrm{T}}$$

# Notation

- **m=1000 genes**
  - row-vectors
  - 10 eigengene ($v_i$) of dimension 10 conditions

- **n=10 conditions (assays)**
  - column vectors
  - 10 eigenconditions ($u_i$) of dimension 1000 genes

$$X = USV^T$$

http://www.gersteinlab.org/courses/452/

# SVD as sum of rank-1 matrices

- $A = USV^T$

- $A = s_1 \boldsymbol{u}_1 \boldsymbol{v}_1{}^T + s_2 \boldsymbol{u}_2 \boldsymbol{v}_2{}^T + \dots + s_n \boldsymbol{u}_n \boldsymbol{v}_n{}^T$

an outer product $(uv^T)$ giving a matrix rather than the scalar of the inner product

- $s_1 \geq s_2 \geq \dots \geq s_n \geq 0$

- What is the rank-r matrix $\hat{A}$ that best approximates $A$ ?

  - Minimize $\sum_{i=1}^{m} \sum_{j=1}^{n} \left( \hat{A}_{ij} - A_{ij} \right)^2$

LSQ approx. If r=1, this amounts to a line fit.

- $\hat{A} = s_1 \boldsymbol{u}_1 \boldsymbol{v}_1{}^T + s_2 \boldsymbol{u}_2 \boldsymbol{v}_2{}^T + \dots + s_r \boldsymbol{u}_r \boldsymbol{v}_r{}^T$

- Very useful for matrix approximation

http://www.gersteinlab.org/courses/452/

# Potential problems of SVD/PCA

If the dataset…

- Lacks Independence
  - **NO PROBLEM**
- Lacks Normality
  - Normality desirable but not essential
- Lacks Precision
  - Precision desirable but not essential
- Lacks Linearity
  - **Problem**: Use other non-linear (kernel) methods
- Many Zeroes in Data Matrix (Sparse)
  - **Problem**: Use Correspondence Analysis

# Conclusion

- SVD is the "absolute high point of linear algebra"

- SVD is difficult to compute; but once we have it, we have many things

- SVD finds the best approximating subspace, using <span style="color:red">linear transformation</span>

- Simple SVD cannot handle translation, non-linear transformation, separation of labeled data, etc.

- Good for exploratory analysis; but once we know what we look for, use appropriate tools and model the structure of data explicitly!

- http://genomicsclass.github.io/book/pages/pca_svd.html

# Goals for lecture

- Multi-omics data
- Machine learning modeling
  - Empirical risk minimization (ERM)
- Multi-layer network clustering
- Dimensionality reduction & Spectral methods
- Decision tree
- Neural network

# Reading list

- What are decision trees?

  – Nat Biotechnol. 2008 Sep; 26(9): 1011–1013.

- Data mining in the Life Sciences with Random Forest: a walk in the park or lost in the jungle?

  – https://academic.oup.com/bib/article/14/3/315/255469

# Decision Trees



A decision tree

- **Classify data by asking questions** that divide data in subgroups
- Keep asking questions until subgroups become homogenous
- Use **tree** of questions to make predictions

Example: Is a picture taken inside or outside?

**a**

| Gene Pair | Interact? | Expression correlation | Shared localization? | Shared function? | Genomic distance |
|---|---|---|---|---|---|
| A-B | Yes | 0.77 | Yes | No | 1 kb |
| A-C | Yes | 0.91 | Yes | Yes | 10 kb |
| C-D | No | 0.1 | No | No | 1 Mb |
| ⋮ | | | | | |

**b**

**A hypothetical example of how a decision tree might predict protein-protein interactions**
Nat Biotechnol. 2008 Sep; 26(9): 1011–1013.

48

# Terminology related to Decision Trees



ROOT Node

Branch/ Sub-Tree

Splitting

Decision Node

A Decision Node

Terminal Node

Decision Node

Terminal Node

B

Terminal Node

C

Terminal Node

Terminal Node

**Note:-** A is parent node of B and C.

49

# What makes a good rule?

- Want resulting groups to be as homogenous as possible



Rule 1

Rule 2

2/3 Groups homogenous
→Good rule

All groups still 50/50
→ Unhelpful rule

Nando de Freitas 2012 University of British Columbia CPSC 340

# Quantifying the value of rules

- ## Decrease in inhomogeneity
  - Most popular metric: Information theoretic entropy

  - Use frequency of classifier characteristic within group as probability

  - Minimize entropy to achieve homogenous group

$$S = -\sum_{i=1}^{m} p_i \log p_i$$

# Algorithm

- ## For each characteristic:
  - Split into subgroups based on each possible value of characteristic

- ## Choose rule from characteristic that maximizes decrease in inhomogeneity

- ## For each subgroup:
  - if (inhomogeneity < threshold):
    - Stop
  - else:
    - Restart rule search (recursion)

# Retrospective Decision Trees

Analysis of the Suitability of 500
M.thermo. proteins
to find optimal sequences purification



[Bertone et al. NAR ('01)]

53

# Retrospective Decision Trees
## Nomenclature

**356 total**

**Not Expressible** ┈┈┈▶ 143

**Expressible** ◀┈┈┈ 213

hydro phobe

53    28

90    185

Has a hydrophobic stretch? (Y/N)

[Bertone et al. NAR ('01)]

# Extensions of Decision Trees

- Decision Trees method is very sensitive to noise in data
- Random forests is an ensemble of decision trees, and is much more effective.



Decision Tree vs Random Forest

DT VS RF

Overfit          Reasonably Smooth

# Exercise

- A Complete Tutorial on Tree Based Modeling from Scratch (in R & Python)

  – https://www.analyticsvidhya.com/blog/2016/04/complete-tutorial-tree-based-modeling-scratch-in-python/

- Random Forests in R

  – https://www.r-bloggers.com/random-forests-in-r/

  – http://dni-institute.in/blogs/random-forest-using-r-step-by-step-tutorial/

# Goals for lecture

- Multi-omics data
- Machine learning modeling
  - Empirical risk minimization (ERM)
- Multi-layer network clustering
- Dimensionality reduction & Spectral methods
- Decision tree
- Neural network

# Reading list

- **Deep learning for computational biology**
  - http://msb.embopress.org/content/12/7/878
- **Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning**
  - https://www.nature.com/articles/nbt.3300
- https://github.com/hussius/deeplearning-biology
- The Incredible Convergence Of Deep Learning And Genomics

# Machine learning and representation learning

**Christof Angermueller et al. Mol Syst Biol 2016;12:878**

# Artificial Neural Network



**A**

Input layer | Hidden layer | Output layer

$w^{(1)}$

$w^{(2)}$

FORWARD PROPAGATION

BACKWARD PROPAGATION

PREDICTED label

TRUE label

$f(x) = 0.7 \neq y = 0.8$

$L(w) = (0.7 - 0.8)^2$ **LOSS**

**B**

Inputs | Weighted sum | Activation function | Output

$1 \times 0.6 +$
$0 \times 0.4 +$
$1 \times 0.2 = 0.8$

$\max(0, 0.8)$

ReLU

**C**

$w \quad \eta \Delta w \quad w' = w + \eta \Delta w$

$L(w)$

Local optimum

Global optimum

60

Christof Angermueller et al. Mol Syst Biol 2016;12:878

# The Incredible Convergence Of Deep Learning And Genomics

Johnny Israeli

# Principles of using neural networks for predicting molecular traits from DNA sequence

**Christof Angermueller et al. Mol Syst Biol 2016;12:878**

# Convolutional Neural Network

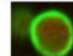# Convolution and pooling operators are stacked, thereby creating a deep network for image analysis

**Christof Angermueller et al. Mol Syst Biol 2016;12:878**

molecular systems biology

# A pre-trained network can be used as a generic feature extractor

**First layer features**

| | In top left? | | In top right? | ... | | In bottom right? |
|---|---|---|---|---|---|---|
| | 0.21 | | 0.24 | | | 0.01 |
| | 0.02 | | 0.01 | | | 0.25 |
| | 0.01 | | 0.03 | | | 0.19 |

**Third layer features**

| | In left? | | In right? | ... | | In bottom? |
|---|---|---|---|---|---|---|
| | 2.51 | | 0.02 | | | 2.92 |
| | 0.03 | | 0.01 | | | 0.02 |
| | 0.02 | | 0.01 | | | 0.01 |

**Christof Angermueller et al. Mol Syst Biol 2016;12:878**

molecular systems biology

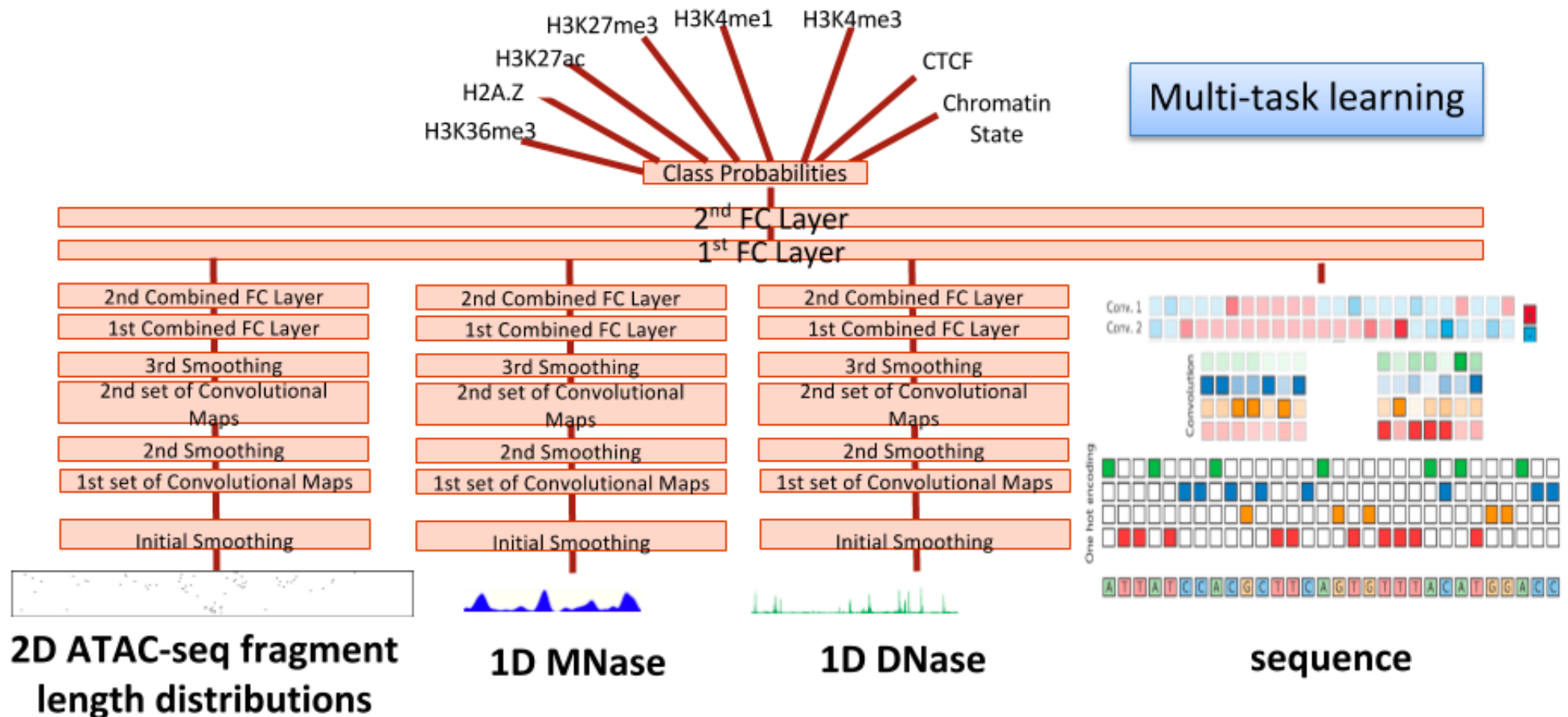# Data normalization for and pre-processing for deep neural networks



Christof Angermueller et al. Mol Syst Biol 2016;12:878

# Overview of existing deep learning frameworks, comparing four widely used software solutions

|  | **Caffe** | **Theano** | **Torch7** | **TensorFlow** |
|---|---|---|---|---|
| Core language | C++ | Python, C++ | LuaJIT | C++ |
| Interfaces | Python, Matlab | Python | C | Python |
| Wrappers | | Lasagne, Keras, sklearn-theano | | Keras, Pretty Tensor, Scikit Flow |
| Programming paradigm | Imperative | Declarative | Imperative | Declarative |
| Well suited for | CNNs, Reusing existing models, Computer vision | Custom models, RNNs | Custom models, CNNs, Reusing existing models | Custom models, Parallelization, RNNs |

# THE CHROMPUTER

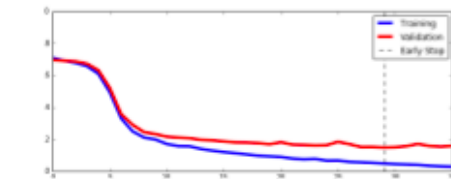Integrating 1D, 2D signals, and sequence to **predict multiple outputs**

Multi-task learning

H3K27me3   H3K4me1   H3K4me3
H3K27ac
H2A.Z                                          CTCF
H3K36me3                                 Chromatin State

Class Probabilities

2nd FC Layer
1st FC Layer

| 2nd Combined FC Layer | 2nd Combined FC Layer | 2nd Combined FC Layer |
| 1st Combined FC Layer | 1st Combined FC Layer | 1st Combined FC Layer |
| 3rd Smoothing | 3rd Smoothing | 3rd Smoothing |
| 2nd set of Convolutional Maps | 2nd set of Convolutional Maps | 2nd set of Convolutional Maps |
| 2nd Smoothing | 2nd Smoothing | 2nd Smoothing |
| 1st set of Convolutional Maps | 1st set of Convolutional Maps | 1st set of Convolutional Maps |
| Initial Smoothing | Initial Smoothing | Initial Smoothing |

**2D ATAC-seq fragment length distributions**

**1D MNase**

**1D DNase**

**sequence**

68

Johnny Israeli

# How to train your DragoNN

- https://drive.google.com/file/d/0B4Yo77Kh_QeeaXZKQUtZWjNrWkE

Johnny Israeli