

# Machine Learning for Healthcare

HST.956, 6.S897

## Lecture 19: Disease progression modeling & subtyping, Part 2

David Sontag



# Recap of goals of disease progression modeling

- Predictive:
  - *What will this patient's future trajectory look like?*
- Descriptive:
  - *Find markers of disease stage and progression, statistics of what to expect when*
  - *Discover new disease subtypes*
- Key challenges we will tackle:
  - Seldom directly observe disease stage, but rather only indirect observations (e.g. symptoms)
  - Data is censored – don't observe beginning to end

# Outline of today's lecture

1. Staging from cross-sectional data
  - Wang, Sontag, Wang, *KDD* 2014
  - Pseudo-time methods from computational biology
2. Simultaneous staging & subtyping
  - Young et al., *Nature Communications* 2018

# Outline of today's lecture

## 1. Staging from cross-sectional data

- Wang, Sontag, Wang, *KDD* 2014
- Pseudo-time methods from computational biology

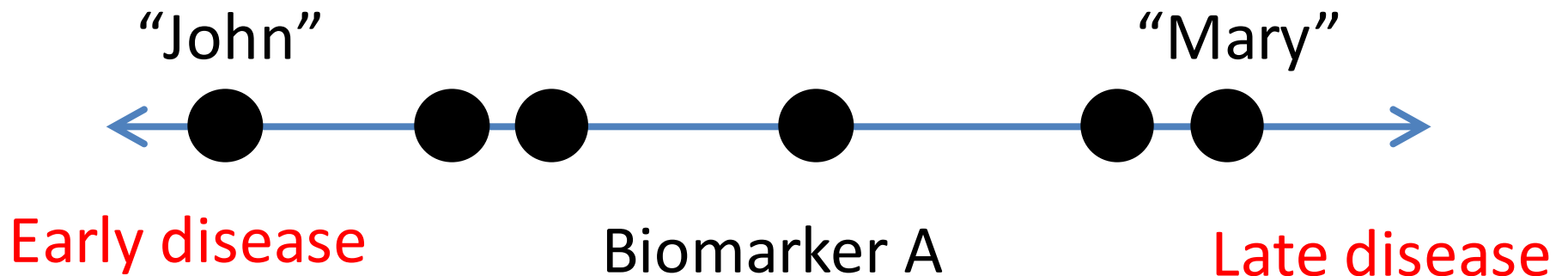
## 2. Simultaneous staging & subtyping

- Young et al., *Nature Communications* 2018

# Stage vs. subtype

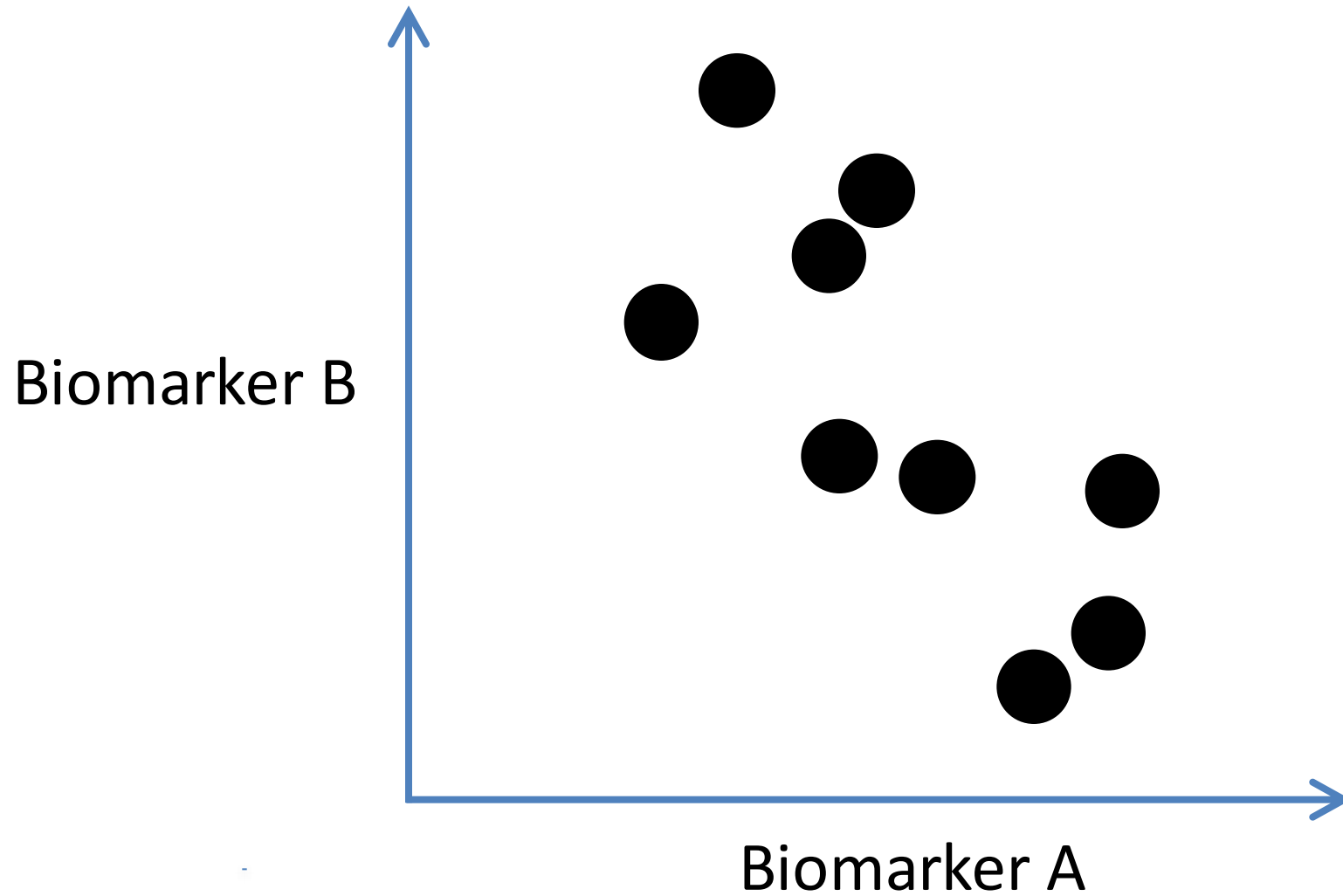
- Staging: sort patients into early-late disease or severity, i.e. discover the trajectory
- Cross-sectional data: only 1 time point observed per patient
  - More generally, censored to be a short window
- Naïve clustering can't differentiate between *stage* and *subtype*
  - Patients assumed to be aligned at baseline
- **Let's build some intuition around how staging from cross-sectional data might be possible...**

In 1-D, might assume that low values correspond to an early disease stage (or vice-versa)

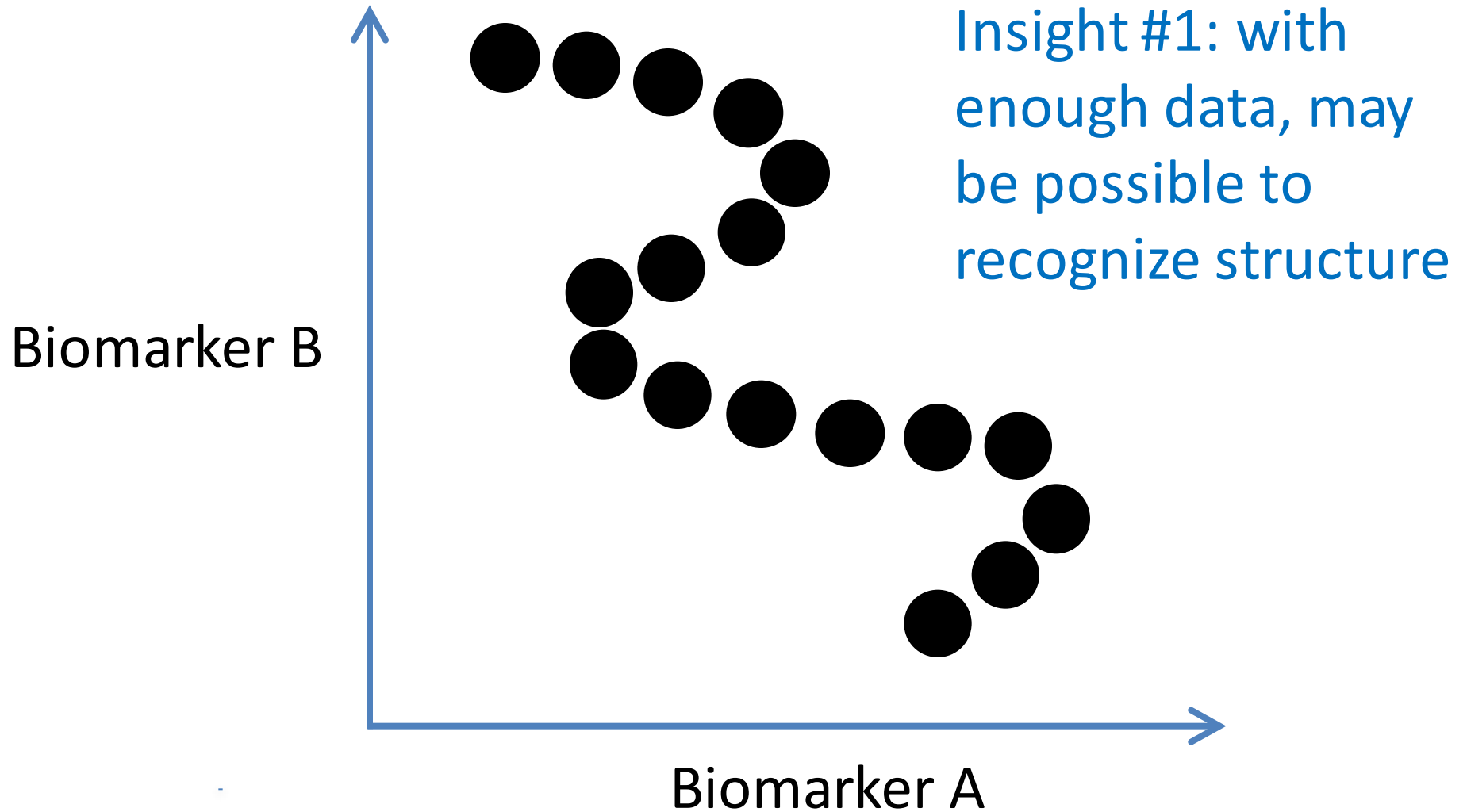


*Assume samples were all taken today*

# What about in higher dimensions?



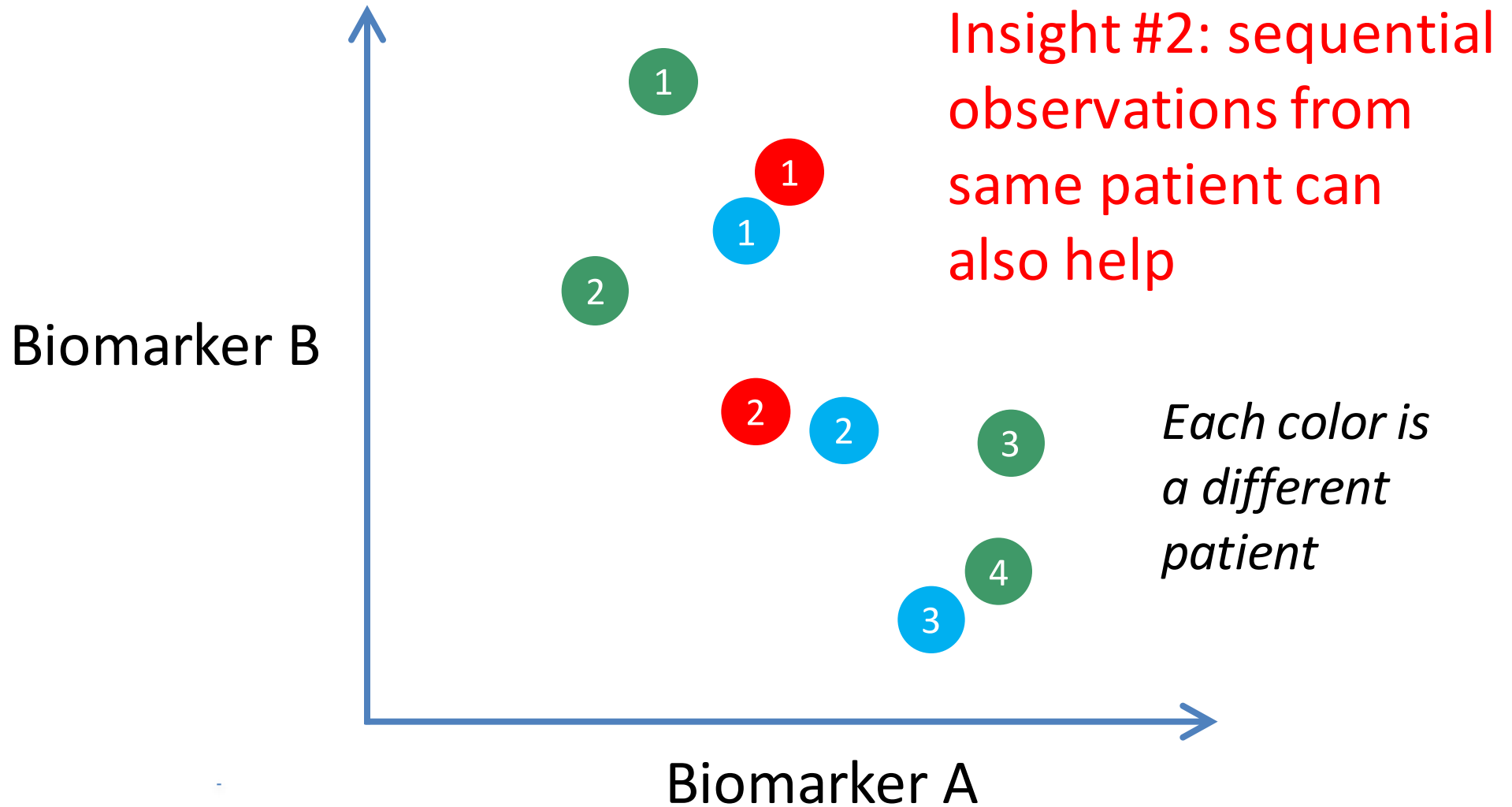
# What about in higher dimensions?



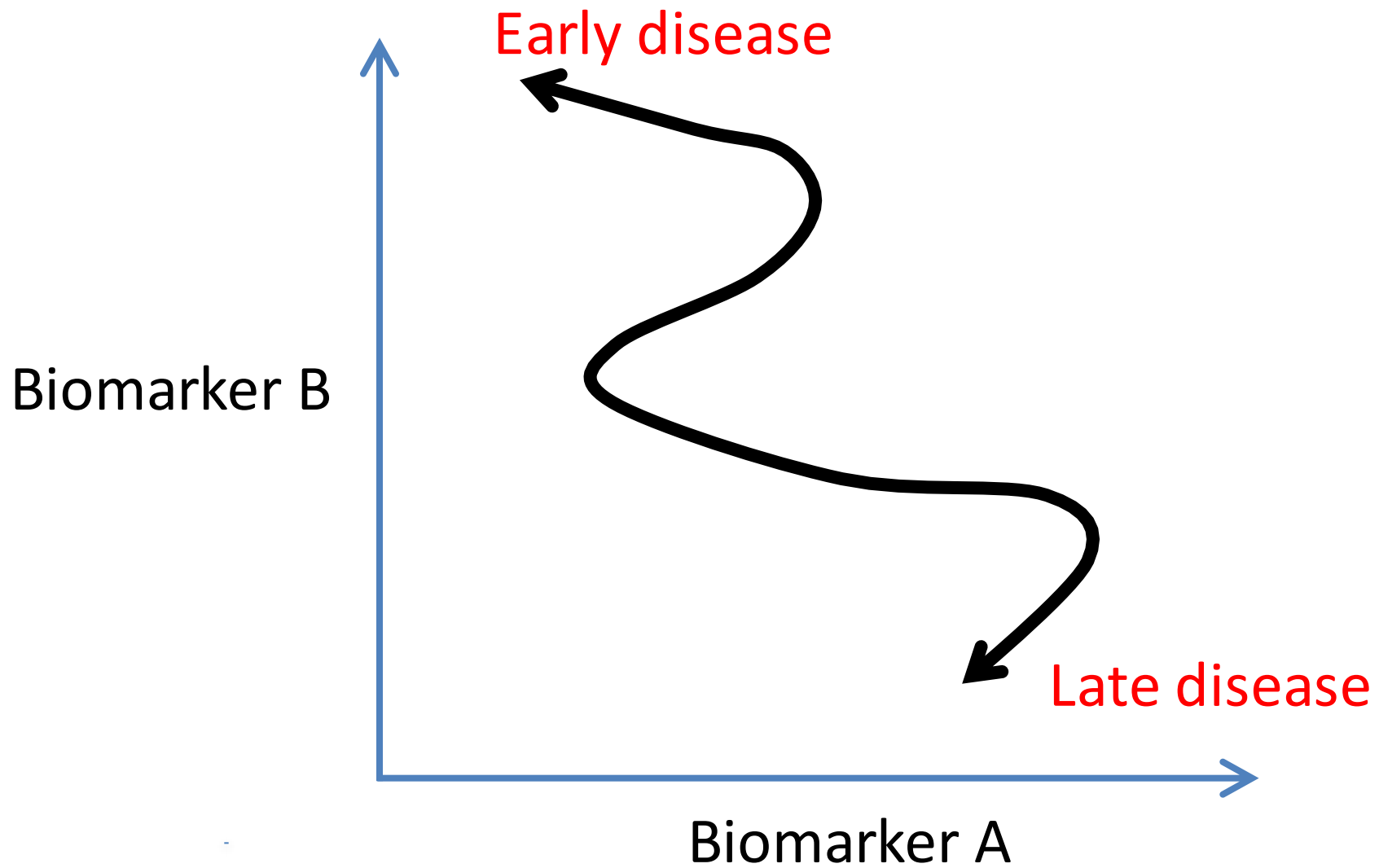
[Bendall et al., Cell 2014 (human B cell development)]



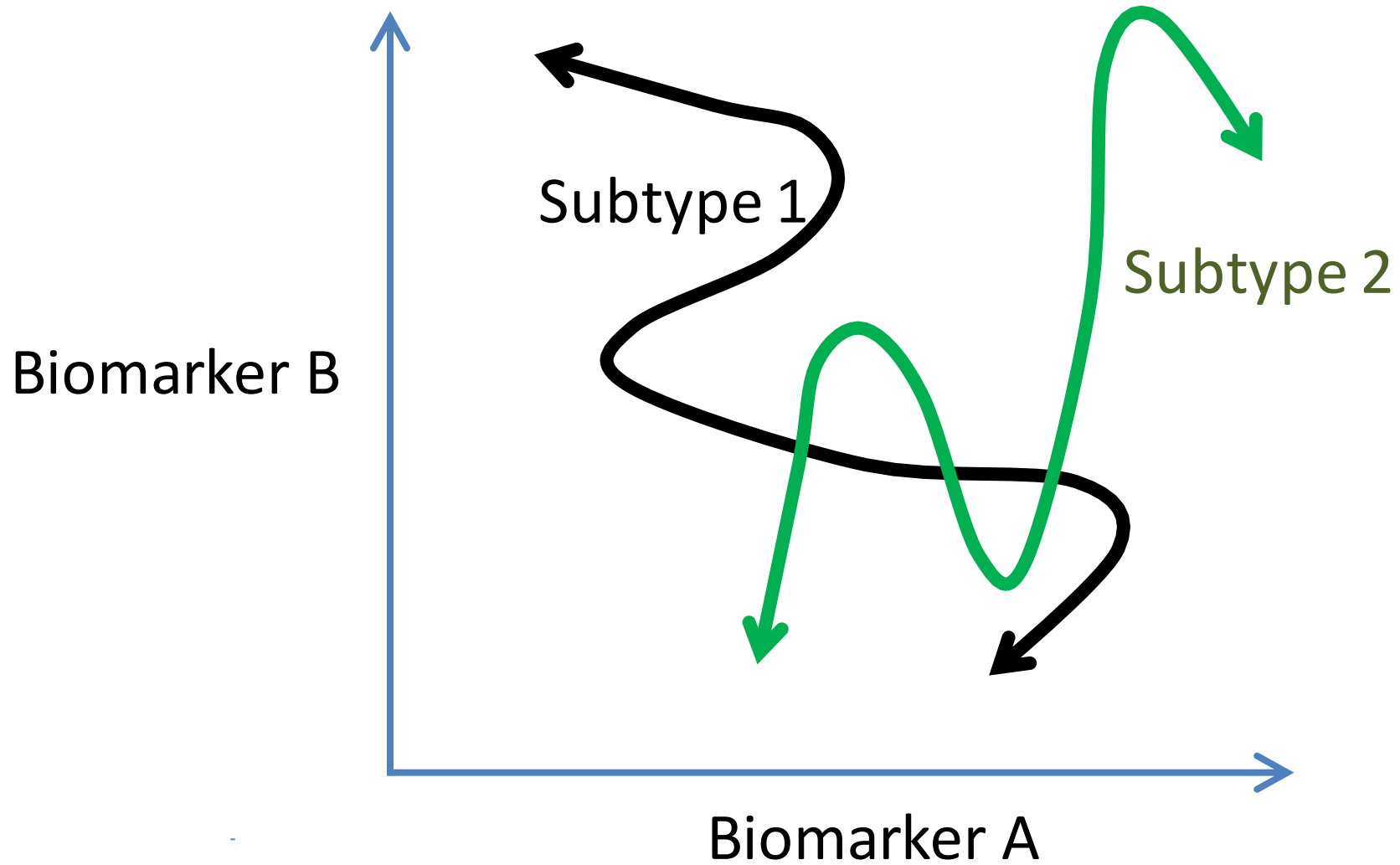
# What about in higher dimensions?



What about in higher dimensions?



May also seek to discover disease subtypes



# Outline of today's lecture

1. Staging from cross-sectional data
  - **Wang, Sontag, Wang, *KDD* 2014**
  - Pseudo-time methods from computational biology
2. Simultaneous staging & subtyping
  - Young et al., *Nature Communications* 2018

# COPD diagnosis & progression

- COPD diagnosis made using a breath test – fraction of air expelled in first second of exhalation < 70%
- Most doctors use GOLD criteria to stage the disease and measure its progression:

	1 (mild)	2 (moderate)	3 (severe)	4 (very severe)
FEV <sub>1</sub> :FVC	<0.70	<0.70	<0.70	<0.70
FEV <sub>1</sub>	≥80% of predicted	50–80% of predicted	30–50% of predicted	<30% of predicted or <50% of predicted plus chronic respiratory failure
Treatment	Influenza vaccination and short-acting bronchodilator* when needed	Influenza vaccination, short-acting and ≥1 long-acting bronchodilator* when needed; consider respiratory rehabilitation	Influenza vaccination and short-acting and ≥1 long-acting bronchodilator* when needed, inhaled glucocorticosteroid if repeated exacerbations; consider respiratory rehabilitation	Influenza vaccination and short-acting and ≥1 long-acting bronchodilator* when needed, inhaled glucocorticosteroid if repeated exacerbations, long-term oxygen if chronic respiratory failure occurs; consider respiratory rehabilitation and surgery

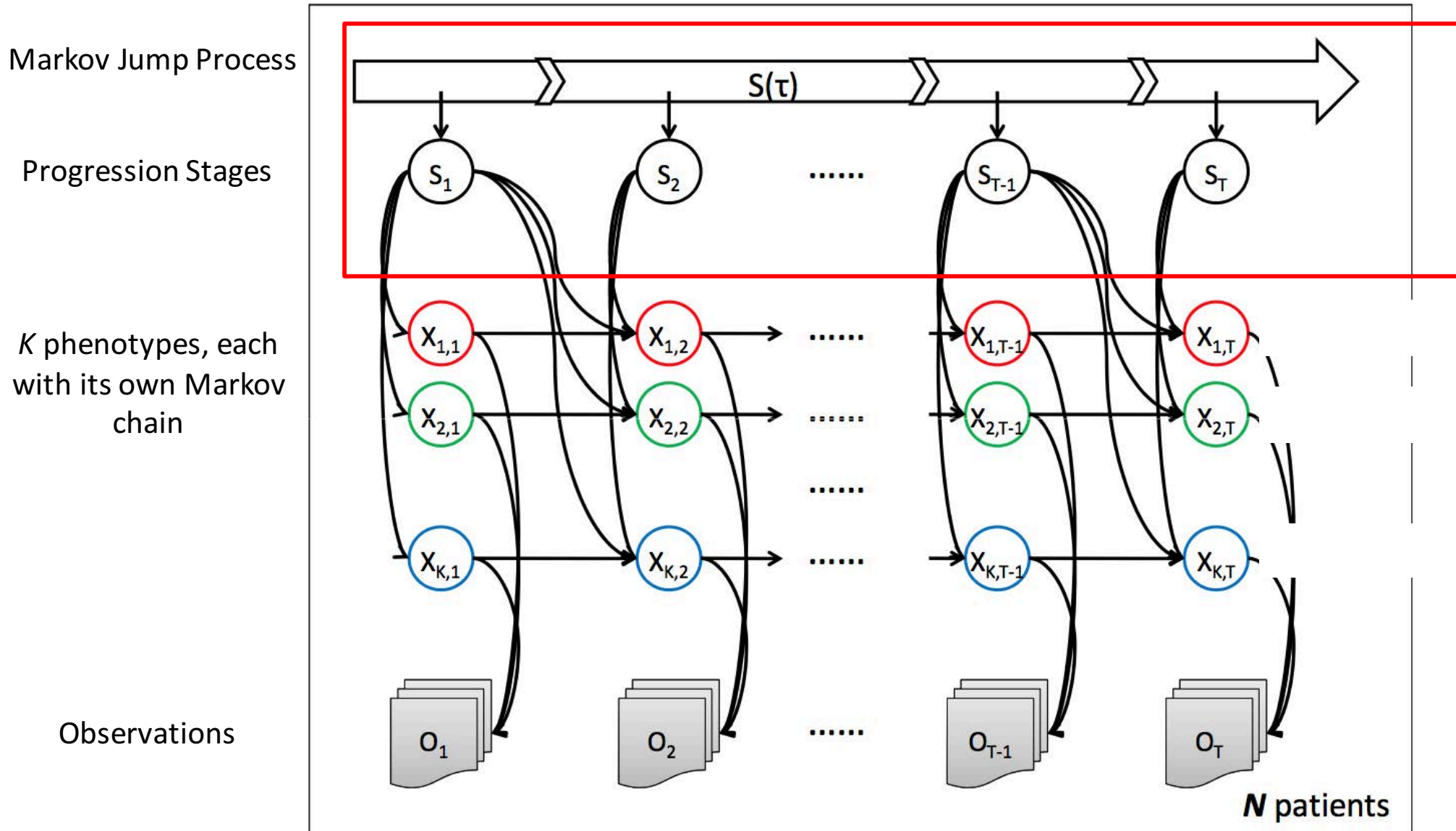
GOLD=Global Initiative on Obstructive Lung Disease. \*β<sub>2</sub> agonists or anticholinergics.

**Table: Therapy at each stage of chronic obstructive pulmonary disease, by GOLD stage<sup>1</sup>**

Chronic obstructive pulmonary disease. *The Lancet*, Volume 379, Issue 9823, Pages 1341 - 1351, 7 April 2012

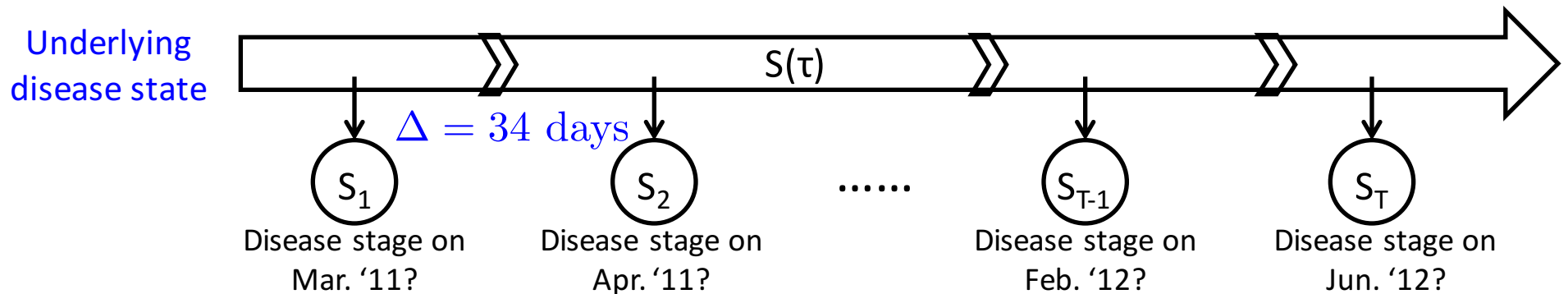
Courtesy of [Elsevier Ltd.](#) Used with permission.

# The big picture: generative model for patient data



[Wang, Sontag, Wang, “Unsupervised learning of Disease Progression Models”, KDD 2014]

# Model for patient's disease progression across time

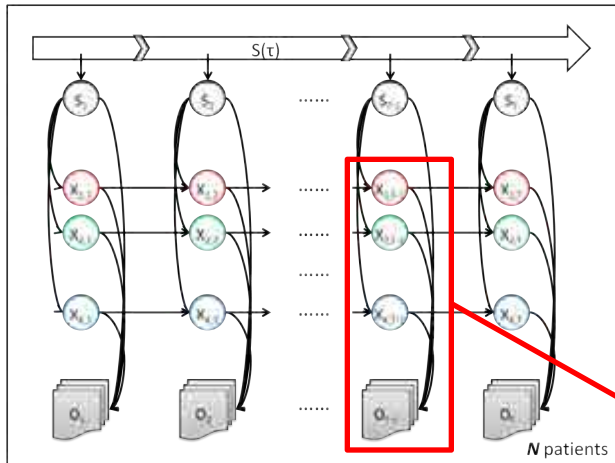


- A continuous-time Markov process with irregular discrete-time observations
- The transition probability is defined by an intensity matrix and the time interval:

$$A_{ij}(\Delta) \triangleq P(S_t = j | S_{t-1} = i, \tau_t - \tau_{t-1} = \Delta; Q) \\ = \text{expm}(\Delta Q)_{ij},$$

**Matrix Q: Parameters to learn**

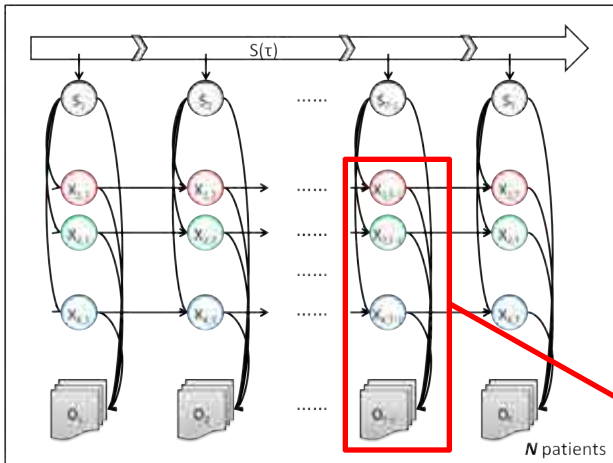
# Model for data at single point in time: Noisy-OR network



Previously used for medical diagnosis, e.g. QMR-DT (Shwe et al. '91)



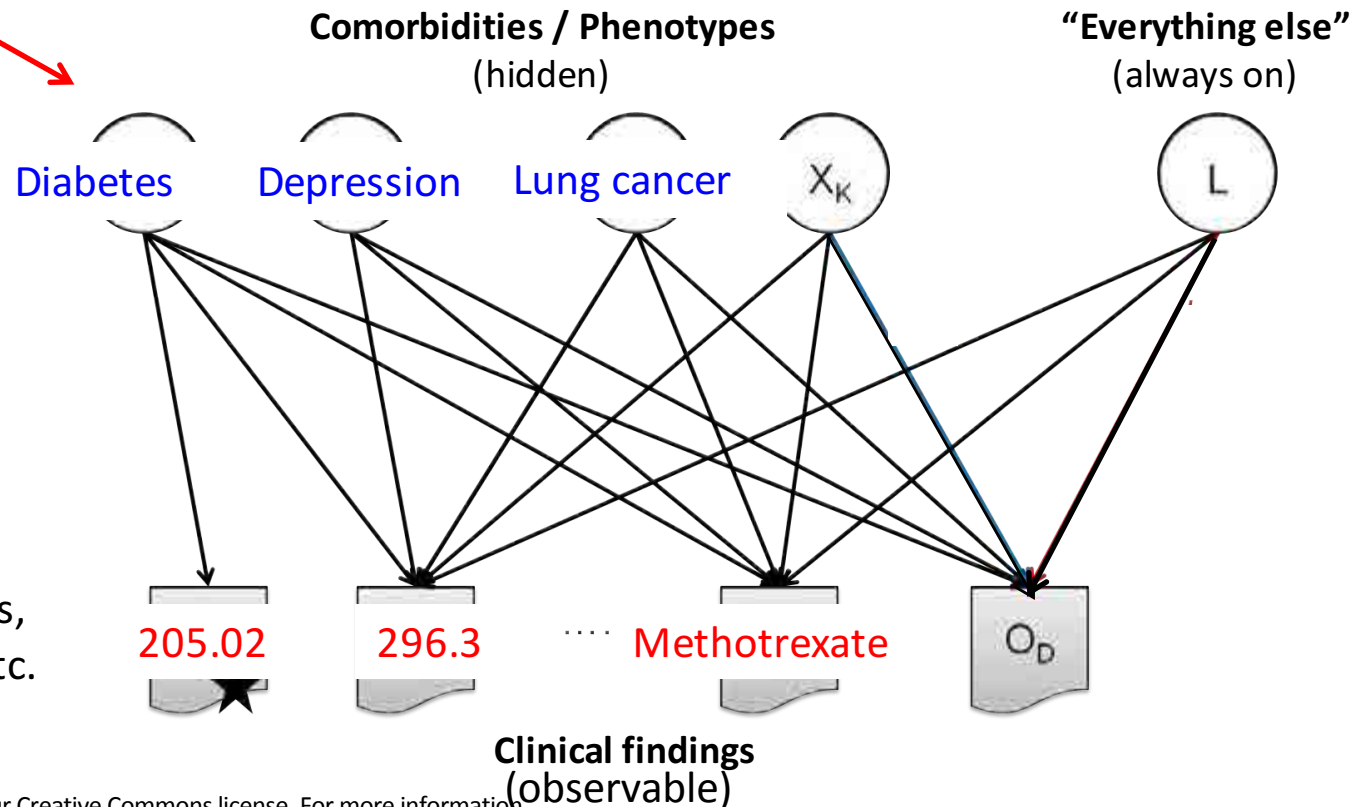
# Model for data at single point in time: Noisy-OR network



Previously used for medical diagnosis, e.g. QMR-DT (Shwe et al. '91)

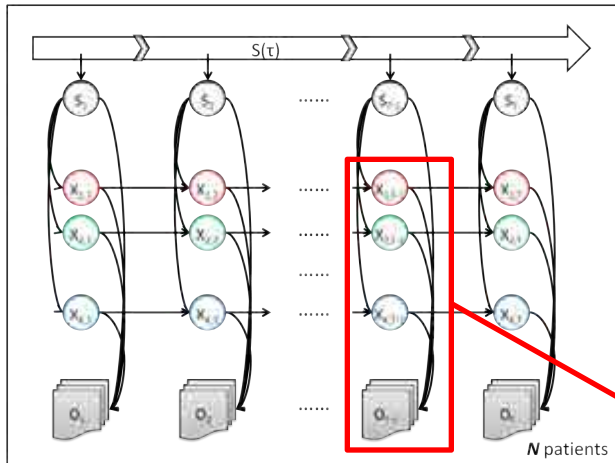
All binary variables

Diagnosis codes,  
medications, etc.



Clinical findings  
(observable)

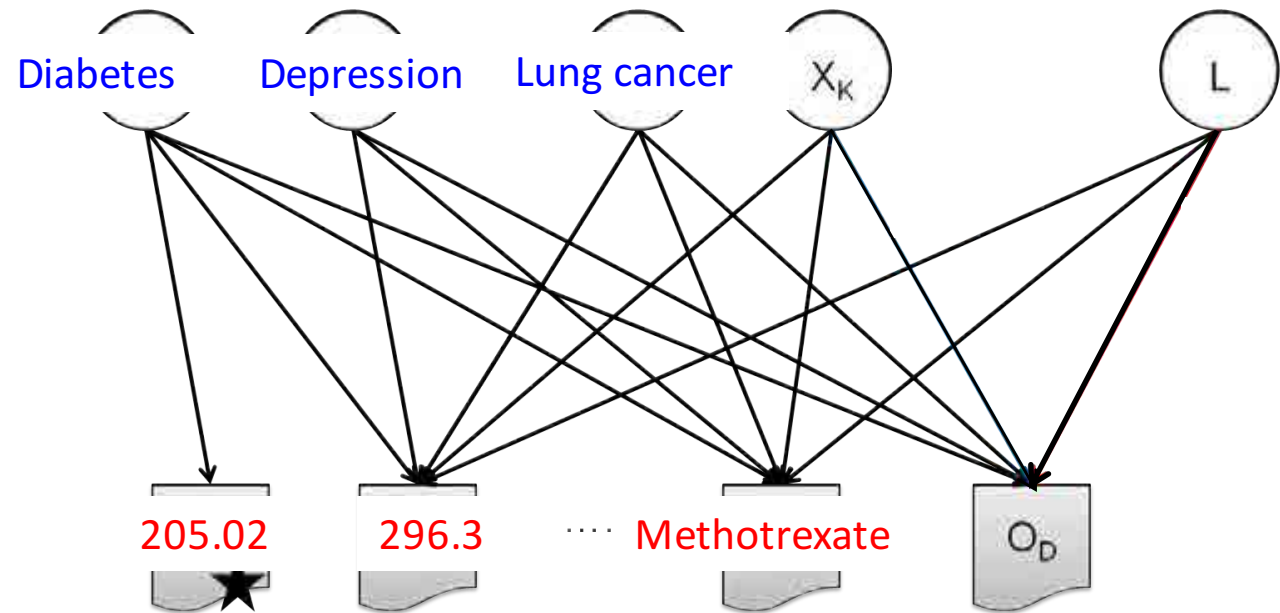
# Model for data at single point in time: Noisy-OR network



Previously used for medical diagnosis, e.g. QMR-DT (Shwe et al. '91)

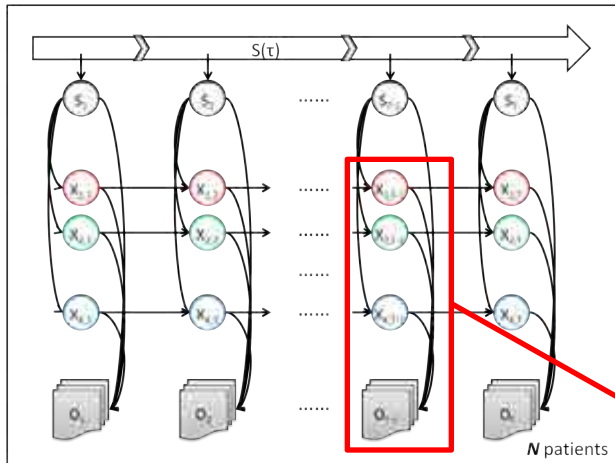
Comorbidities / Phenotypes  
(hidden)

"Everything else"  
(always on)

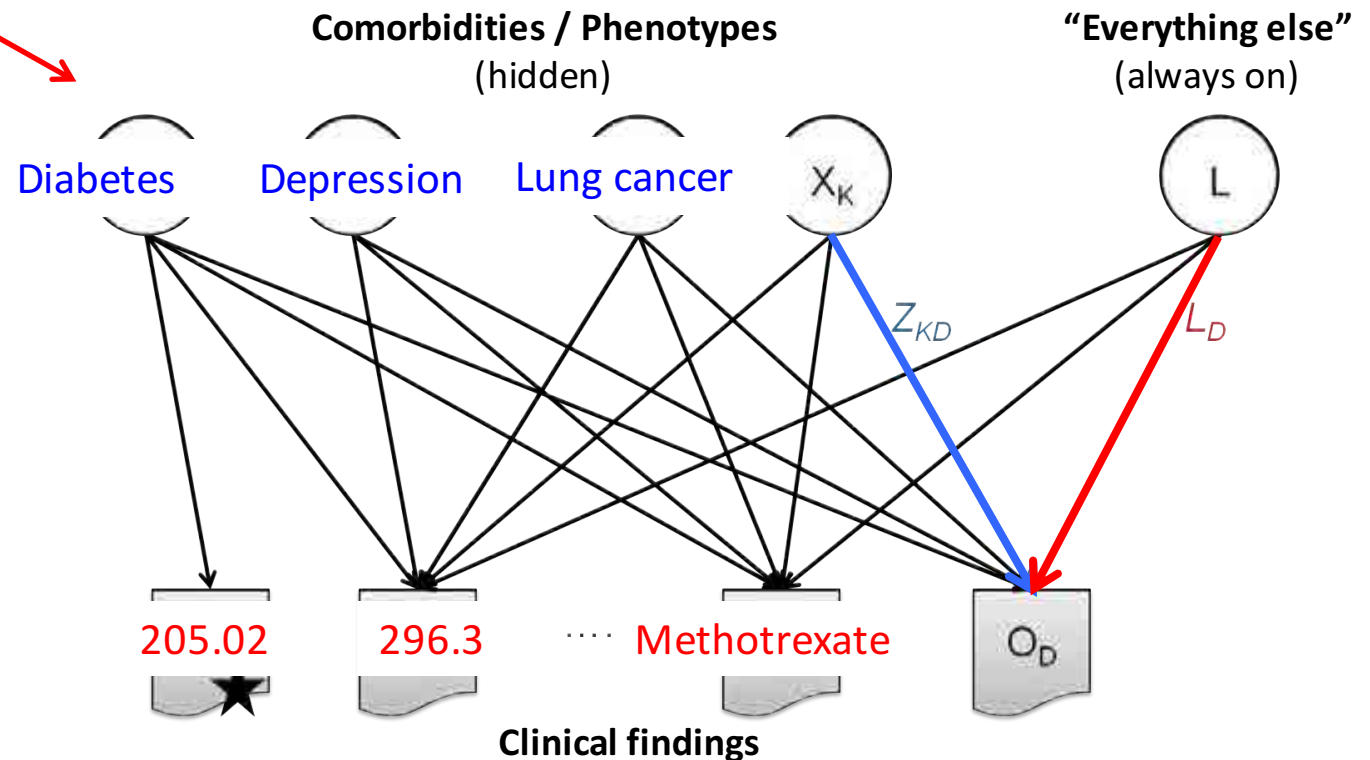


We also learn which edges exist

# Model for data at single point in time: Noisy-OR network



Previously used for medical diagnosis, e.g. QMR-DT (Shwe et al. '91)

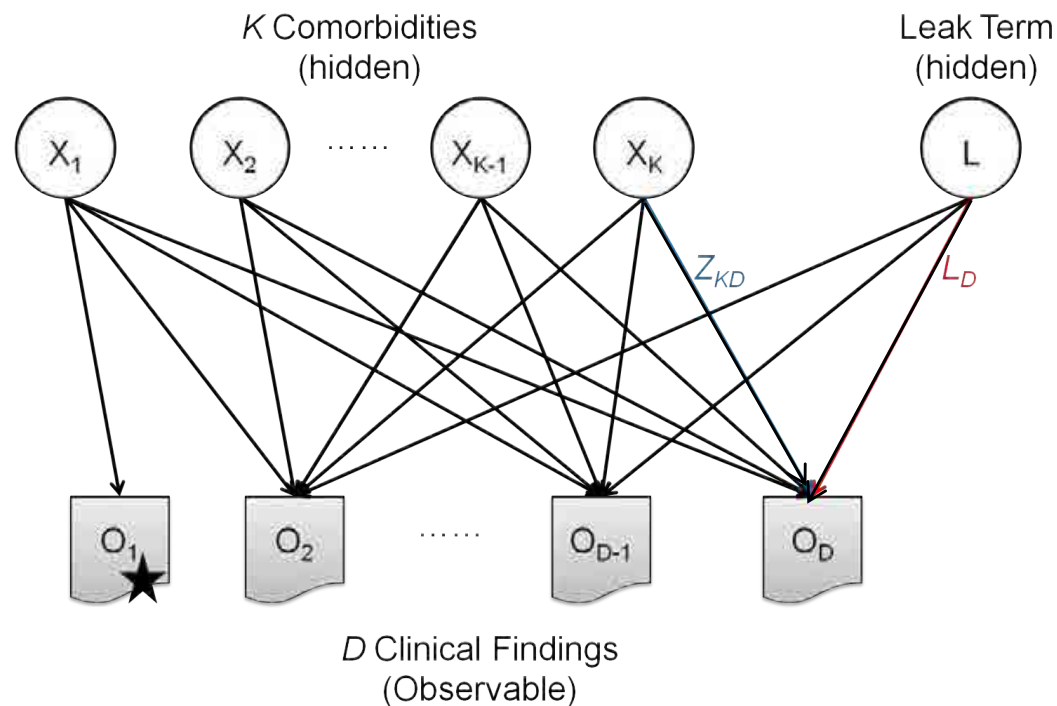


We also learn which edges exist

Associated with each edge is a *failure probability*

# Using anchors to ground the hidden variables

- An *anchor* is a finding that can only be caused by a single comorbidity (discussed in Lecture 8)



Y. Halpern, YD Choi, S. Horng, D. Sontag. Using Anchors to Estimate Clinical State without Labeled Data. To appear in the American Medical Informatics Association (AMIA) Annual Symposium, Nov. 2014

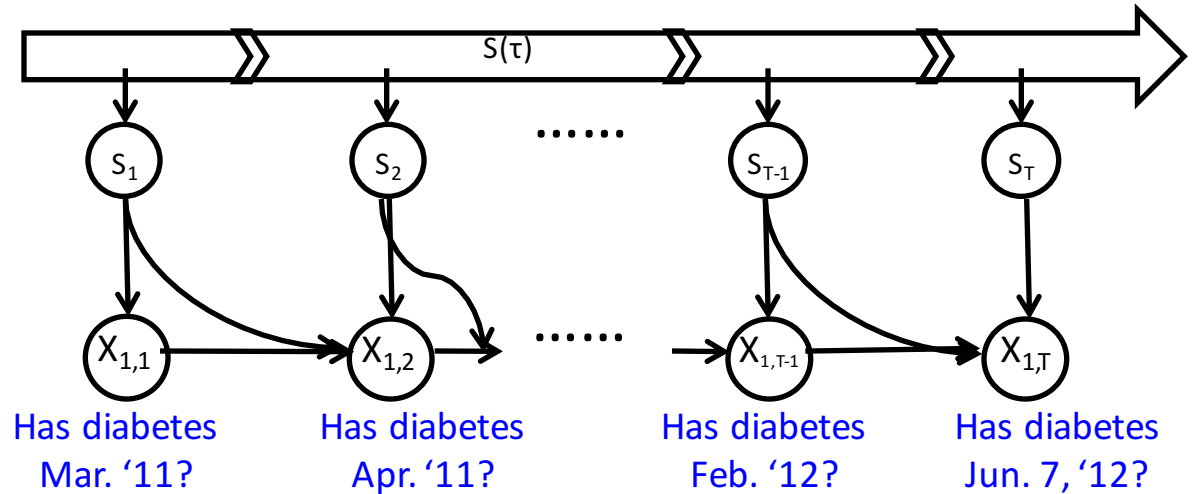
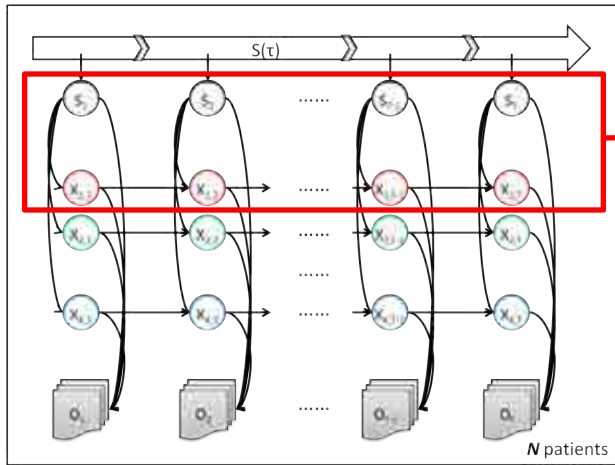
# Using anchors to ground the hidden variables

- Provide anchors for each of the comorbidities:

Comorbidity	Representative Conditions (Anchor ICD-9 Codes)
COPD	Chronic Bronchitis (491), Emphysema (492, 518), Chronic Airway Obstruction (496)
Asthma	Asthma (493)
Cardiovascular	Hypertension (401), Congestive Heart Failure (428), Arrhythmia (427), Ischemic Heart Disease (414)
Lung Infection	Pneumonia (481, 485, 486)
Lung Cancer	Malignant Neoplasm of Upper/Lower Lobe, Bronchus or Lung (162)
Diabetes	Diabetes with Different Types and Complications (250)
Musculoskeletal	Spinal Disorders (724), Soft Tissue Disorders (729), Osteoporosis (733)
Kidney	Acute Kidney Failure (584), Chronic Kidney Disease (585), Renal Failure (586)
Psychological	Anxiety (300), Depression (296, 311)
Obesity	Morbid Obesity (278)

- Can be viewed as a type of weak supervision, using clinical domain knowledge
- Without these, the results are less interpretable

# Model of comorbidities across time



- Presence of comorbidities depends on value at previous time step and on disease stage
- Later stages of disease = more likely to develop comorbidities
- Make the assumption that once patient has a comorbidity, likely to always have it

# Experimental evaluation

- We create a COPD cohort of 3,705 patients:
  - At least one COPD-related diagnosis code
  - At least one COPD-related drug
- Removed patients with too few records
- Clinical findings derived from 264 diagnosis codes
  - Removed ICD-9 codes that only occurred to a small number of patients
- Combined visits into 3-month time windows
- 34,976 visits, 189,815 positive findings

# Inference

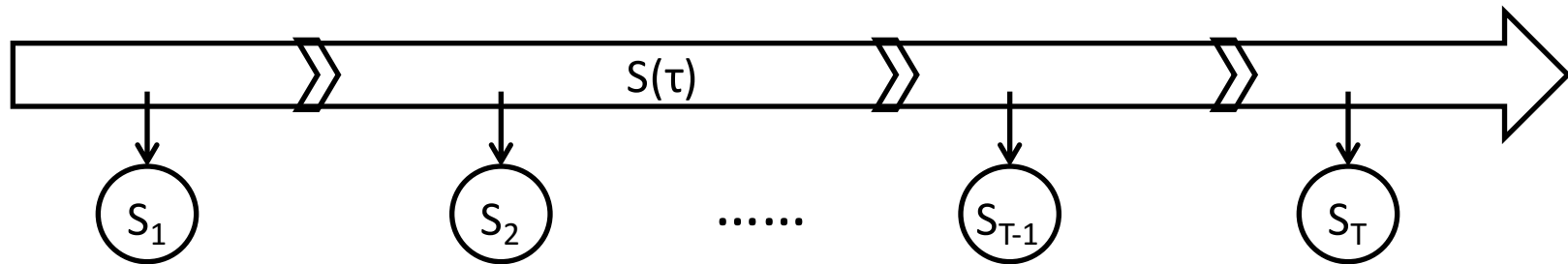
- Outer loop
  - EM
  - Algorithm to estimate the Markov Jump Process is borrowed from recent literature in physics
- Inner loop
  - Gibbs sampler used for approximate inference
  - Perform block sampling of the Markov chains, improving the mixing time of the Gibbs sampler
- **If I were to do it again... would do variational inference with a recognition network (as in VAEs)**

*P. Metzner, I. Horenko, and C. Schutte. Generator estimation of markov jump processes based on incomplete observations nonequidistant in time. Physical Review E, 76(6):066702, 2007.*



# Customizations for COPD

- Enforce monotonic stage progression, i.e.  $S_{t+1} \geq S_t$ :



- Enforce monotonicity in distributions of comorbidities in first time step, e.g.  $\Pr(X_{j,1} | S_1 = 2) \geq \Pr(X_{j,1} | S_1 = 1)$ 
  - To do this, we solve a tiny convex optimization problem within EM
- Enforce that transitions in  $X$  can only happen at the same time as transitions in  $S$
- Edge weights given a  $\text{Beta}(0.1, 1)$  prior to encourage sparsity

# Edges learned for *kidney disease*

<u>Diagnosis code</u>	<u>Weight</u>	
*585.3	0.20	Chronic Kidney Disease, Stage Iii (Moderate)
285.9	0.15	Anemia, Unspecified
*585.9	0.10	Chronic Kidney Disease, Unspecified
599.0	0.08	Urinary Tract Infection, Site Not Specified
*585.4	0.08	Chronic Kidney Disease, Stage Iv (Severe)
*584.9	0.07	Acute Renal Failure, Unspecified
*586	0.07	Renal Failure, Unspecified
782.3	0.06	Edema
*585.6	0.05	End Stage Renal Disease
593.9	0.04	Unspecified Disorder Of Kidney And Ureter
272.4	0.04	Other And Unspecified Hyperlipidemia
272.2	0.03	Mixed Hyperlipidemia

# Edges learned for *kidney disease*

<u>Diagnosis code</u>	<u>Weight</u>	
<b>*585.3</b>	<b>0.20</b>	<b>Chronic Kidney Disease, Stage Iii (Moderate)</b>
285.9	0.15	Anemia, Unspecified
<b>*585.9</b>	<b>0.10</b>	<b>Chronic Kidney Disease, Unspecified</b>
599.0	0.08	Urinary Tract Infection, Site Not Specified
<b>*585.4</b>	<b>0.08</b>	<b>Chronic Kidney Disease, Stage Iv (Severe)</b>
<b>*584.9</b>	<b>0.07</b>	<b>Acute Renal Failure, Unspecified</b>
<b>*586</b>	<b>0.07</b>	<b>Renal Failure, Unspecified</b>
782.3	0.06	Edema
<b>*585.6</b>	<b>0.05</b>	<b>End Stage Renal Disease</b>
593.9	0.04	Unspecified Disorder Of Kidney And Ureter
272.4	0.04	Other And Unspecified Hyperlipidemia
272.2	0.03	Mixed Hyperlipidemia

# Edges learned for *kidney disease*

<u>Diagnosis code</u>	<u>Weight</u>		
*585.3	0.20	Chronic Kidney Disease, Stage Iii (Moderate)	
<b>285.9</b>	<b>0.15</b>	<b>Anemia, Unspecified</b>	<p><b>Why do people with kidney disease get anemia?</b></p> <p>Your kidneys make an important hormone called <i>erythropoietin (EPO)</i>. Hormones are secretions that your body makes to help your body work and keep you healthy. EPO tells your body to make red blood cells. When you have kidney disease, your kidneys cannot make enough EPO. This causes your red blood cell count to drop and anemia to develop.</p>
*585.9	0.10	Chronic Kidney Disease,	
<b>599.0</b>	<b>0.08</b>	<b>Urinary Tract Infection,</b>	
*585.4	0.08	Chronic Kidney Disease,	
*584.9	0.07	Acute Renal Failure, Uns	
*586	0.07	Renal Failure, Unspecifie	
<b>782.3</b>	<b>0.06</b>	<b>Edema</b>	
*585.6	0.05	End Stage Renal Disease	
<b>593.9</b>	<b>0.04</b>	<b>Unspecified Disorder Of</b>	
<b>272.4</b>	<b>0.04</b>	<b>Other And Unspecified  </b>	
<b>272.2</b>	<b>0.03</b>	<b>Mixed Hyperlipidemia</b>	

[WWW.KIDNEY.ORG](http://WWW.KIDNEY.ORG)

# Edges learned for *lung cancer*

<u>Diagnosis code</u>	<u>Weight</u>	
*162.9	0.60	Malignant Neoplasm Of Bronchus And Lung
518.89	0.15	Other Diseases Of Lung, Not Elsewhere Classified
*162.8	0.15	Malignant Neoplasm Of Other Parts Of Lung
*162.3	0.15	Malignant Neoplasm Of Upper Lobe, Lung
786.6	0.15	Swelling, Mass, Or Lump In Chest
793.1	0.10	Abnormal Findings On Radiological Exam Of Lung
786.09	0.07	Other Respiratory Abnormalities
*162.5	0.06	Malignant Neoplasm Of Lower Lobe, Lung
*162.2	0.04	Malignant Neoplasm Of Main Bronchus
702.0	0.03	Actinic Keratosis
511.9	0.03	Unspecified Pleural Effusion
*162.4	0.03	Malignant Neoplasm Of Middle Lobe, Lung

# Edges learned for *lung cancer*

<u>Diagnosis code</u>	<u>Weight</u>	
<b>*162.9</b>	<b>0.60</b>	<b>Malignant Neoplasm Of Bronchus And Lung</b>
518.89	0.15	Other Diseases Of Lung, Not Elsewhere Classified
<b>*162.8</b>	<b>0.15</b>	<b>Malignant Neoplasm Of Other Parts Of Lung</b>
<b>*162.3</b>	<b>0.15</b>	<b>Malignant Neoplasm Of Upper Lobe, Lung</b>
786.6	0.15	Swelling, Mass, Or Lump In Chest
793.1	0.10	Abnormal Findings On Radiological Exam Of Lung
786.09	0.07	Other Respiratory Abnormalities
<b>*162.5</b>	<b>0.06</b>	<b>Malignant Neoplasm Of Lower Lobe, Lung</b>
<b>*162.2</b>	<b>0.04</b>	<b>Malignant Neoplasm Of Main Bronchus</b>
702.0	0.03	Actinic Keratosis
511.9	0.03	Unspecified Pleural Effusion
<b>*162.4</b>	<b>0.03</b>	<b>Malignant Neoplasm Of Middle Lobe, Lung</b>

# Edges learned for *lung cancer*

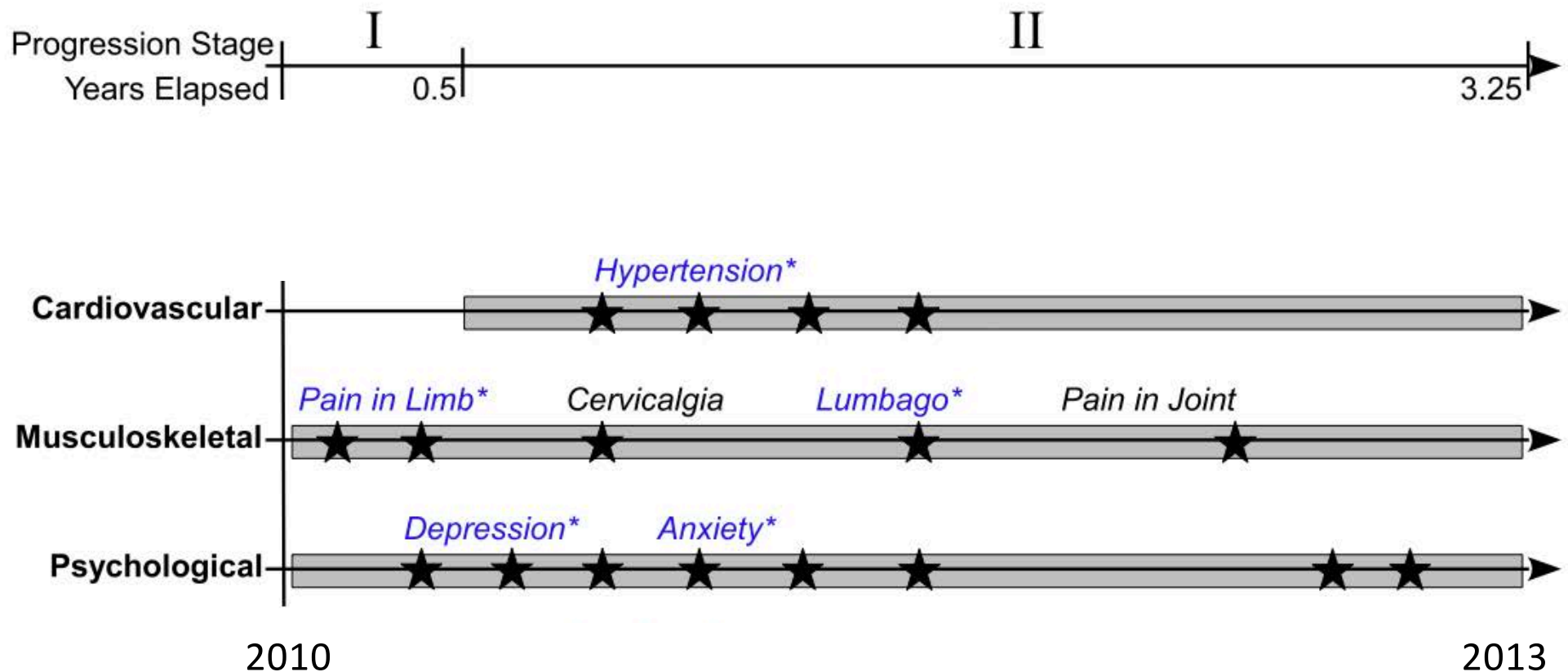
<u>Diagnosis code</u>	<u>Weight</u>	
*162.9	0.60	Malignant Neoplasm Of Bronchus And Lung
<b>518.89</b>	<b>0.15</b>	<b>Other Diseases Of Lung, Not Elsewhere Classified</b>
*162.8	0.15	Malignant Neoplasm Of Other Parts Of Lung
*162.3	0.15	Malignant Neoplasm Of Upper Lobe, Lung
<b>786.6</b>	<b>0.15</b>	<b>Swelling, Mass, Or Lump In Chest</b>
<b>793.1</b>	<b>0.10</b>	<b>Abnormal Findings On Radiological Exam Of Lung</b>
<b>786.09</b>	<b>0.07</b>	<b>Other Respiratory Abnormalities</b>
*162.5	0.06	Malignant Neoplasm Of Lower Lobe, Lung
*162.2	0.04	Malignant Neoplasm Of Main Bronchus
<b>702.0</b>	<b>0.03</b>	<b>Actinic Keratosis</b>
<b>511.9</b>	<b>0.03</b>	<b>Unspecified Pleural Effusion</b>
*162.4	0.03	Malignant Neoplasm Of Middle Lobe, Lung

# Edges learned for *lung infection*

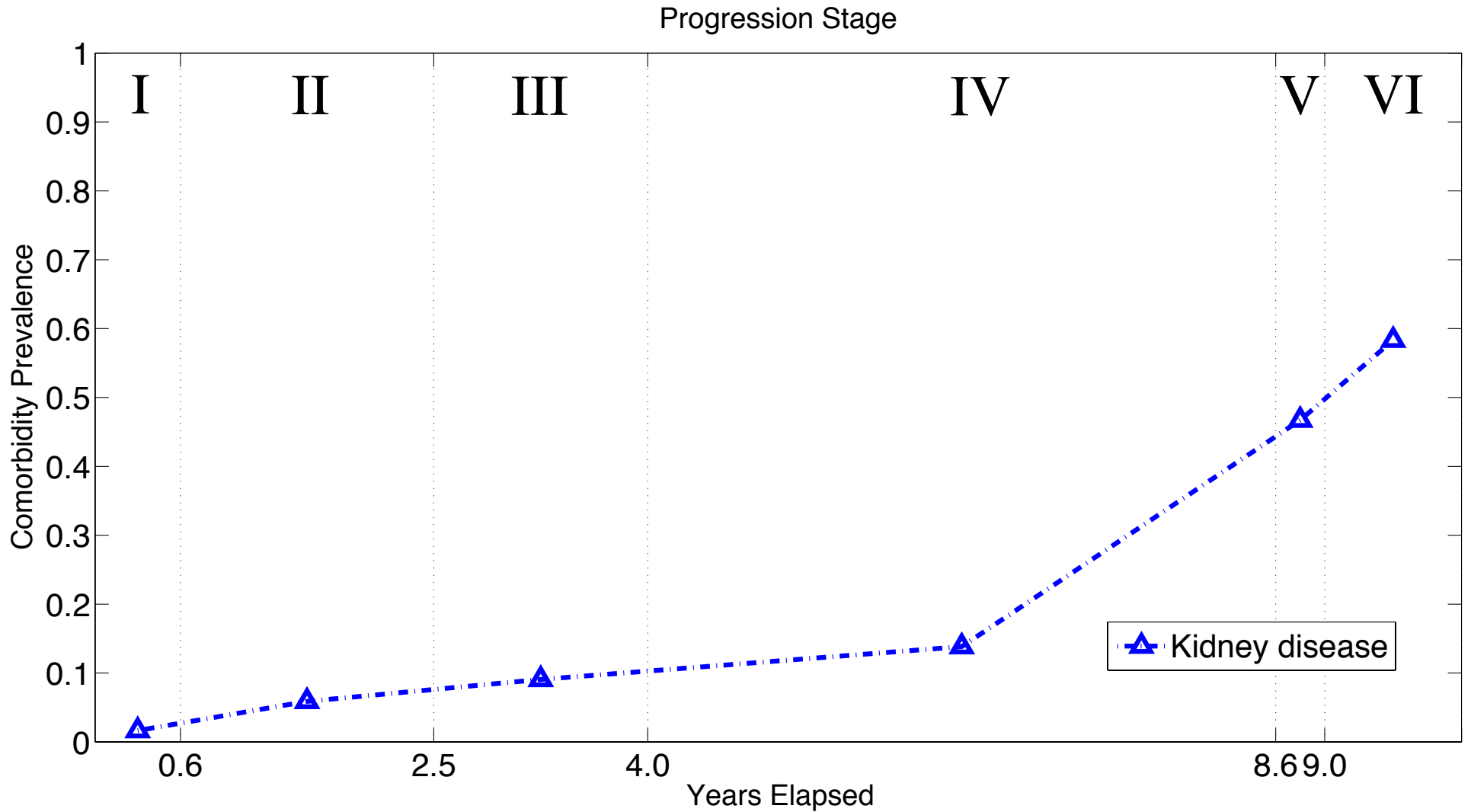
<u>Diagnosis code</u>	<u>Weight</u>	
<b>*486</b>	<b>0.30</b>	<b>Pneumonia, Organism Unspecified</b>
786.05	0.10	Shortness Of Breath
786.09	0.10	Other Respiratory Abnormalities
786.2	0.10	Cough
793.1	0.06	Abnormal Findings On Radiological Exam Of Lung
285.9	0.05	Anemia, Unspecified
518.89	0.05	Other Diseases Of Lung, Not Elsewhere Classified
466.0	0.05	Acute Bronchitis
799.02	0.05	Hypoxemia
599.0	0.04	Urinary Tract Infection, Site Not Specified
V58.61	0.04	Long-Term (Current) Use Of Anticoagulants
786.50	0.04	Chest Pain, Unspecified



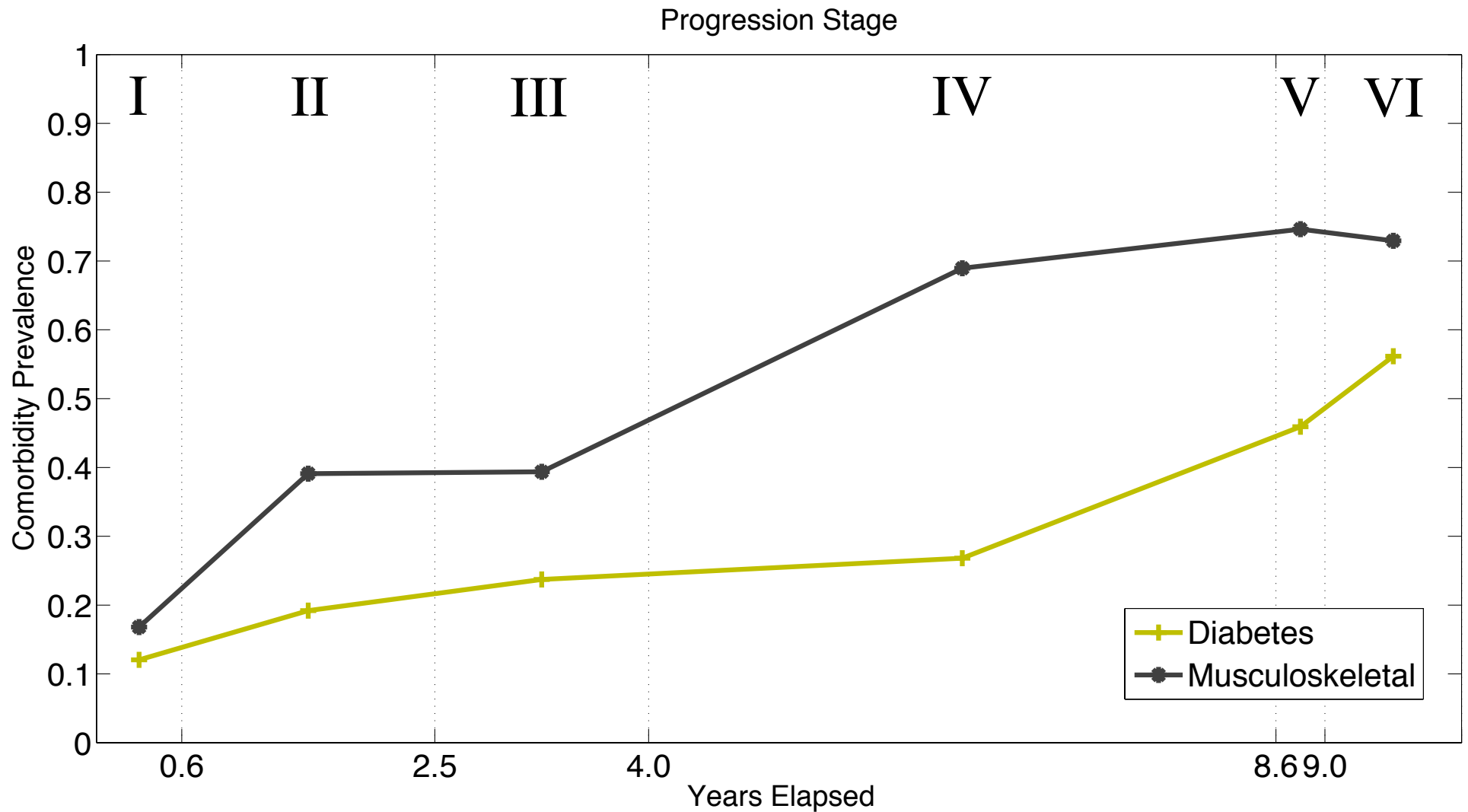
# Progression of a single patient



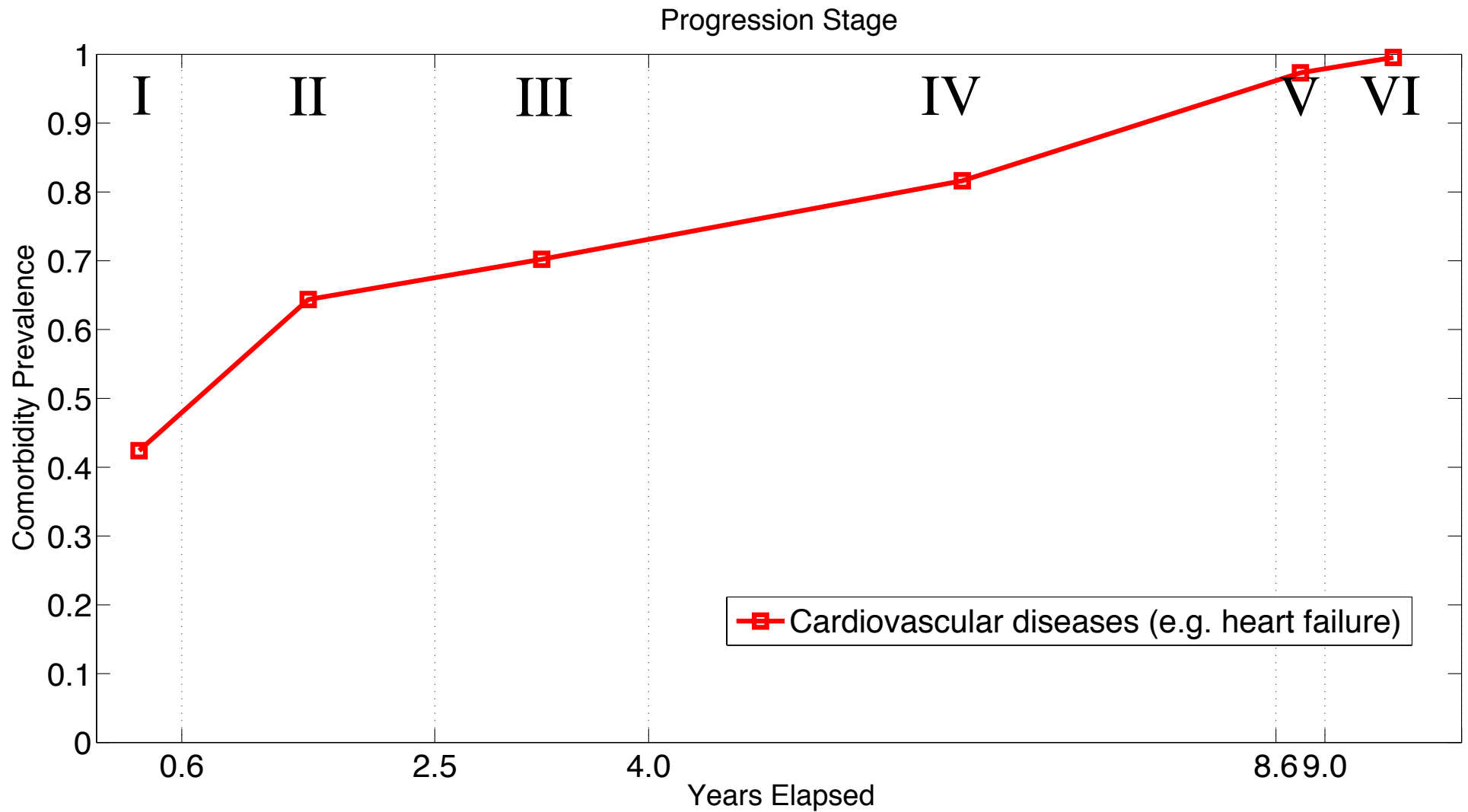
# Prevalence of comorbidities across stages (Kidney disease)



# Prevalence of comorbidities across stages (Diabetes & Musculoskeletal disorders)



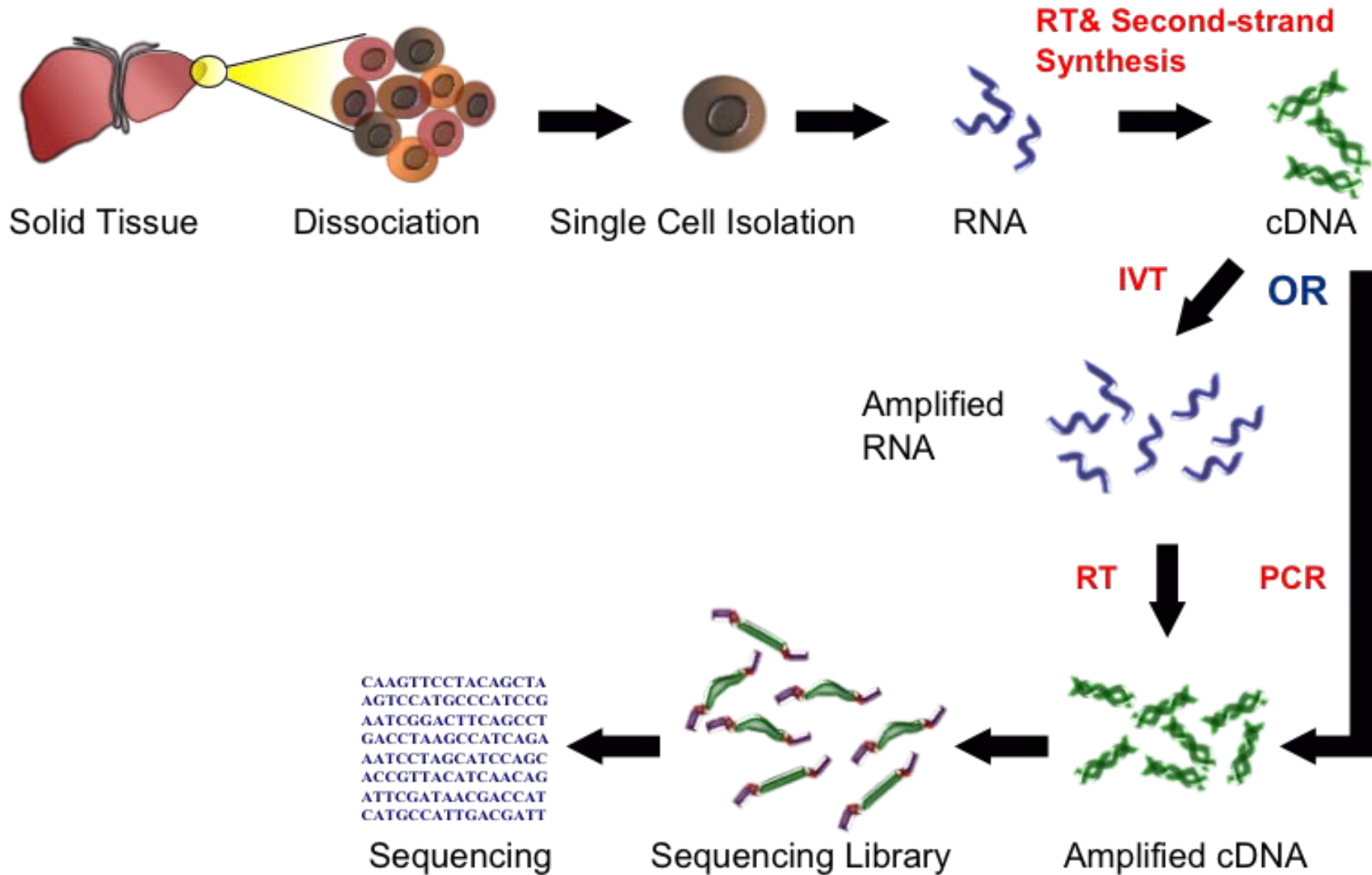
# Prevalence of comorbidities across stages (Cardiovascular disease)



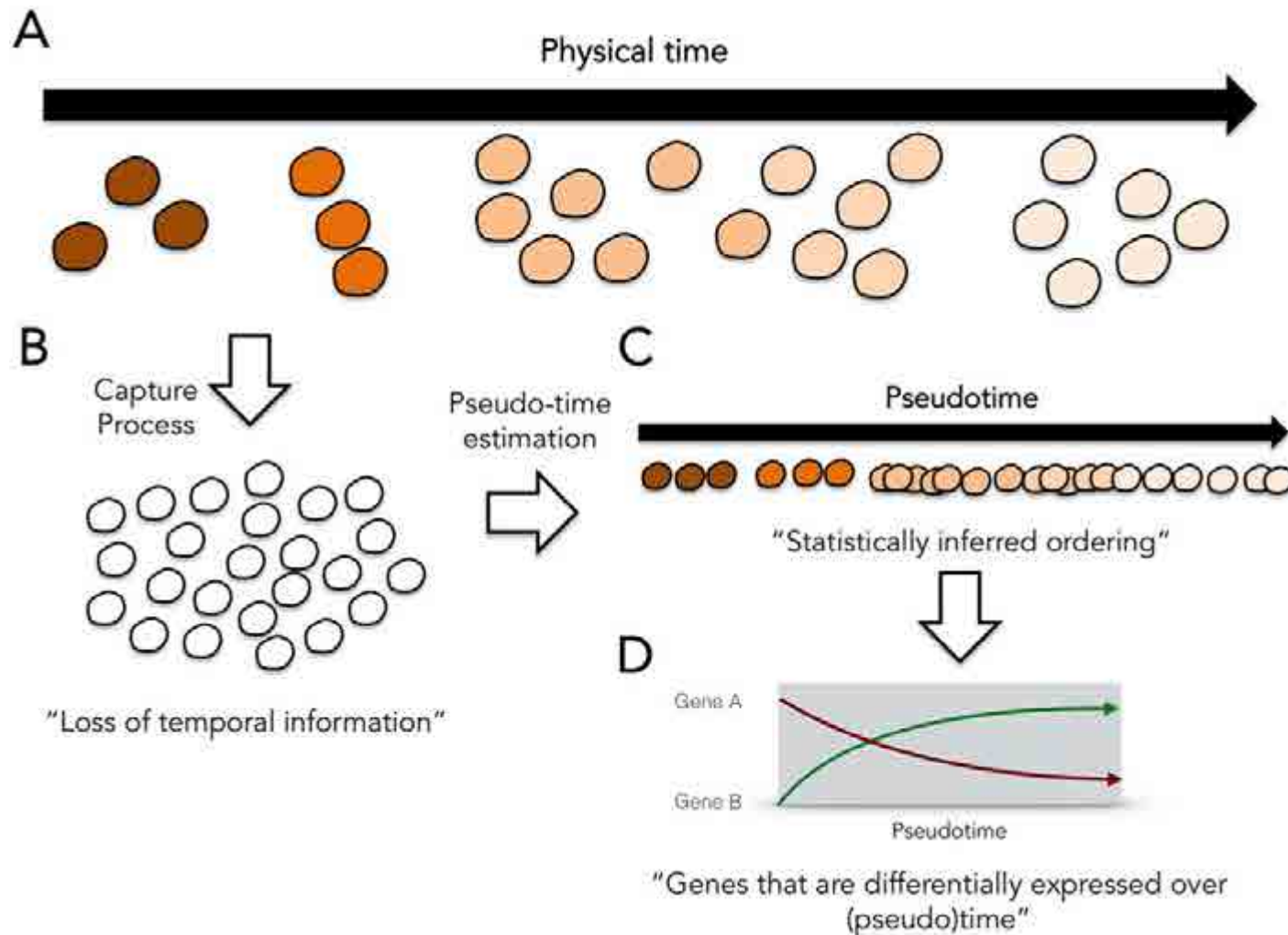
# Outline of today's lecture

1. Staging from cross-sectional data
  - Wang, Sontag, Wang, *KDD* 2014
  - **Pseudo-time methods from computational biology**
2. Simultaneous staging & subtyping
  - Young et al., *Nature Communications* 2018

# Single-cell sequencing

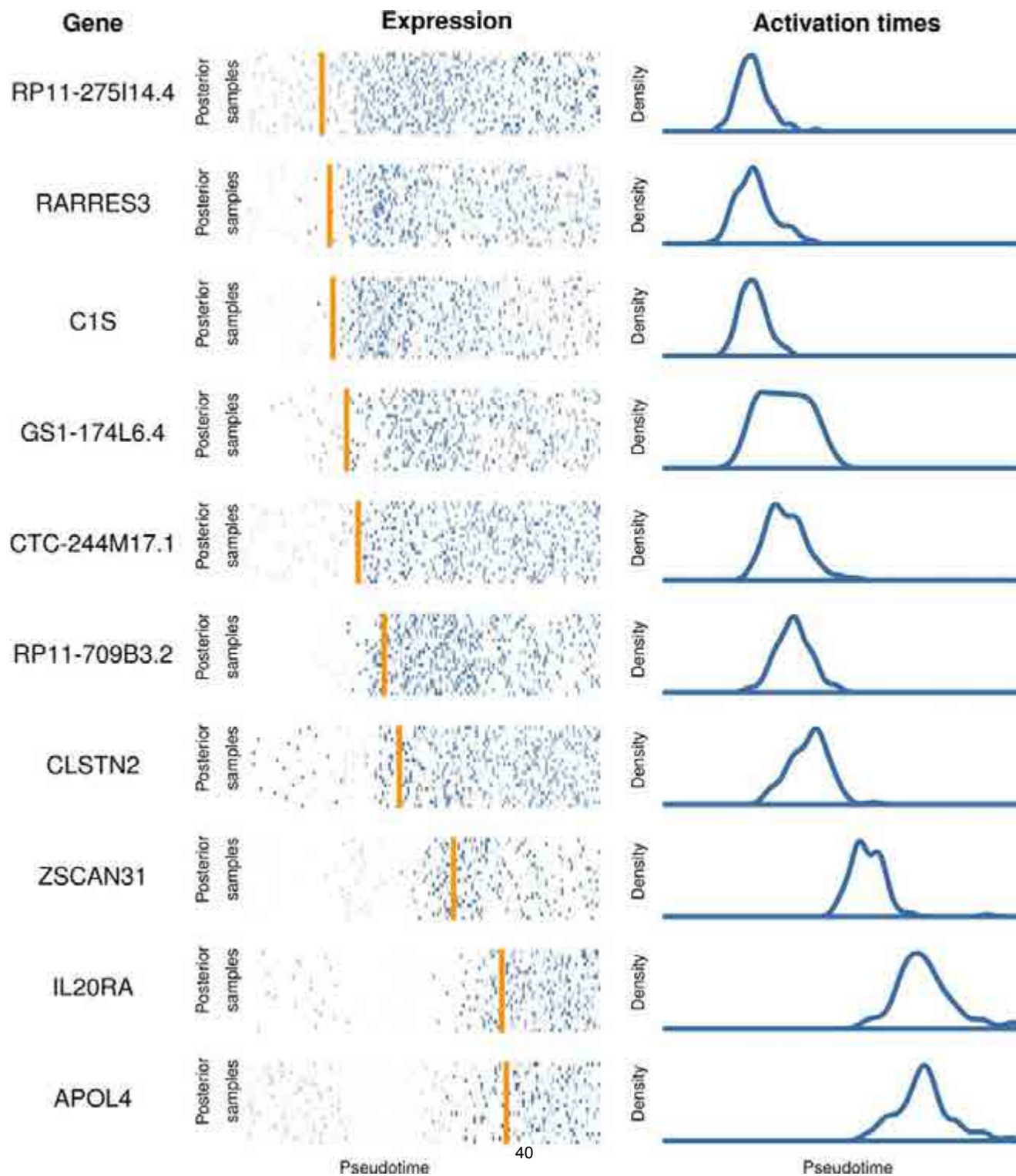


# Inferring original trajectory from single-cell data



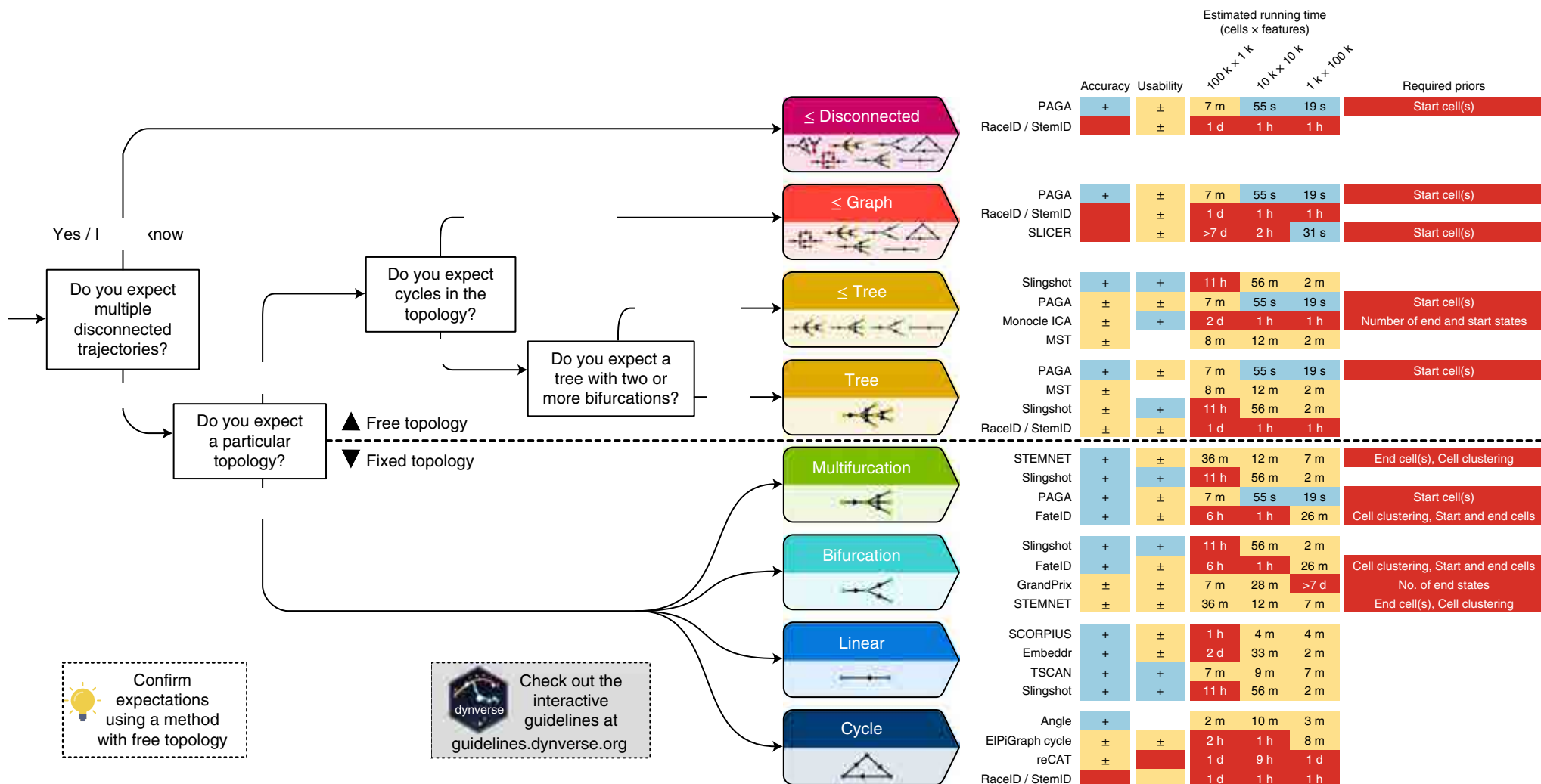
**Fig 1. The single cell pseudotime estimation problem.** (A) Single cells at different stages of a temporal process. (B) The temporal labelling information is lost during single cell capture. (C) Statistical pseudotime estimation algorithms attempt to reconstruct the relative temporal ordering of the cells but cannot fully reproduce physical time. (D) The pseudotime estimates can be used to identify genes that are differentially expressed over (pseudo)time.

[Figure from: Campbell & Yau, PLOS Computational Biology, 2016]



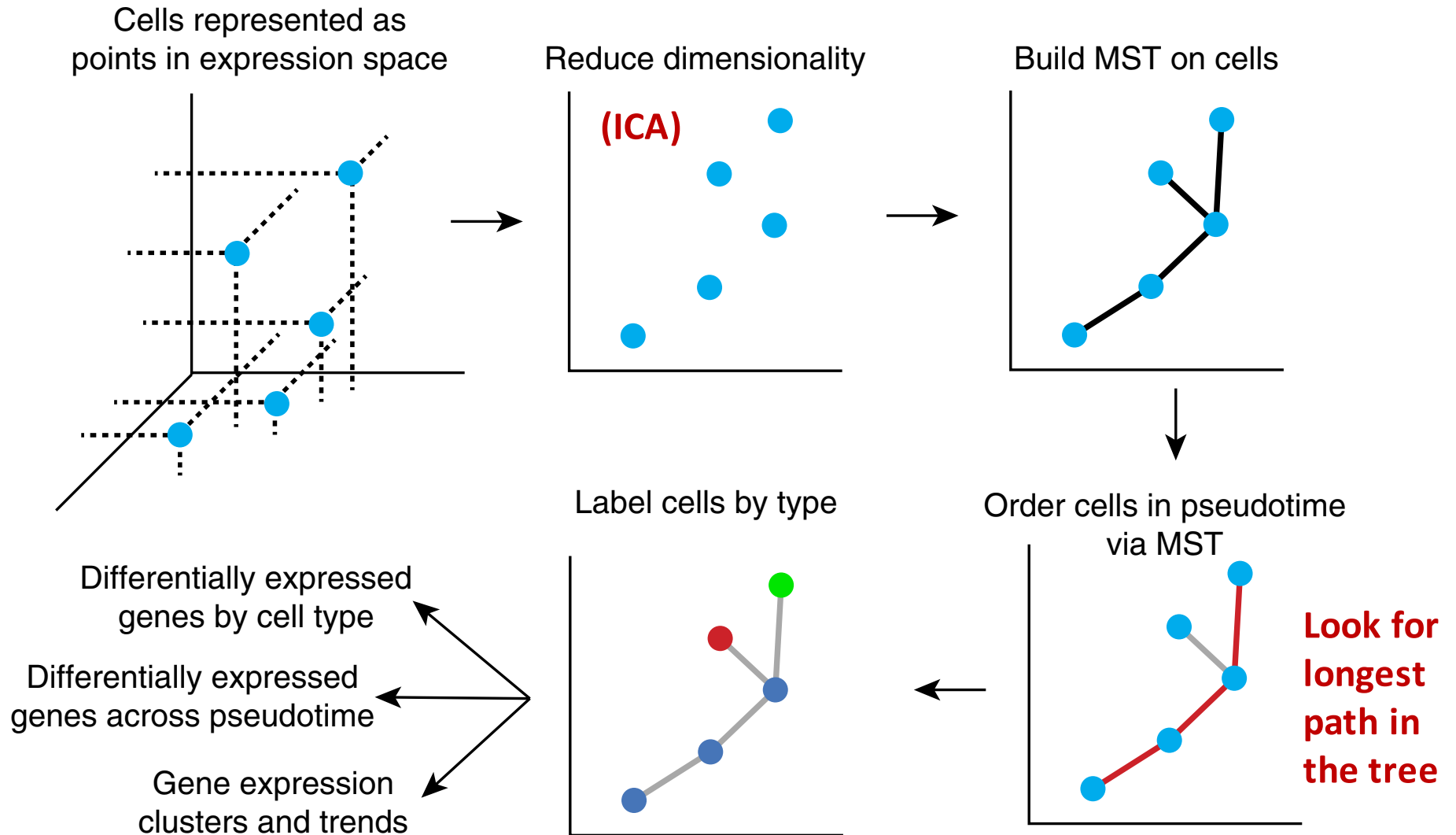
[Campbell & Yau, PLOS Computational Biology, 2016]





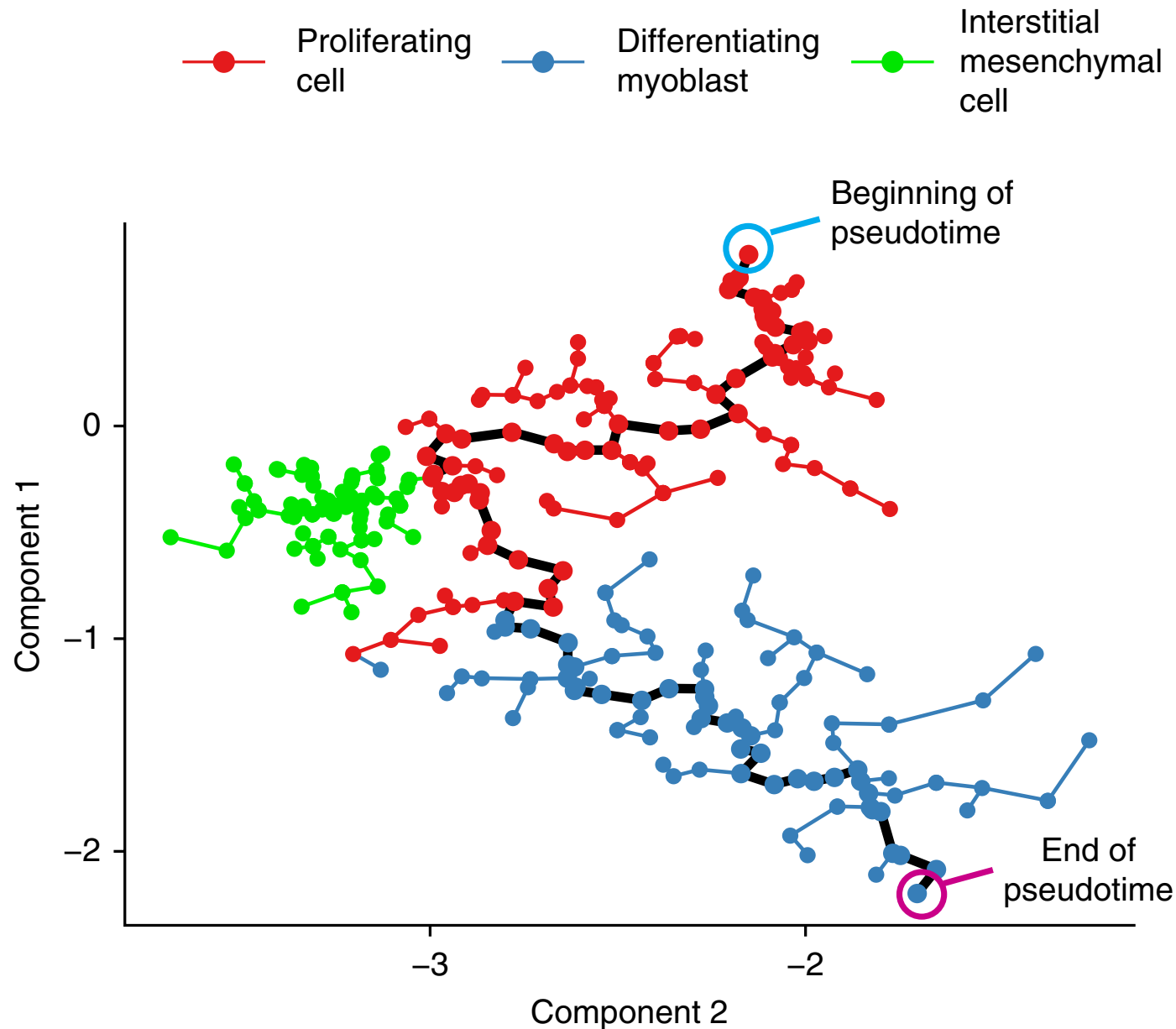
[Saelens, Cannoodt, Todorov, Saeys. A comparison of single-cell trajectory inference methods. *Nature Biotechnology*, 2019]

# MST-based approach (Monocle)



[Magwene et al., *Bioinformatics*, 2003; Trapnell et al., *Nature Biotechnology*, 2014]

# MST-based approach (Monocle)



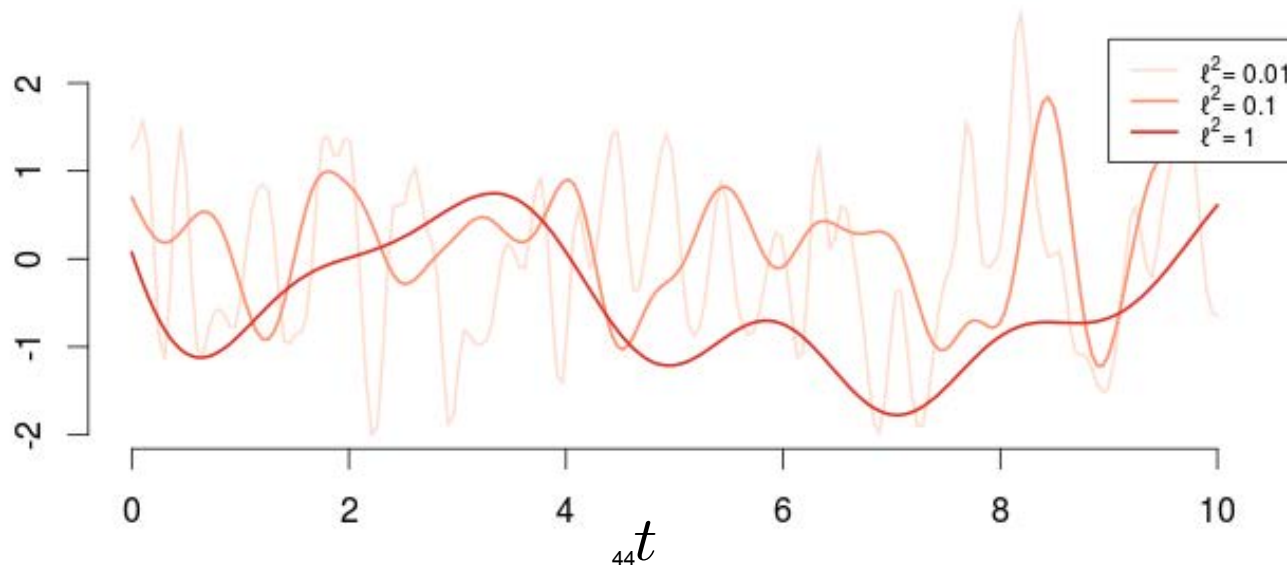
# Statistical model for probabilistic pseudotime

## Definition

$\mu$  is a Gaussian process if for any collection  $\mathbf{T} = \{t_i, i = 1, \dots, N\}$ ,

$$\begin{pmatrix} \mu(t_1) \\ \vdots \\ \mu(t_N) \end{pmatrix} \sim \mathcal{N}(0, K(\mathbf{T}, \mathbf{T}))$$

$$k(t_{i_1}, t_{i_2}) = \tau^2 \exp\left(-\frac{\|t_{i_1} - t_{i_2}\|^2}{2\ell^2}\right) \text{ (squared exponential)}$$



# Statistical model for probabilistic pseudotime

$$t_i \sim \text{TruncNormal}_{[0,1)}(\mu_t, \sigma_t^2), \quad i = 1, \dots, N,$$

$$\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_P^2)$$

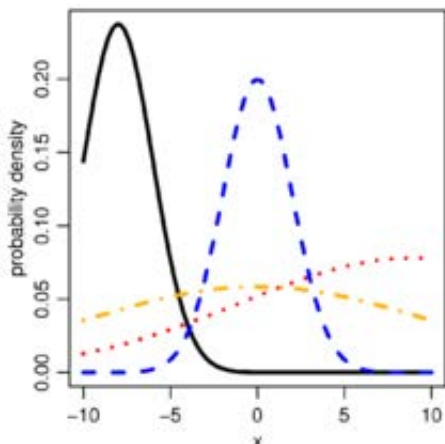
**P: dimension (e.g. 2)**

$$K^{(j)}(t, t') = \exp(-\lambda_j(t - t')^2), \quad j = 1, \dots, P,$$

$$\mu_j \sim \text{GP}(0, K^{(j)}), \quad j = 1, \dots, P, \quad \text{GP: Gaussian Process (1-D)}$$

$$\mathbf{x}_i \sim \text{MultiNorm}(\boldsymbol{\mu}(t_i), \Sigma), \quad i = 1, \dots, N.$$

**N: number of data points**



**Truncated normal  
distribution**

[Campbell & Yau, PLOS Computational Biology, 2016]

# Outline of today's lecture

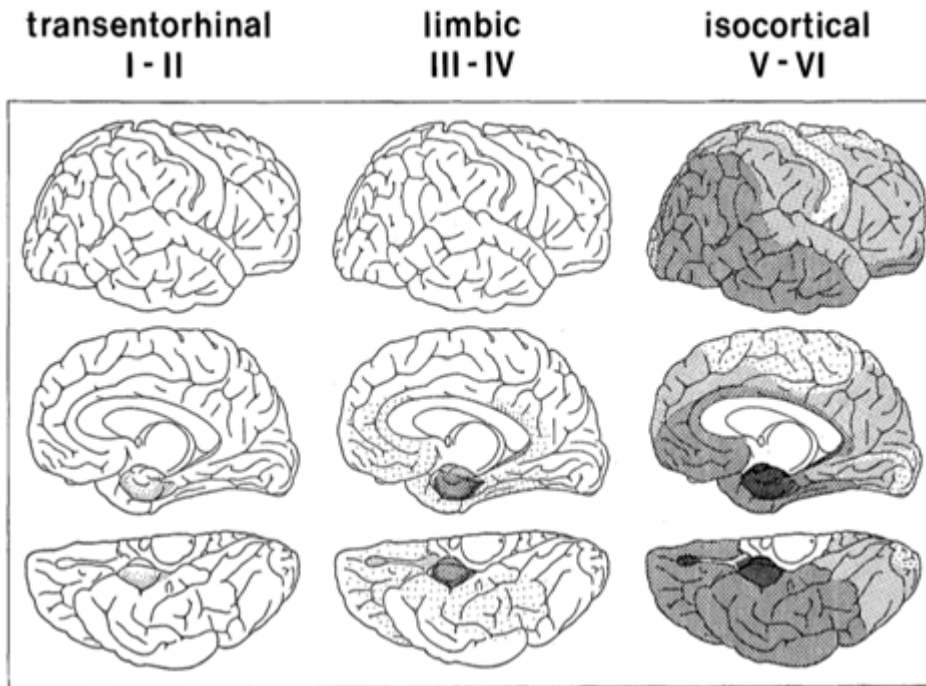
1. Staging from cross-sectional data
  - Wang, Sontag, Wang, *KDD* 2014
  - Pseudo-time methods from computational biology
- 2. Simultaneous staging & subtyping**
  - **Young et al., *Nature Communications* 2018**

Acknowledgement: Subsequent slides adapted from Daniel Alexander

# Temporal heterogeneity

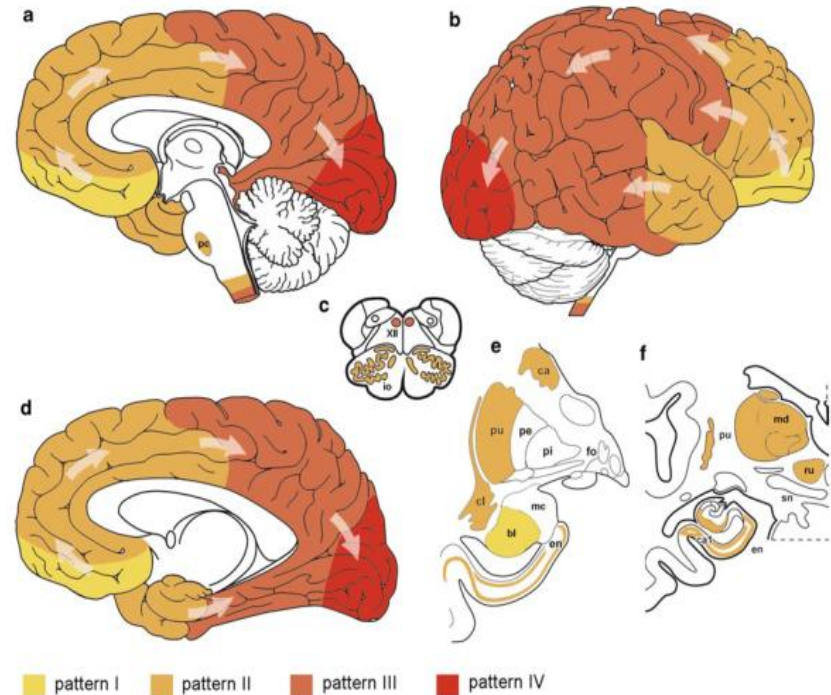
Patients show various disease stages through which patterns of pathology evolve

Alzheimer's disease



[Braak and Braak 1991](#)

Frontotemporal dementia



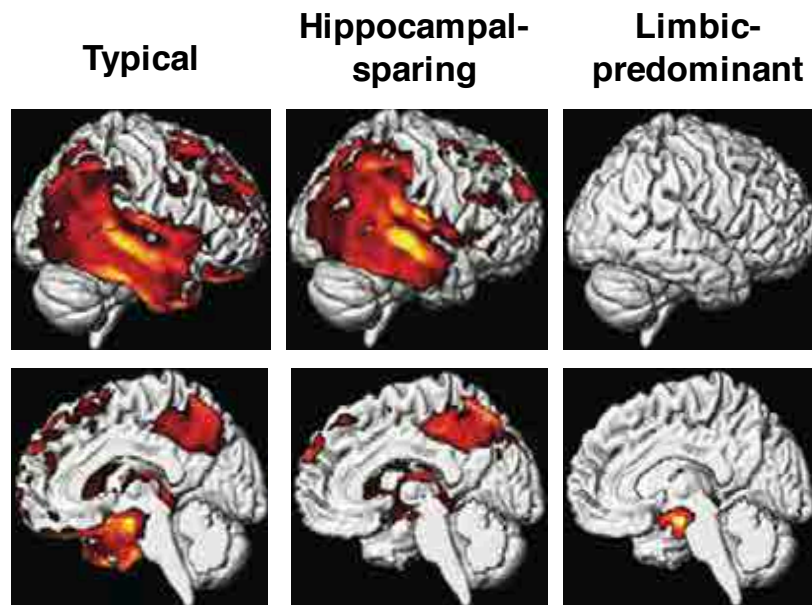
[Brettschneider et al. 2014](#)



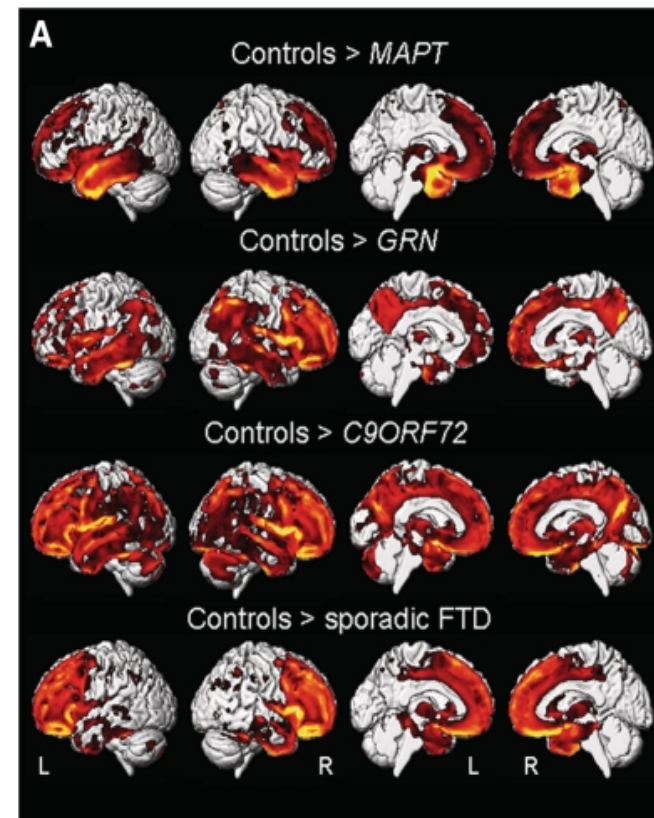
# Phenotypic heterogeneity

Individuals have different disease subtypes with distinct patterns of pathology

Alzheimer's disease



Frontotemporal dementia

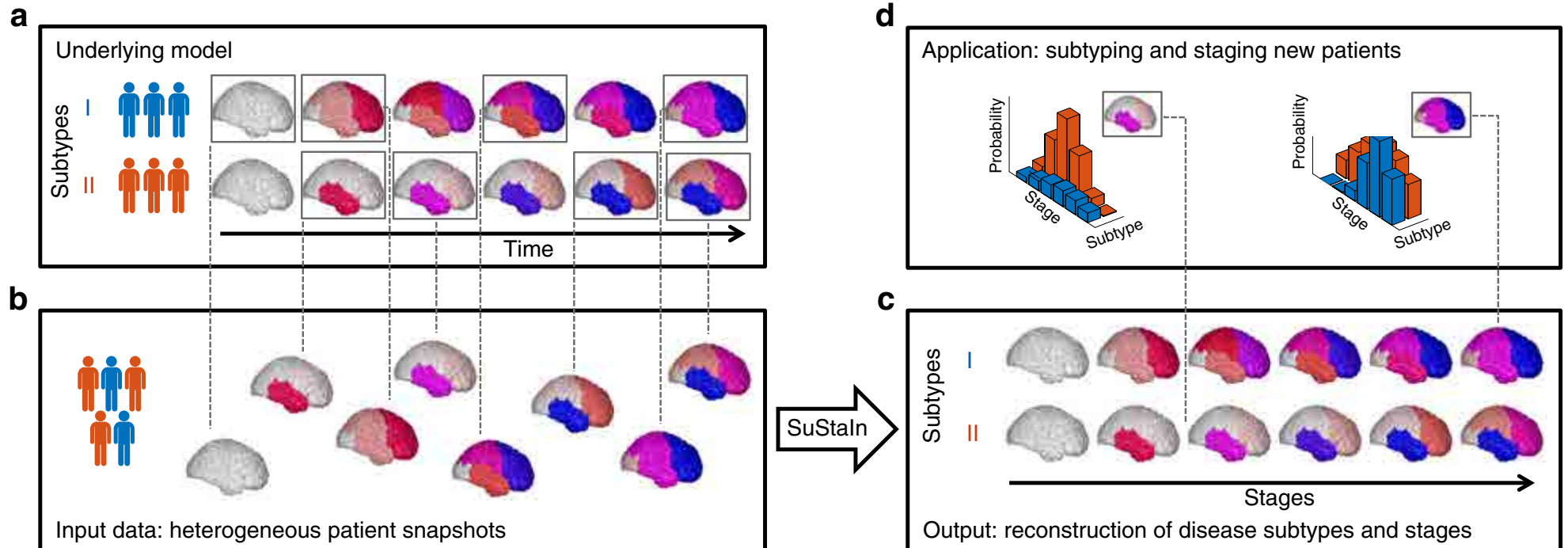


Murray et al. 2011, Whitwell et al. 2012

Whitwell et al. 2012



# Subtype and Stage Inference (SuStaln)



[Young et al., *Nature Communications* 2018]

Courtesy of [Springer Nature](https://www.springer.com). Used under CC BY.

# Subtype and Stage Inference (SuStain)

- Generative model for a data point:
  - Sample subtype  $c \sim \text{Categorical}(f_1, \dots, f_C)$
  - Sample stage  $t \sim \text{Categorical}(\text{uniform})$
  - For each biomarker  $i$ , sample  $x_i \sim \mathcal{N}(g_{c,i}(t), \sigma_i)$
- Means are enforced to be monotonically increasing and piece-wise linear:

$$g(t) = \begin{cases} \left( \frac{z_1}{t_{E_{z_1}}} t, 0 < t < t_{E_{z_1}} \right. \\ \left. z_1 + \frac{z_2 - z_1}{t_{E_{z_2}} - t_{E_{z_1}}} (t - t_{E_{z_1}}), t_{E_{z_1}} < t < t_{E_{z_2}} \right. \\ \vdots \\ \left. z_{R-1} + \frac{z_R - z_{R-1}}{t_{E_{z_R}} - t_{E_{z_{R-1}}}} (t - t_{E_{z_{R-1}}}), t_{E_{z_{R-1}}} < t < t_{E_{z_R}} \right. \\ \left. z_R + \frac{z_{\max} - z_R}{1 - t_{E_{z_R}}} (t - t_{E_{z_R}}), t_{E_{z_R}} < t < 1 \right) \end{cases}$$

**Shown here for one choice of  $c, i$  – no parameter sharing across biomarkers or subtypes**

[Young et al., *Brain* 2014; Young et al., *Nature Communications* 2018]

MIT OpenCourseWare

<https://ocw.mit.edu>

**6.S897 / HST.956 Machine Learning for Healthcare**

Spring 2019

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>