



Bayesian Hierarchical Models to Quantify and Communicate an Individual's Disease *State, Trajectory, and Likely Treatment Benefits* from Multivariate *Longitudinal Data*

Scott L. Zeger

John C. Malone Professor of Biostatistics and Medicine
Co-Director, Hopkins *inHealth*
Johns Hopkins University Bloomberg School of Public Health

And

Biostatistics: Jisoo Kim, Zhenke Wu, Yizhen Xu, Shannon Wonbgvibulsin, Gege Gui, Zitong Wang, Emily Scott, Joe Sartini, MaryGrace Bowring, Tianxu Wang, Ning Meng, Zixing Liu,...

Clinical Sciences: Antony Rosen, Ami Shah, Brian Garibaldi, John Robinson,...

Johns Hopkins 1st Year Biostatistics Student Seminar
(slides from sz@jhu.edu)

September 19, 2024

What is Johns Hopkins *in*Health Precision Medicine?

<https://www.hopkinsmedicine.org/inhealth/>

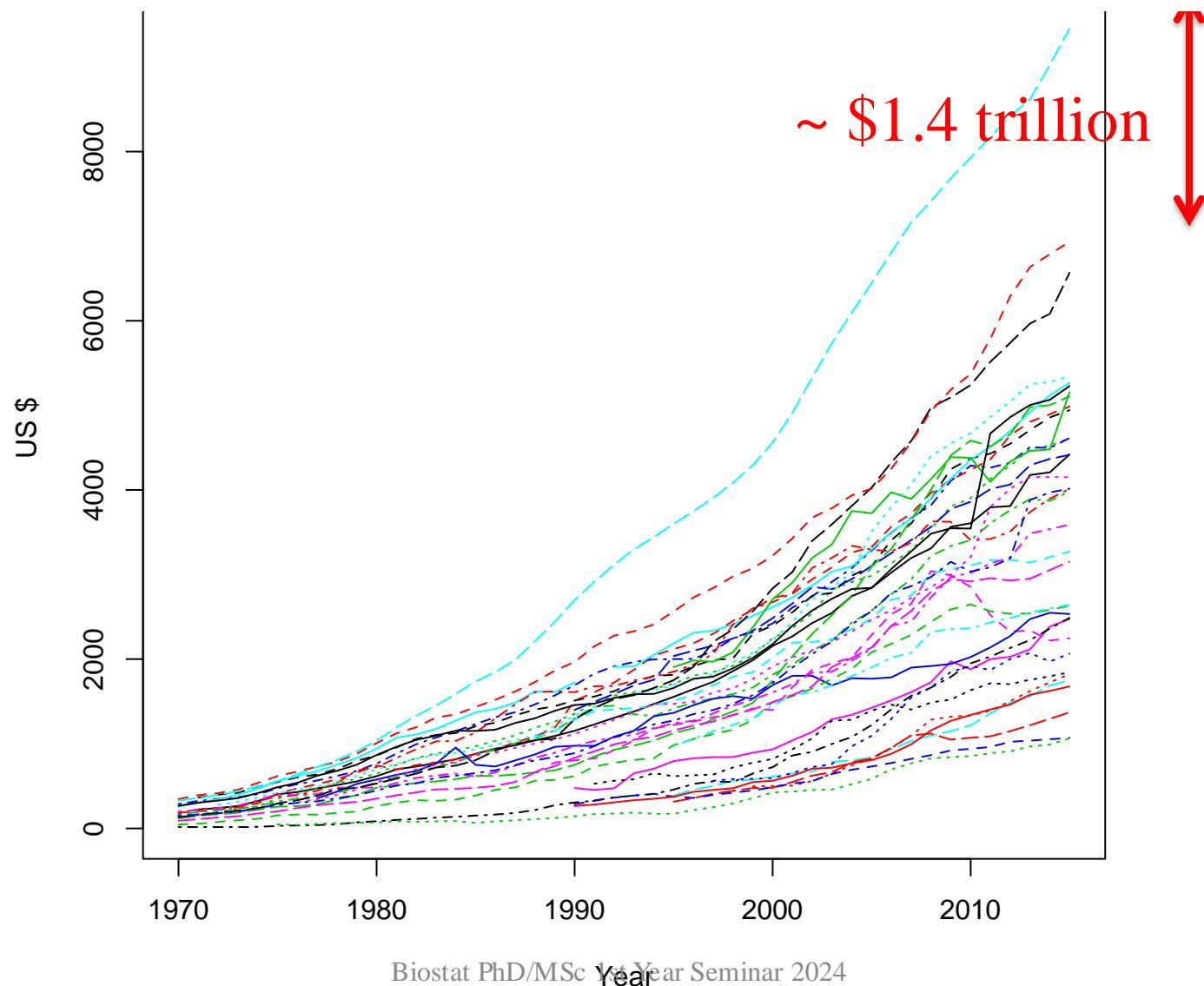
Use clinical data to answer three questions

1. What is this patient's current disease state?
2. What is their disease trajectory?
3. Do they belong to an identified subgroup?
4. Which available intervention is more likely to produce a better future state?

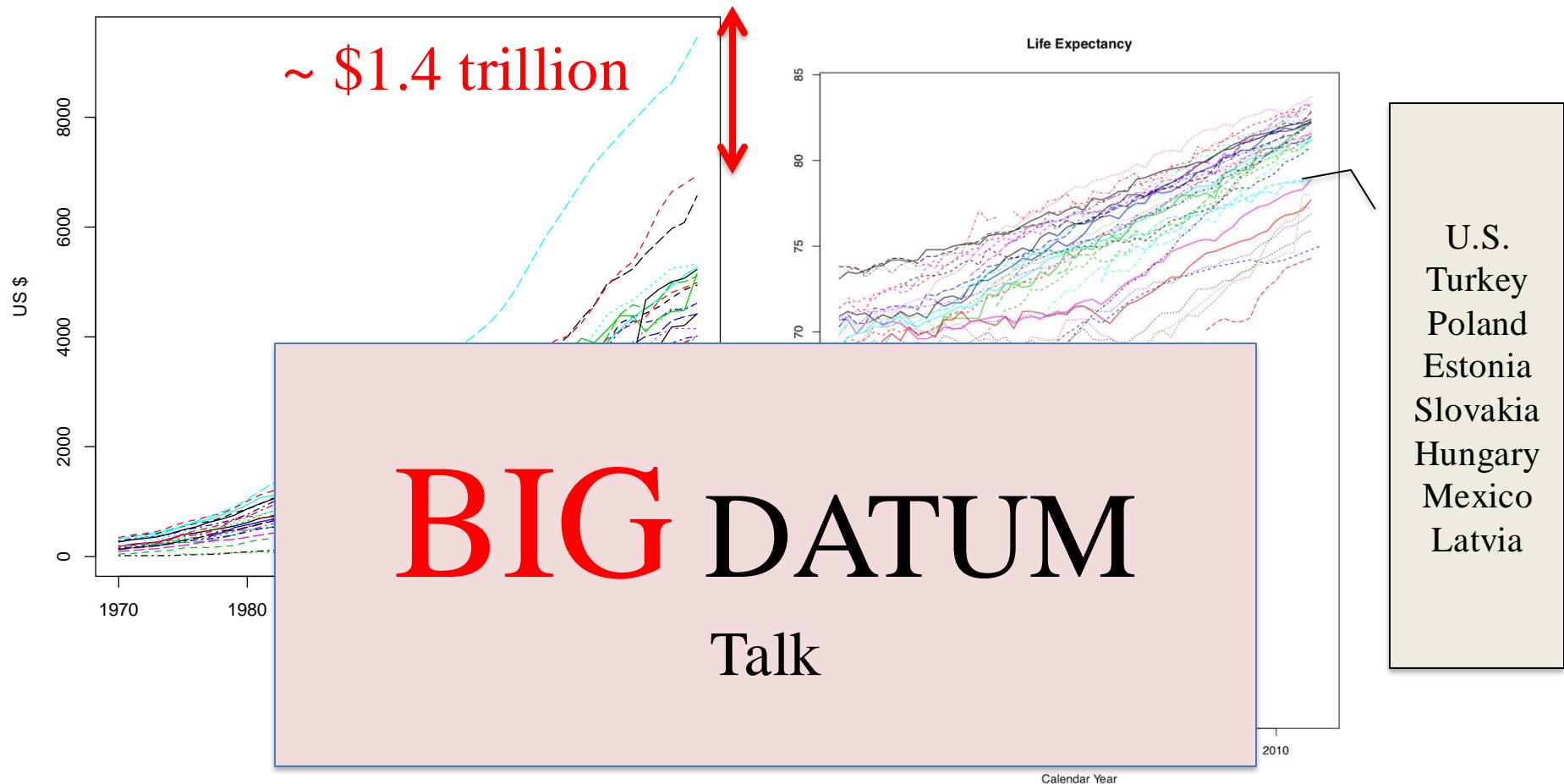
Outline

- **Context:** waste in American healthcare
- **Perspective:** individual's health
- **Statistical role:** Bayesian hierarchical models
- **Specific applications**
 - Autoimmune diseases
 - Covid-19
- Research opportunities

Per Capita Annual Medical Expenditures - OECD Countries



Per Capita Annual Medical Expenditures (left) and Life Expectancy (right) - OECD Countries



What can statisticians do about this situation (in addition to voting)?

Learning in Healthcare



Individual's Perspective: PLAY A CLINICIAN for a moment

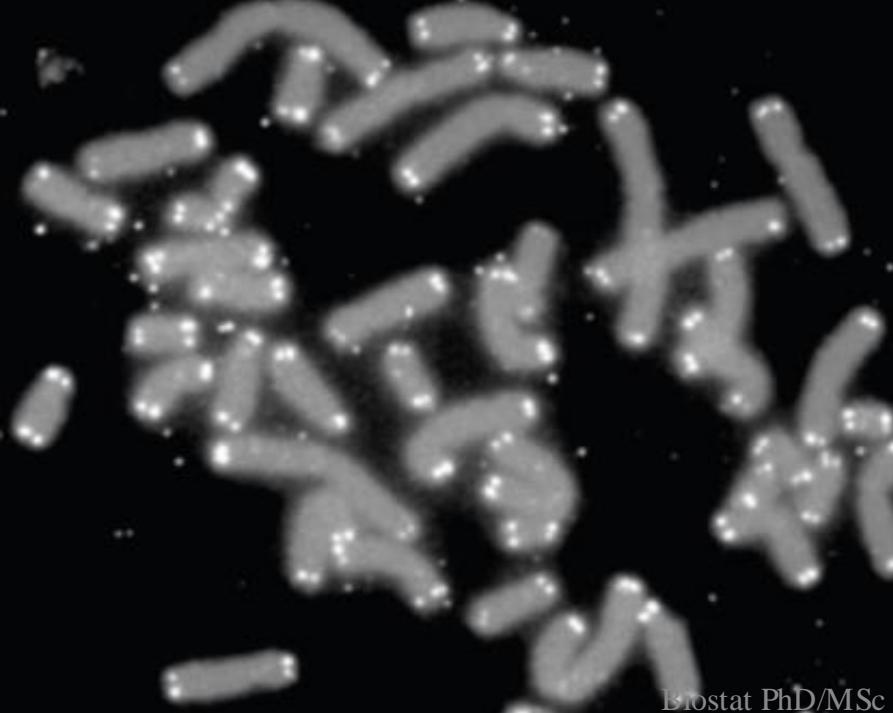
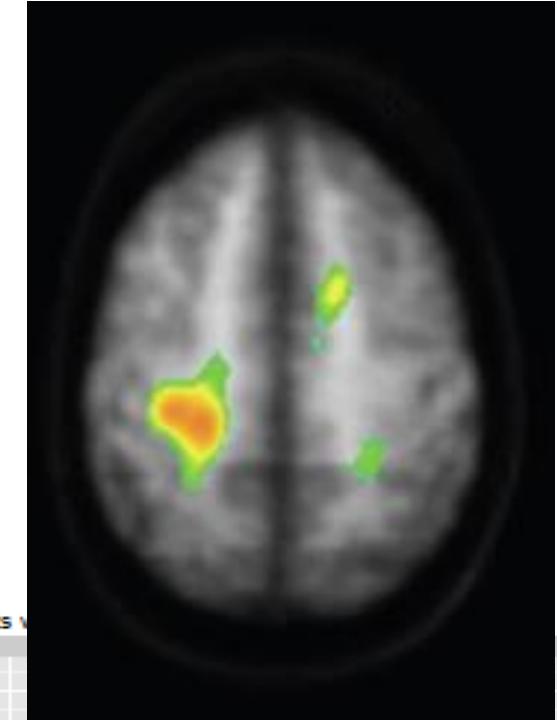
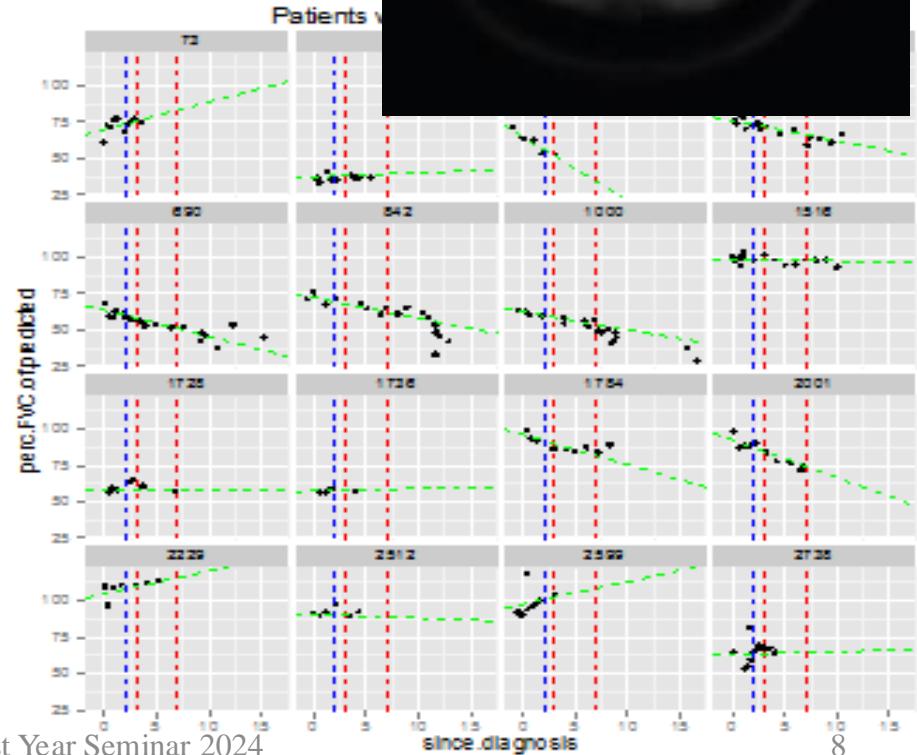
40 year old man, no family history, tests positive for a life-threatening disease in a routine screen, so comes to your office asking: what is my disease state; what action do you recommend?

Data from prior population of similar people

	True disease status		
Exam result	Yes	No	Total
Positive	15	985	1,000
Negative	5	8,995	9000
Total	20	9,980	10,000

Biostat PhD/MSc 1st Year Seminar 2024

119679	tttgaatcaa	aatttacata	gttttcttt	tagactaagc	tcc
119739	cccattcctc	attaccattg	aaatgtctca	tgagcatgtc	aca
119799	atccaggatg	acagttagt	tcttttaat	ccaattggaa	gcct
119859	gaacctaaag	aaaggttaag	atacatttat	tccctgggt	aagg
119919	gttttcctaa	gggtcatatt	tcaattttaga	ttttttttta	tag
119979	tcccctttgc	aatatgaat	atgttagtctt	ttaaaaaat	tctt
120039	aaaaaaaaatt	aatttggct	attcagtttg	tttagcactta	ccat
120099	ctctactttt	gtattingta	acattttccc	tactacaggg	cagg
120159	tagatattag	caccaaataa	ataggcaaaa	aaaatctatt	atg
120219	tgcttggcag	tgcacatcg	actagatgg	gaagaatag	aaa
120279	gtttccctgg	tcttttggaa	acaactagag	agttttgtt	ttg
120339	tcctgtttaa	tgctttcatt	ctatgattgt	taagaatatg	tca
120399	gtttctttat	gtcttcctt	ctgtttgtt	attagaat	cctt
120459	tagtacgta	gatatgtat	atattccat	aattacactg	ctgt
120519	taatttttagg	gcagctttat	gacagtttgt	ttatgtttt	ggg
120579	agcattgaaa	tctgggttatt	aagcacactg	ttttctatgt	ggta
120639	gccctgagaa	aatggaaaat	aaaaatattt	ttcccttttta	ccat
120699	cactctatca	taaactgcat	aaatcttata	actctaaaaac	attt
120759	gaacttgc	actcaattgc	ttctatatac	acccaaatattt	ttt
120819	gtccttgaaa	atattttgtt	ctactcaata	gaagcagttt	agg
120879	aaaccgttag	gaaataattt	tatattatga	tgactagacc	agtc
120939	gttattgttc	cattagtaaa	tattataatt	atttctgaga	ttt
120999	gttggcaatg	ccagcattat	taacactcct	ctagttgaa	caaa
121059	aaaacataat	aatagccaaa	taaagagtg	cttagaatgt	acaa
121119	gagtaattcg	attatttcta	ggaaatacac	ttttgtgcta	gaad
121179	gctaatttct	gggtttttt	tcattttgaa	ttaacttgaa	tttc
121239	tttttacaga	tacagtgcatt	agaagctcg	tgatacaatg	agaa
121299	aaaatgccat	tagattttc	atcggtatac	tatctgata	gtg
121359	tataccctcat	tatagtactt	cctaatgtaa	tttcttaatt	taaa
121419	tttttttata	taaacttaag	tactgtttaa	tatthaaggc	----
121479	acttgtgtat	atcttattcc	aagcatattt	gttctctcc	
121539	tcatttccaa	aattgtttta	ctcacaactg	tttggttttt	



Two *in*Health Statistical Goals

- Create the inferential analogue of the 2x2 table for any measurements relevant to a major health decision

Population \Leftrightarrow Individual

- Build capacity to make tables for ever narrower sets of “otherwise similar” individuals

Subset, Subset, Subset

Bayesian Hierarchical Models

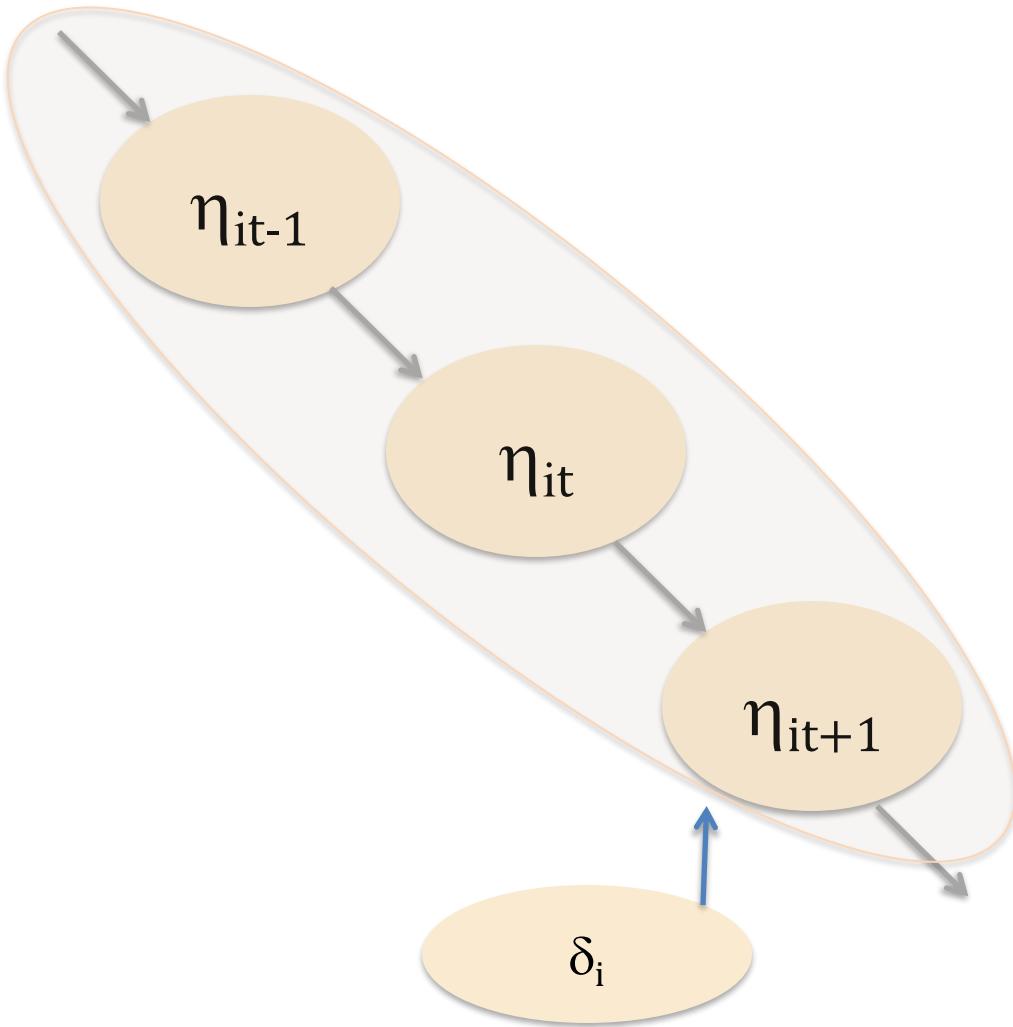
A Really Useful Model for inHealth

<https://www.youtube.com/watch?v=dVmLTApnwag>

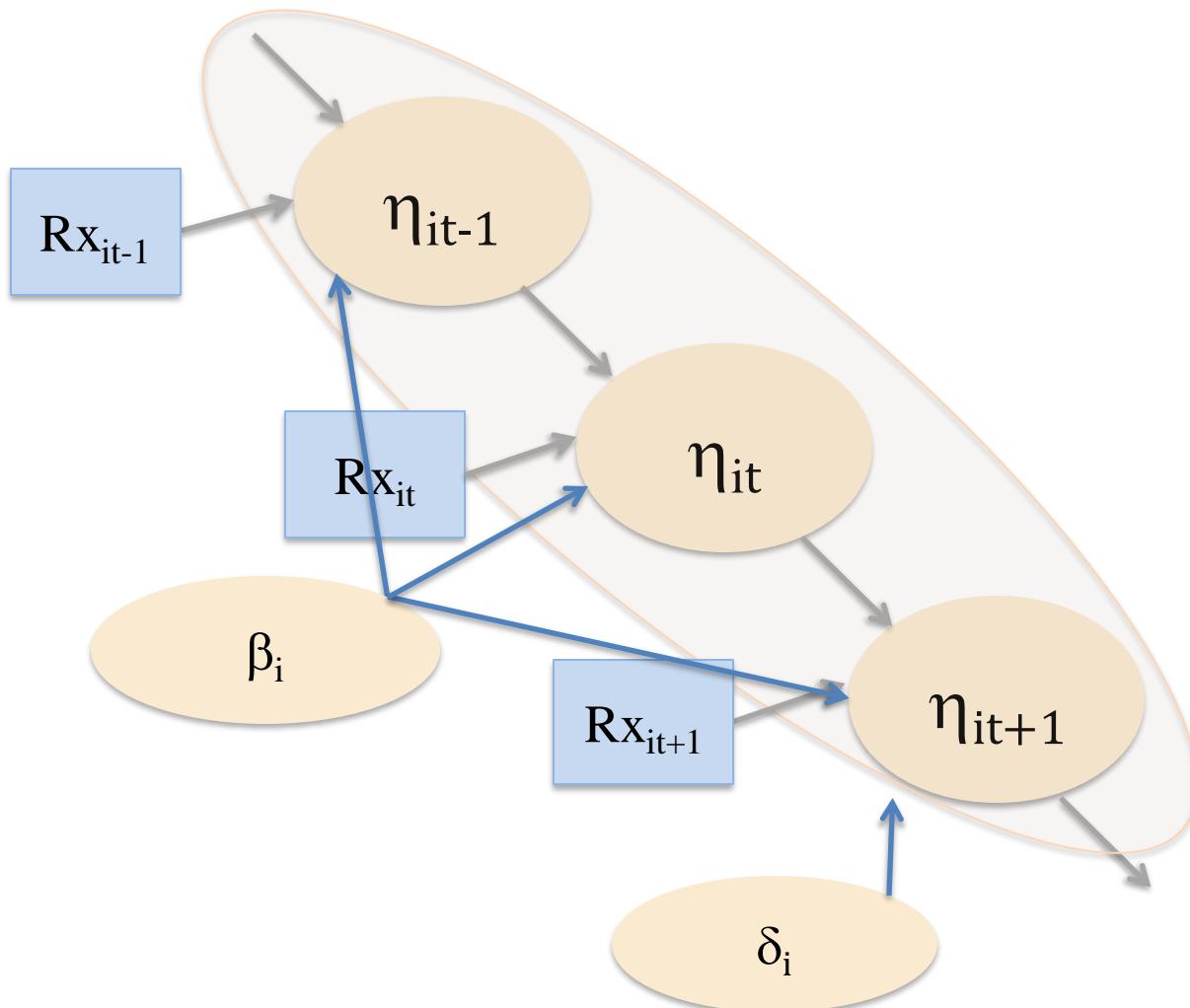
What is this train's name?



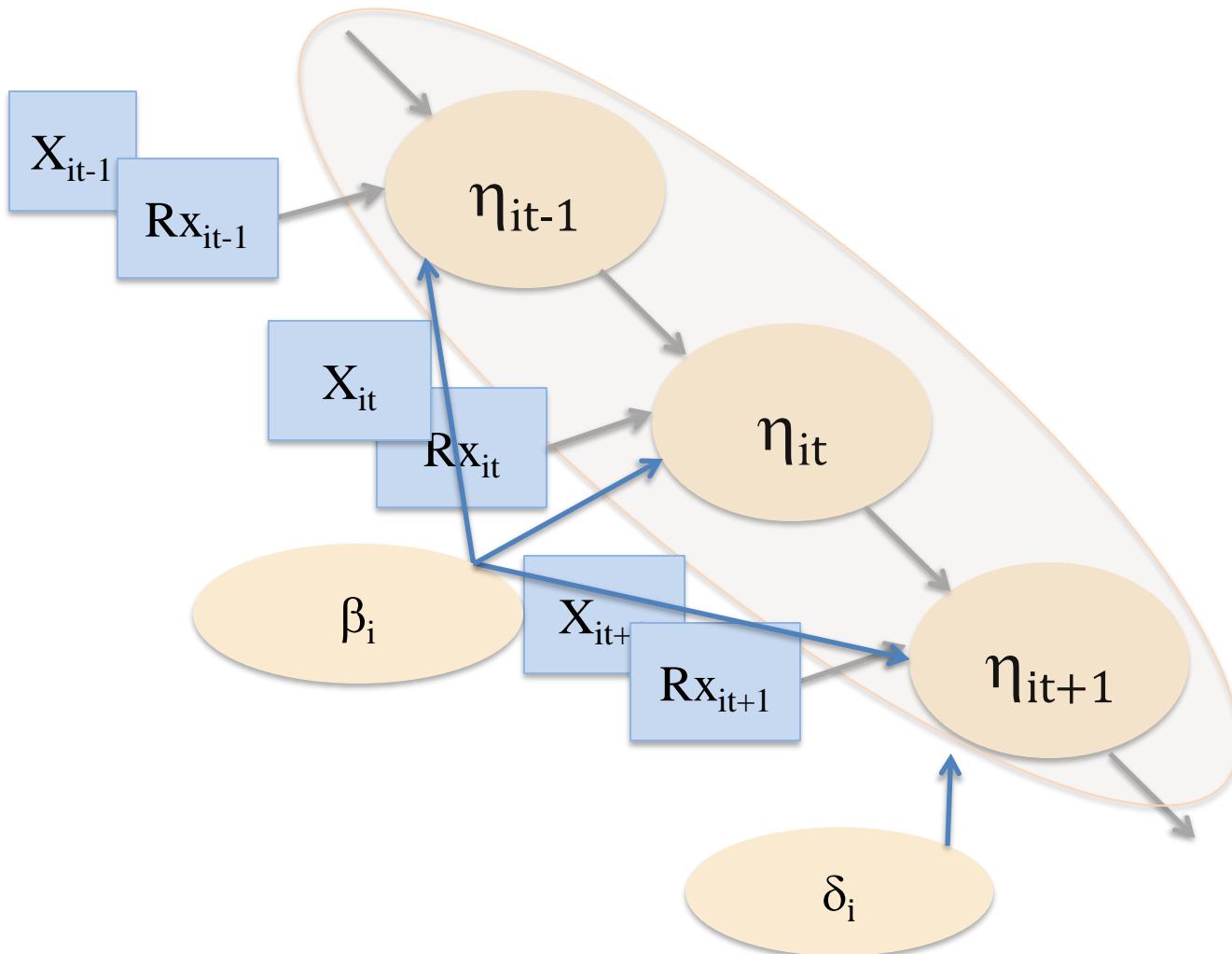
Health State/Trajectory (η_{it}) with Person-specific Indicator (δ_i)



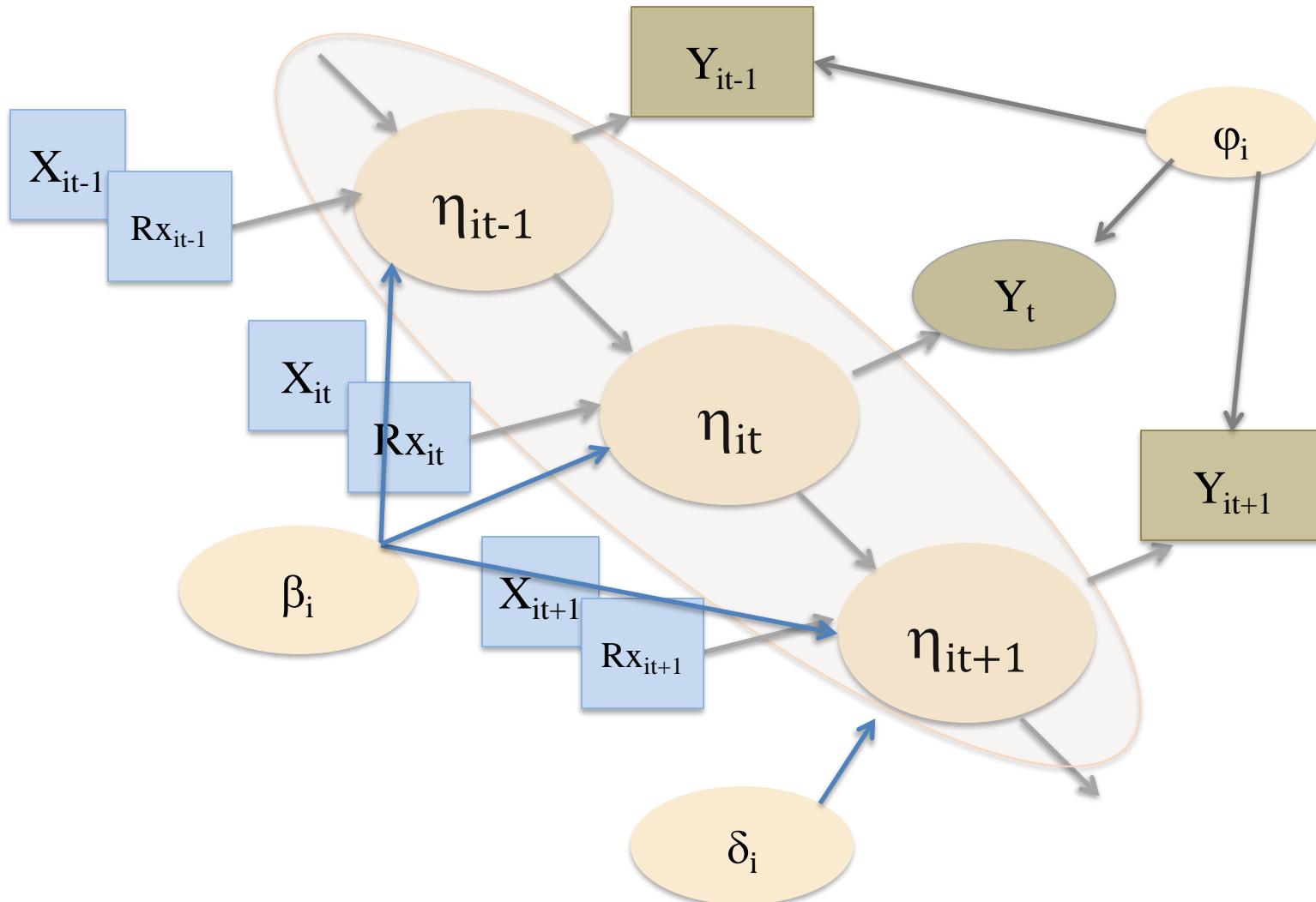
Effects of Exogenous (X) and Endogenous (Rx) Covariates on Health State/Trajectory with Person-specific Regression Coefficients (β_i)



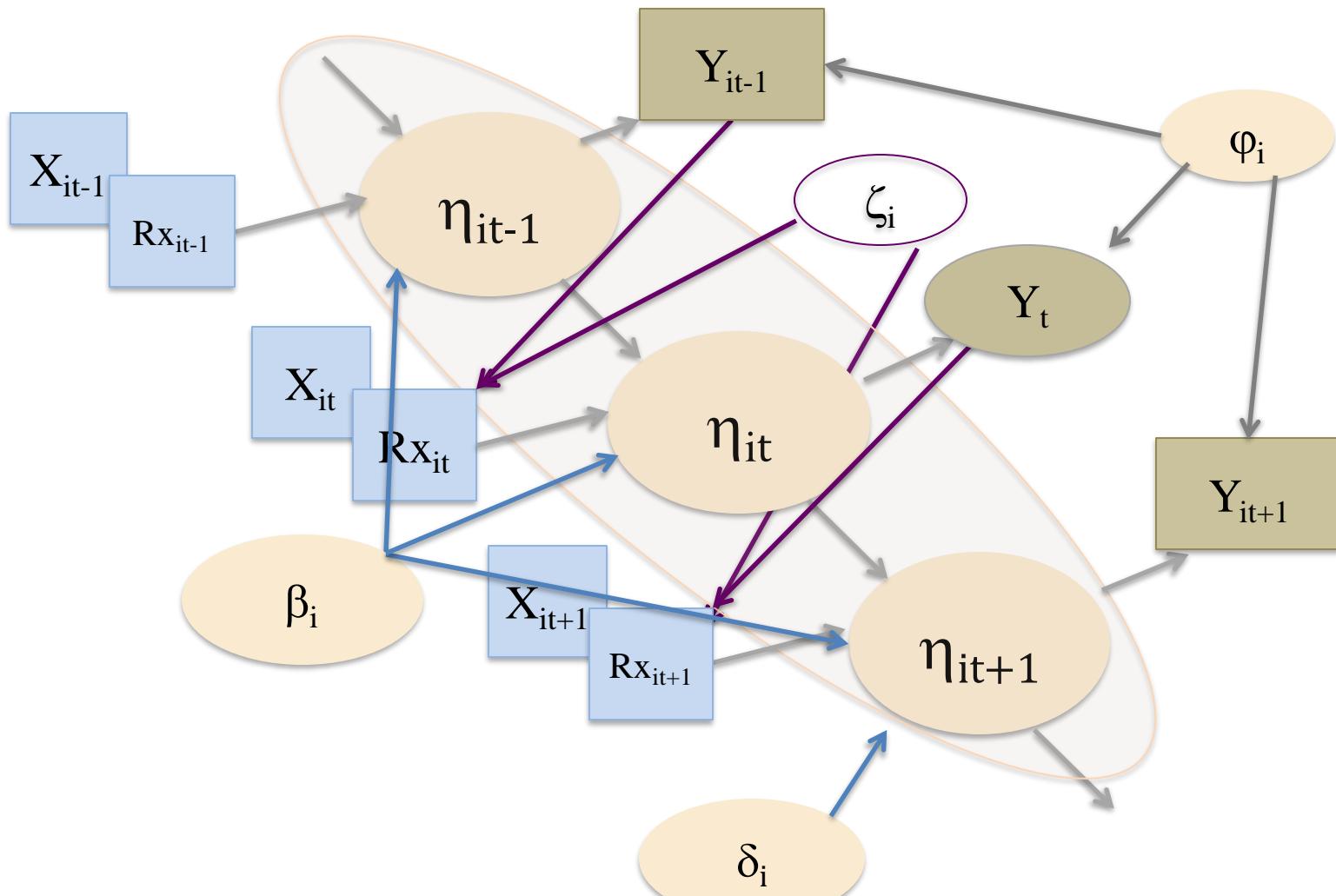
Effects of Exogenous (X) and Endogenous (Rx) Covariates on Health State/Trajectory with Person-specific Regression Coefficients (β_i)

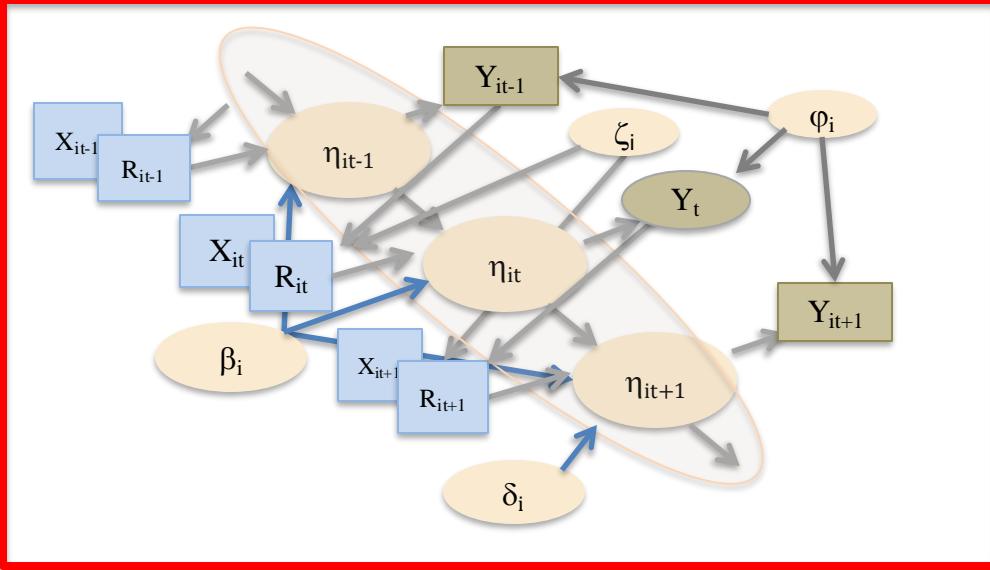


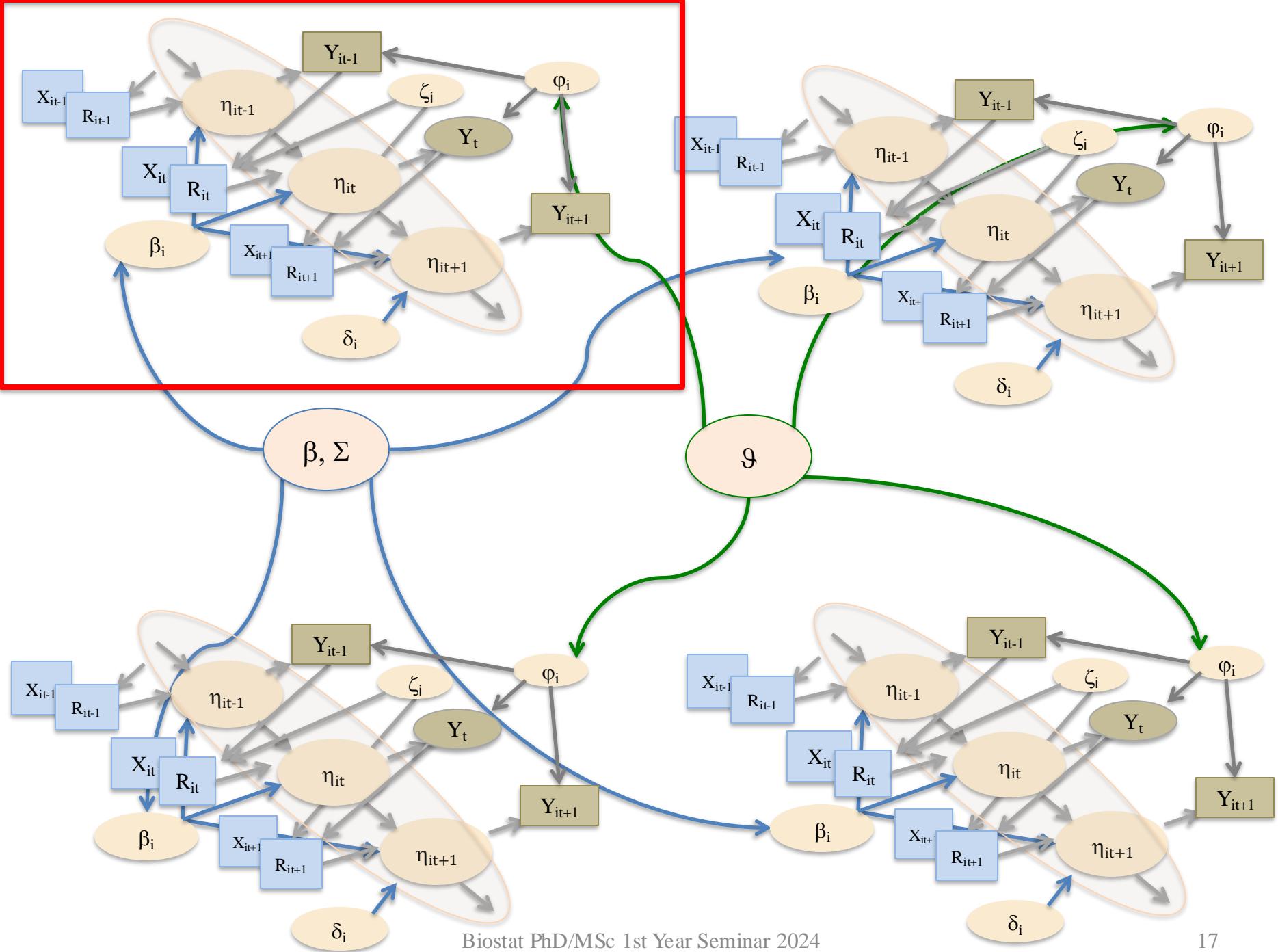
Observations (Y) that Inform about Health State through Coefficients (φ_i)

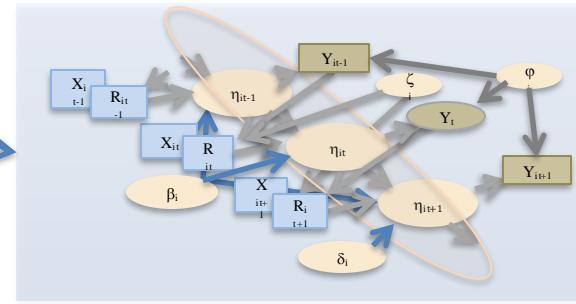
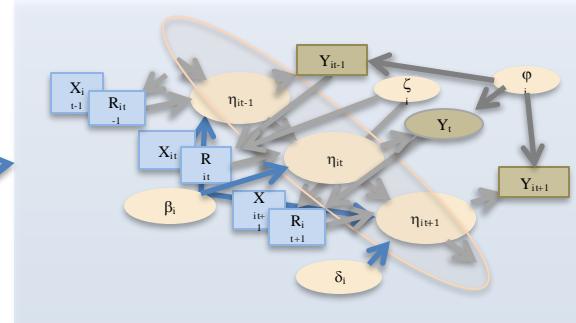
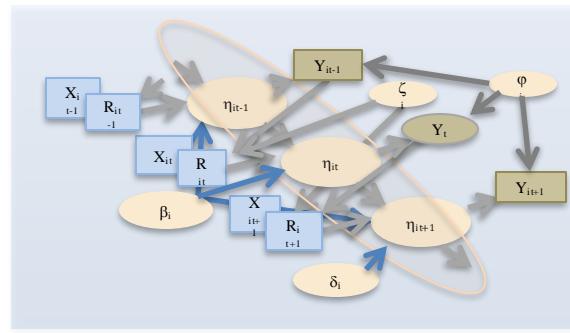
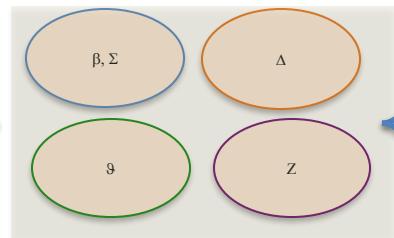
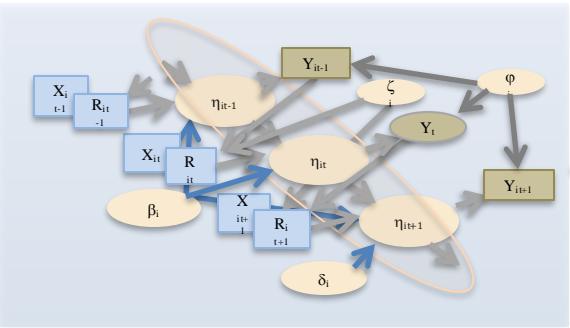


Treatment Decisions Depend on Past Measured Outcomes through Parameters (ζ_i)

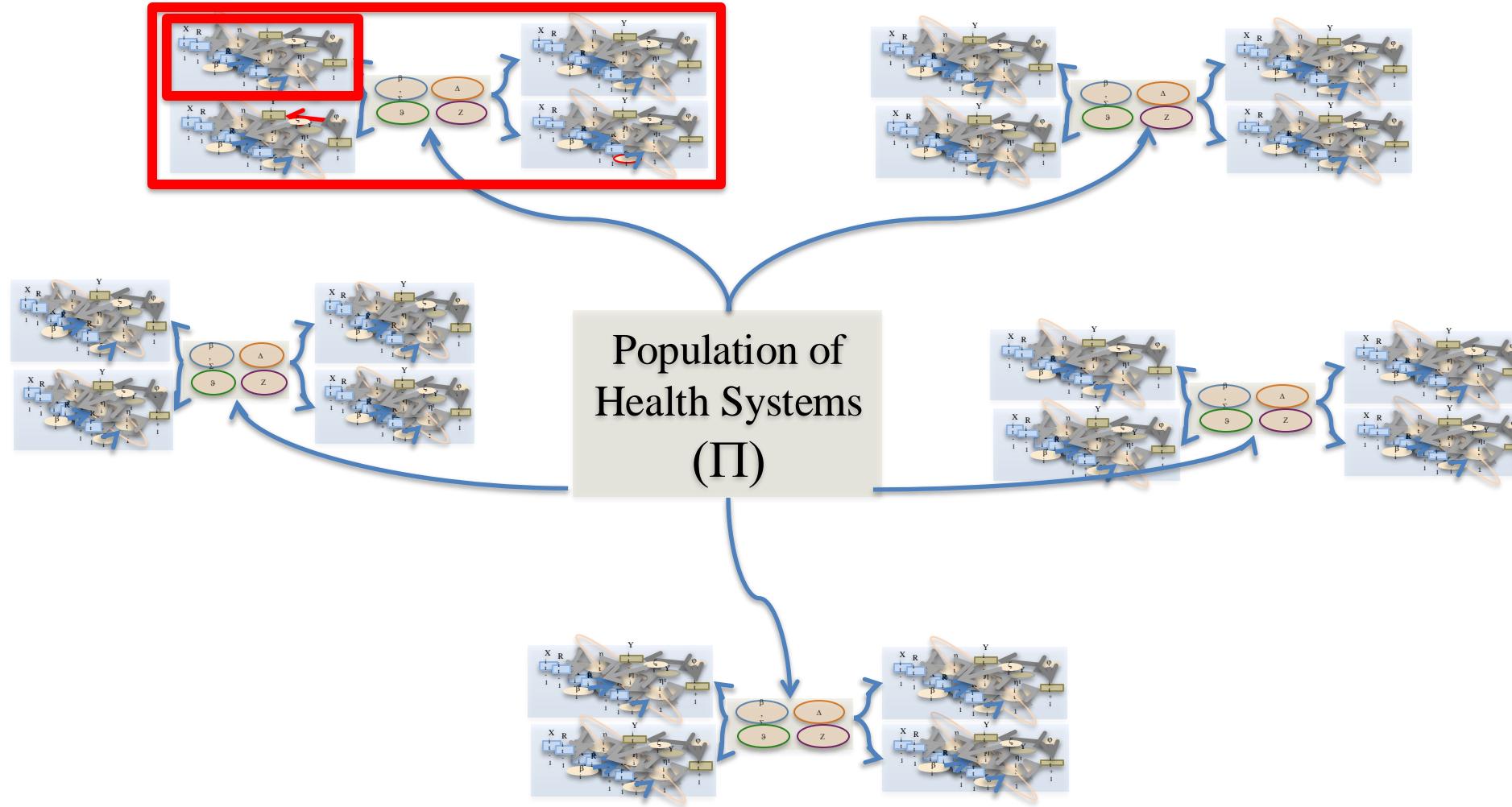








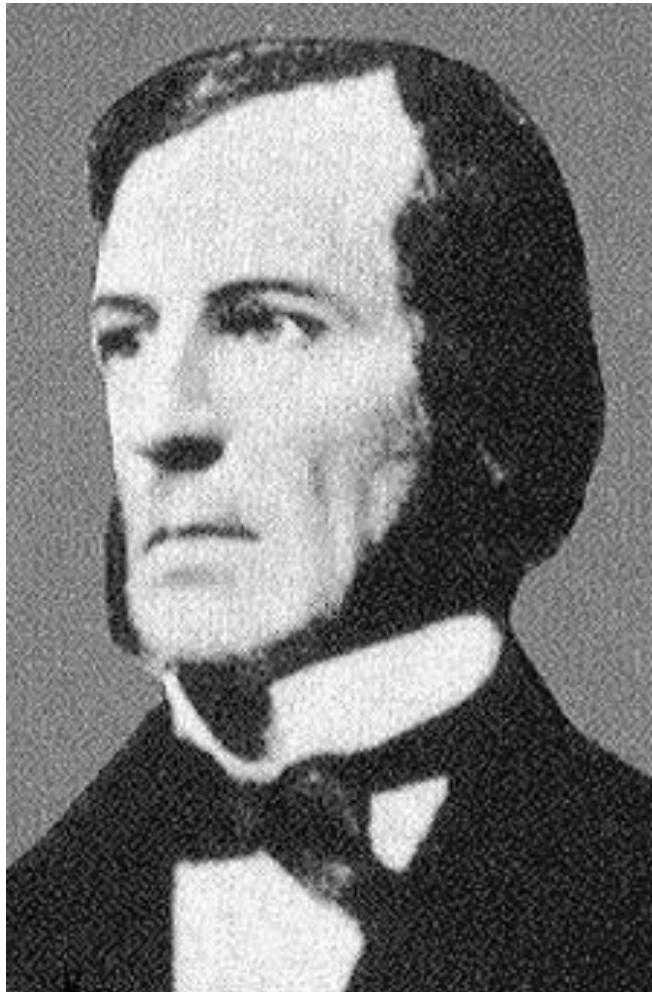
Population of Health Systems (Π)



Comments

- Formalism to make inferences about a single individual's health state, trajectory, and likely treatment effect.
- Address questions by combining individual's data with data on a reference population of others within and beyond her system.
- Make predictions and observe outcomes for model improvement
- Explicit reliance on informative priors where appropriate!
- A positive step from current practice?

Boole -to- Bayes



George Boole



Thomas Bayes



Bouillabaisse



Bayes to von Neumann ?

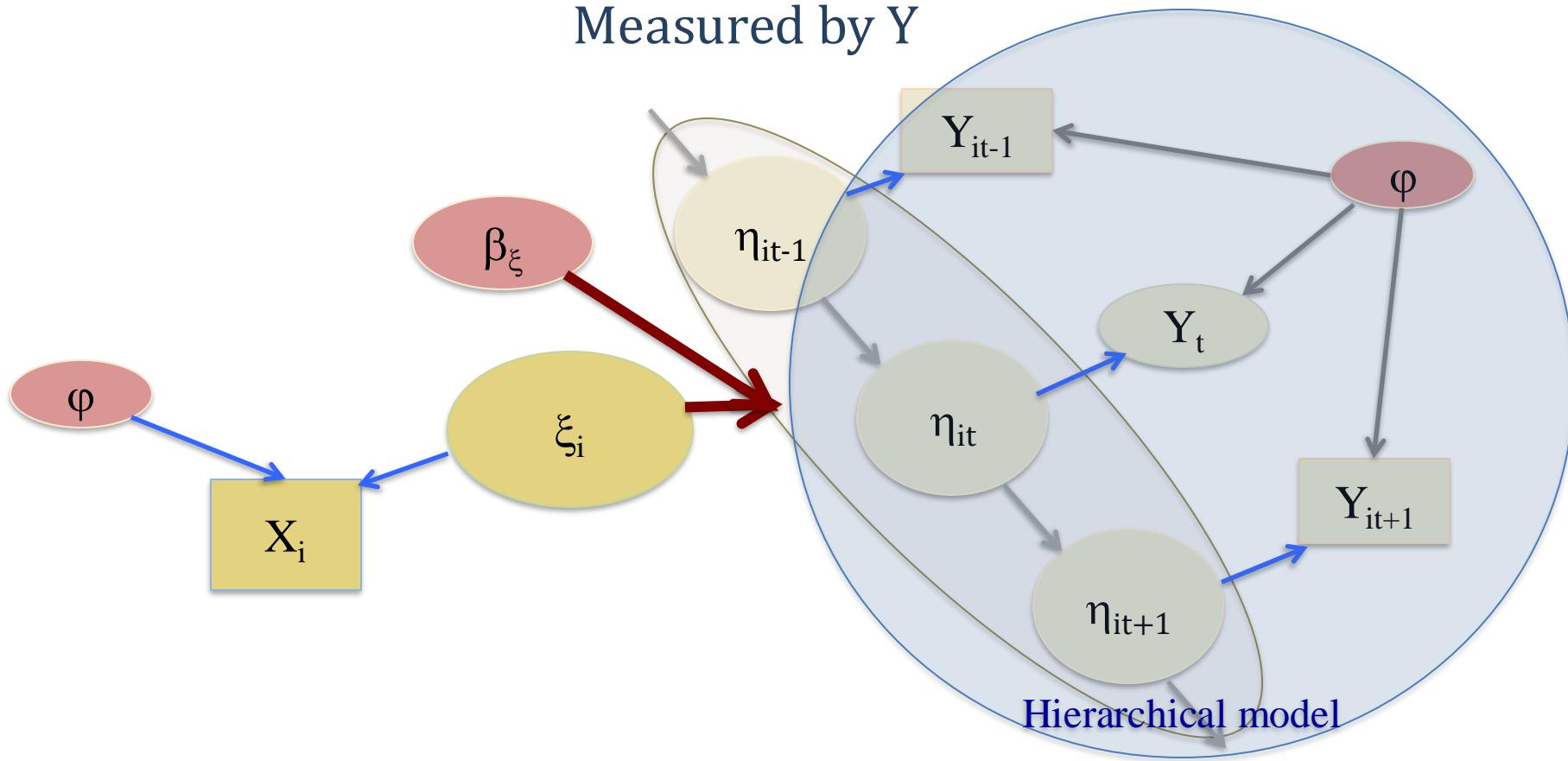


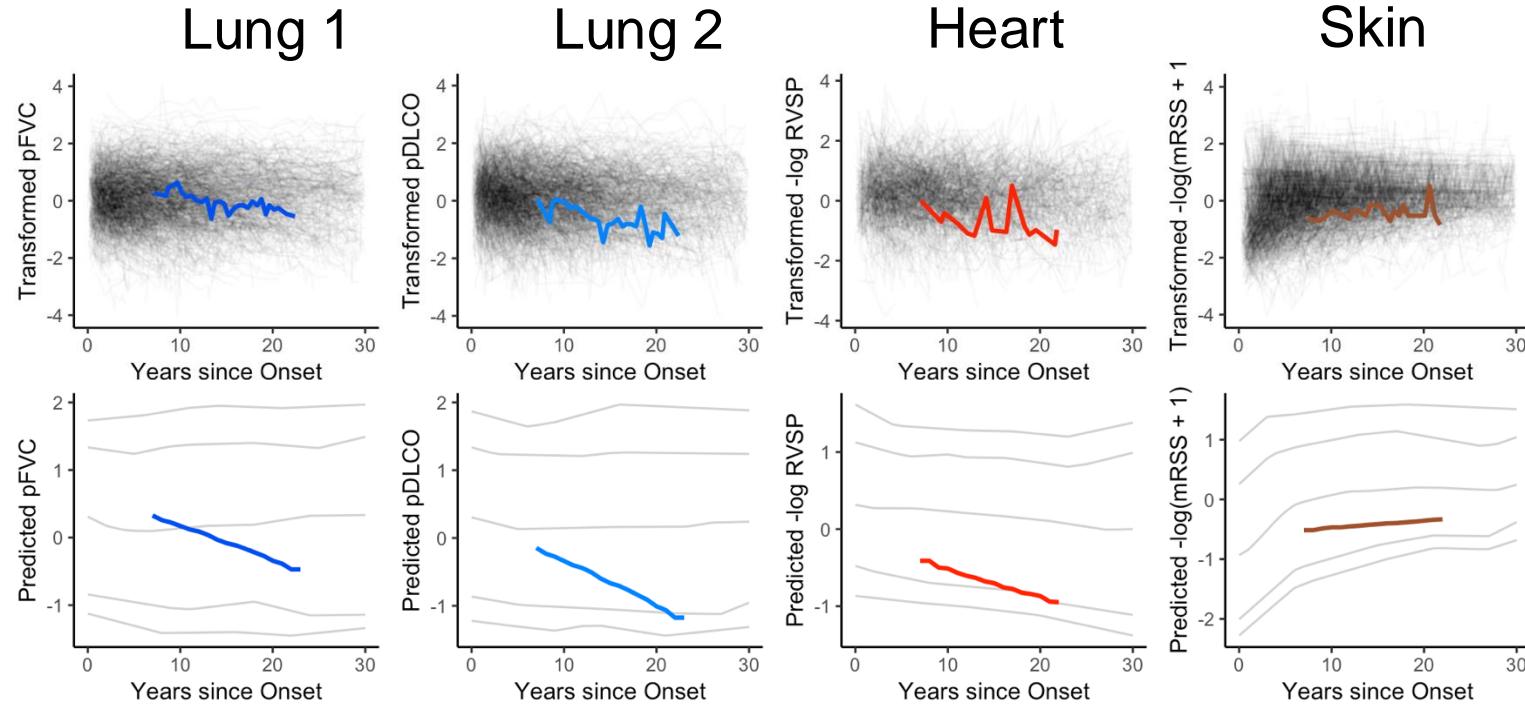
John von Neumann



Thomas Bayes

Disease Mechanism ξ , Measured by X, Causes Health Trajectory η , Measured by Y







Scleroderma CoE

Patient Lookup

@paslanb1

Patient Information Name: "Marie Curie" EMRN: EXXXXXXXX DOB: 02/19/1960 Race: White Sex: Female

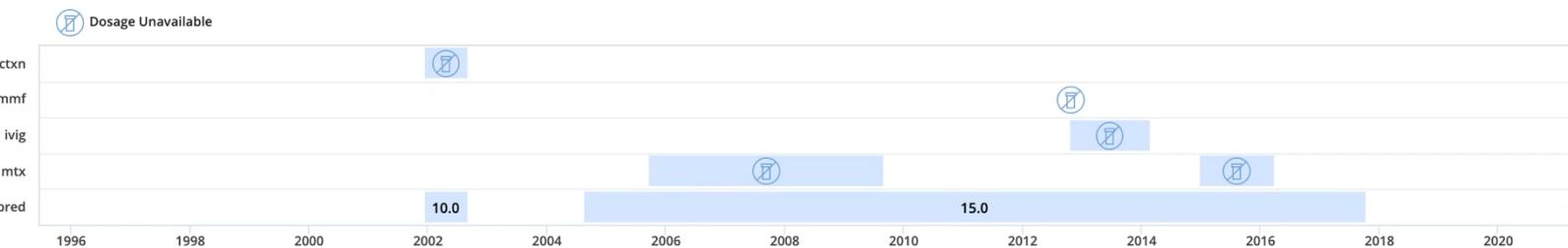
Cohort Filter

Match Patient

Summary

Chart View: All Scleroderma

Medications



Patient Events

Ejection Fraction

RVSP

FVC



Causal Inference using Multivariate Generalized Linear Mixed-Effects Models with Longitudinal Data

Yizhen Xu, Jisoo Kim, Laura K. Hummers, Ami A. Shah, Scott Zeger

March 3, 2023|

2 Notation and Model

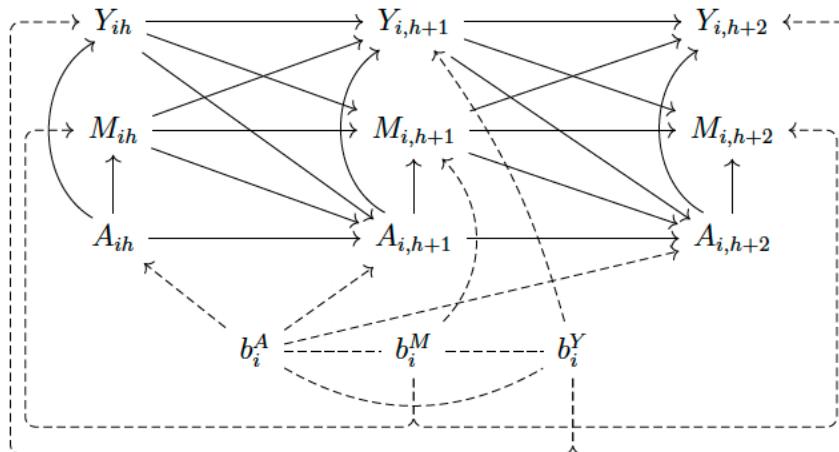


Figure 1: Directed acyclic graph (DAG) for the generalized linear mixed model displaying temporal order of the observed variables and time-invariant unmeasured heterogeneity in both treatment assignment and biomarker dynamics. Baseline characteristics V_i is excluded from the figure for simplicity.

$$Y_{it}|(M_{it} = 1) = \phi_1(\mathcal{H}_{it})\beta_1^Y + \phi_2(\mathcal{H}_{it})\phi_A(\bar{A}_{i,0:t})^T\beta_2^Y + b_{i0}^Y + e_{it}^Y$$

$$\text{logit}\{P(M_{it} = 1)\} = \phi_1(\mathcal{H}_{it})\beta_1^M + \phi_2(\mathcal{H}_{it})\phi_A(\bar{A}_{i,0:t})^T\beta_2^M + b_{i0}^M$$

$$\text{logit}\{P(A_{it} = 1|A_{i,t-1} = 0)\} = \phi_1(\mathcal{H}_{it})\beta_1^A + b_{i0}^A$$

where

$$(b_{i0}^Y, b_{i0}^M, b_{i0}^A)^T \sim N(0, G),$$

$$\phi_1(\mathcal{H}_{it}) = \{1, \tilde{Y}_{i,t-1}, V_i, B_i, ns(S_{it}, \nu_s), B_i \times ns(S_{it}, \nu_s)\},$$

$$\phi_2(\mathcal{H}_{it})\phi_A(\bar{A}_{i,0:t})^T = \{D(\bar{A}_{i,0:t}), V_i \times D(\bar{A}_{i,0:t}), B_i \times D(\bar{A}_{i,0:t}), I_{it} \times D(\bar{A}_{i,0:t})\},$$

and we assume $\nu_s = 4$. Note that both the outcomes Y_{it} and confounders M_{it} are multivariate, i.e. Y_{it} and $M_{it} \in \mathbb{R}^3$. Specifically, $(b_{i0}^Y, b_{i0}^M, b_{i0}^A) \in \mathbb{R}^7$ and the covariance matrix $G \in \mathbb{R}^{7 \times 7}$. Model

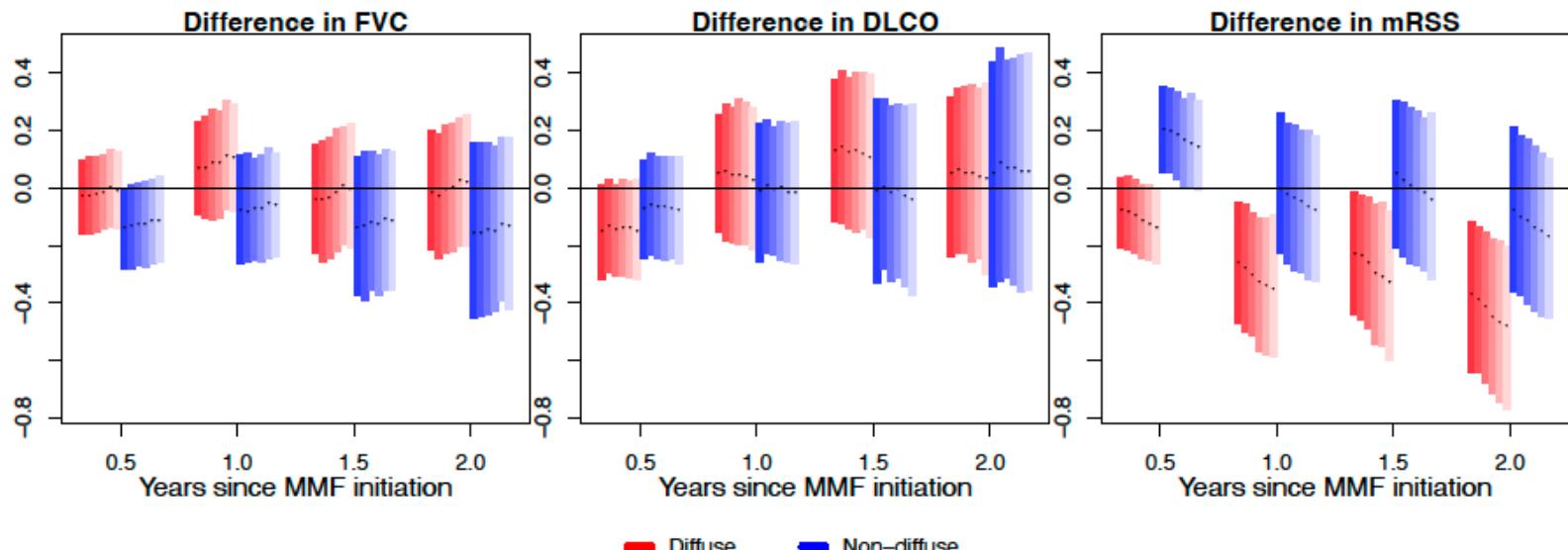


Figure 8: Application causal estimation by subgroup diffuse versus nondiffuse.

Covid-19

Department of Biostatistics All Hands on Deck

John Muschelli, Karen Bandeen Roche, **Mary Grace Bowring**, Martina Fu, **Yizhen Xu**, **Zitong Wang**, Mei-Cheng Wang; Jiyang Wen,
Jamie Perin, **Shannon Wongvibulsin**, Grant Schumock; Josh Betz; **Scott Zeger**

Many Clinical Scientific Colleagues

Brian Garibaldi, Amita Gupta, Matt Robinson; Antony Rosen; Eileen Scully;...

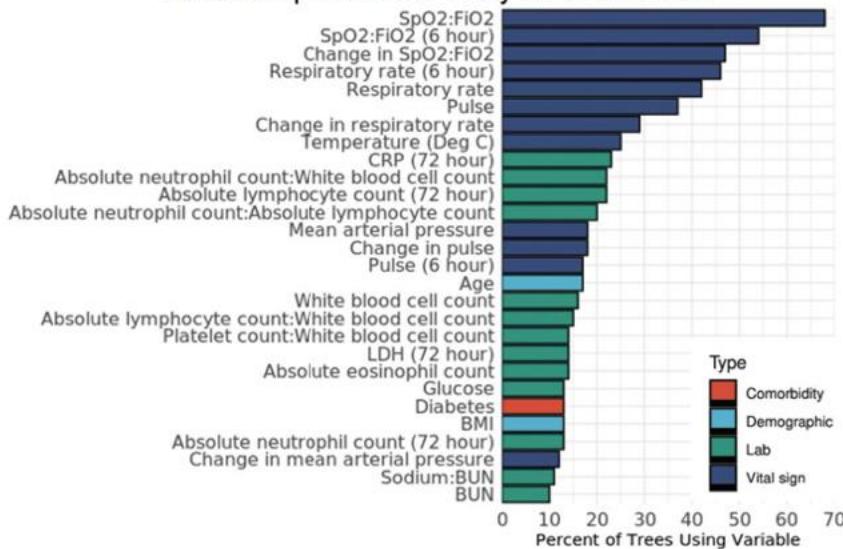
Question: *What is my risk of having severe disease given my current condition?*

- Predict risk of severe outcome from baseline and all intervening measures prior to t with 6-hourly updates
- Random Forest for Survival, Longitudinal and Multivariate (RF-SLAM) outcomes (Wongvibulsin, et al, 2019)
- Approximating tree to explain predictions to clinicians and patients (Wongvibulsin, et al, 2020)
- Severe Covid 19 Adaptive Risk Predictor (SCARP) (Wongvibulsin, et al, 2021)
- SCARP-lite implemented in JH Epic HER (Robinson, et al in 2021)



Results

Variable Importance for 1-Day Risk Predictions



Variable Importance for 1-Week Risk Predictions

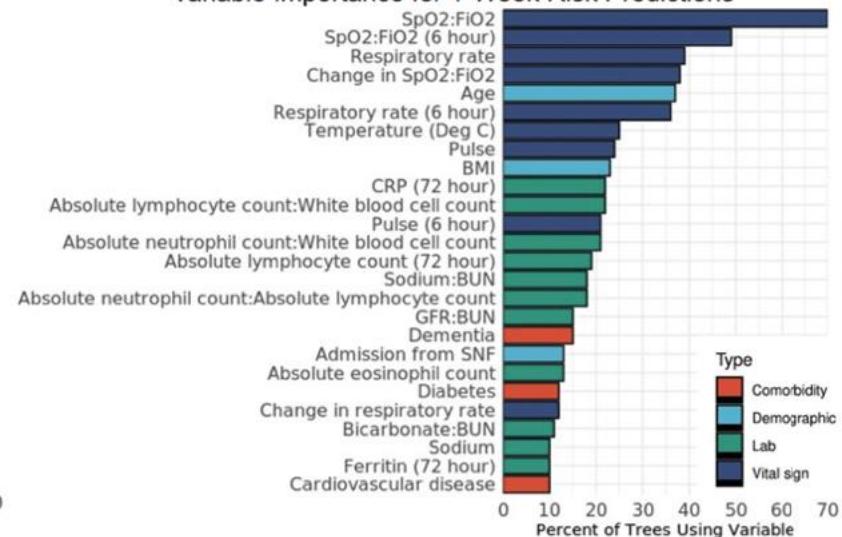
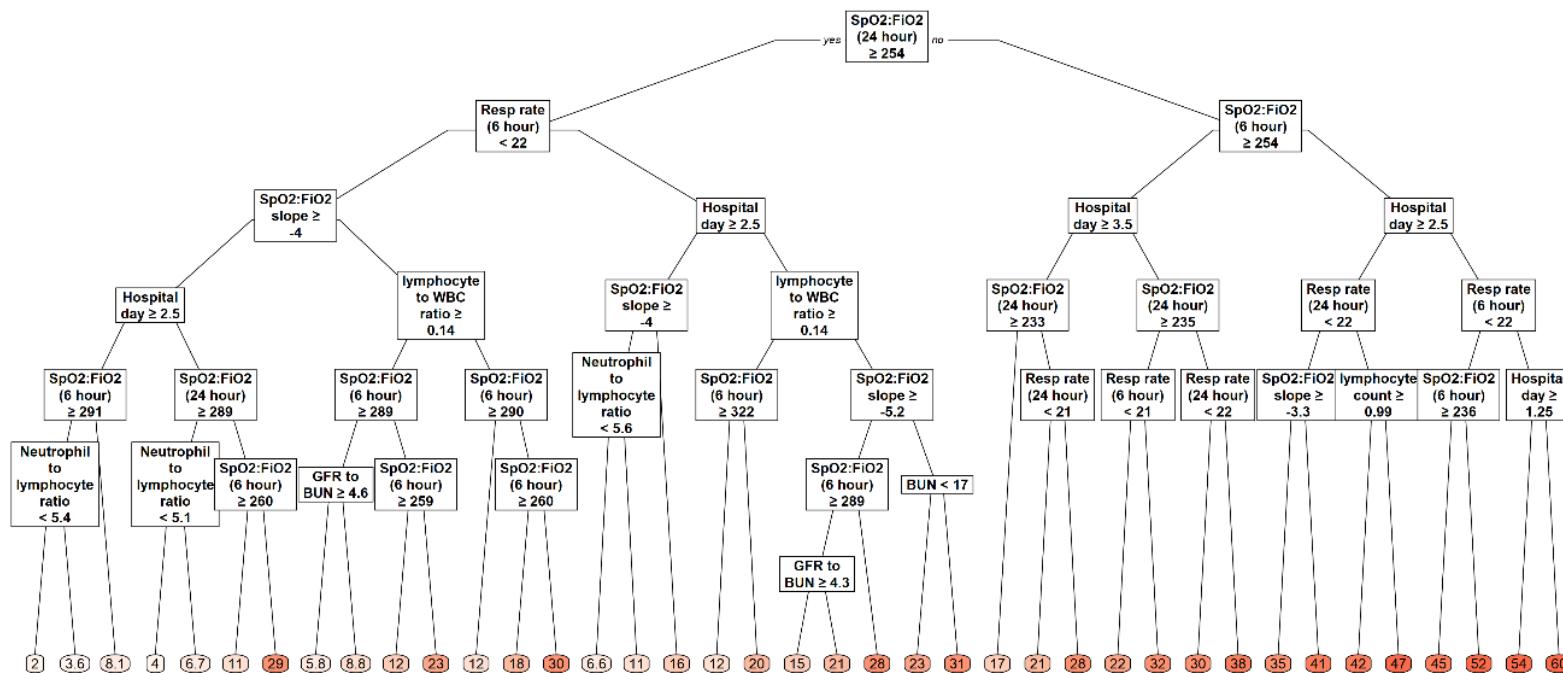
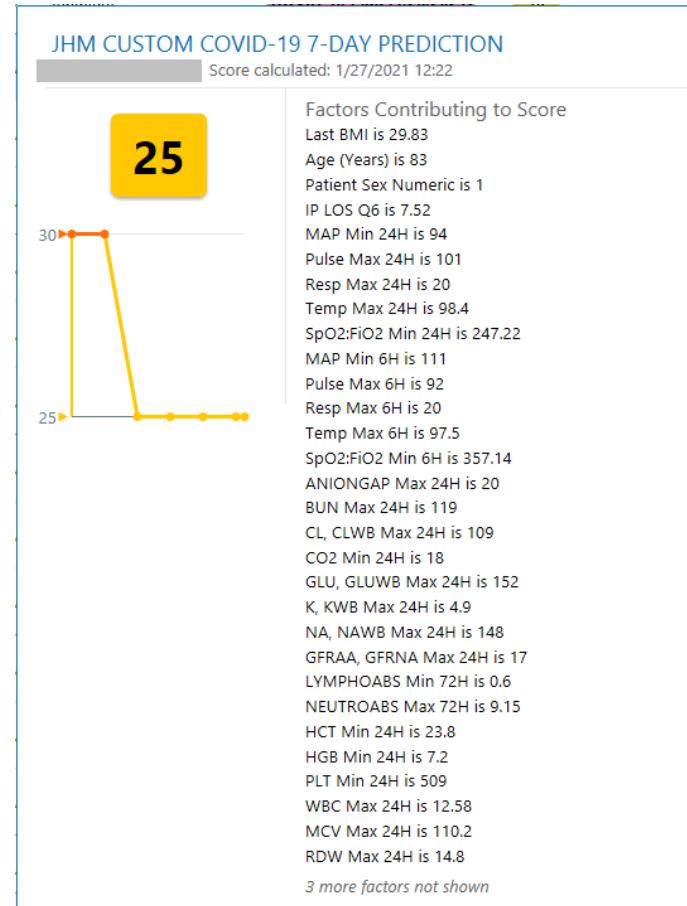
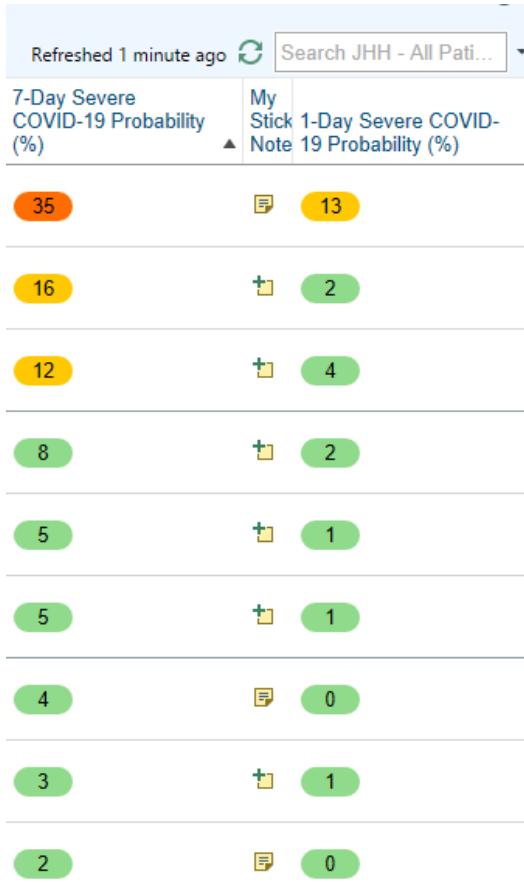


Figure 2: 1-Day risk prediction variable importance plot and 1-week risk prediction variable importance plot: The percentage of trees incorporating each of the variables is used as a simple and interpretable measure of variable importance. The variables used by 10% or greater of the trees are shown in the plots. Note: values for labs and vital signs correspond to values in the past 24 hours unless otherwise specified (e.g., 6 hour indicates that the value corresponds to the past 6 hours).

Summary of RF-SLAM predictions



Summary tree of RF-SLAM predictions of 1-week risk of severe disease or death. The predicted probabilities are expressed in the terminal nodes and shaded according to lowest risk (0%) to highest risk (100%) prediction



Severe COVID-19 Adaptive Risk Predictor (SCARP)

The COVID-19 adaptive risk predictor (SCARP)* is an online tool that calculates the 1-day and 7-day risk of progression to severe disease or death for patients hospitalized with COVID-19.

Instructions

Enter the information for the patient below into the orange box. Inputs will be entered sequentially (additional boxes will appear as you enter information). The sequential inputs are determined adaptively based on the information entered in order to tailor the calculator to the individual patient. The 1-day and 7-day risk predictions and visual displays of summary decision trees appear at each step. Additional information regarding the development of SCARP can be found [reference to manuscript].

Clinical predictors

Respiratory rate (highest in past 6 hours)



Days since hospital admission



SaO₂:FiO₂ (24-hour min): 228

Enter the supplemental oxygen and pulse oximetry recorded at the most hypoxic moment in the past 24 hours

Supplemental Oxygen Delivery (L/min) (24)



Oxygen Saturation by Pulse Oximeter (24)



GFR: 41

Risk of severe illness or death

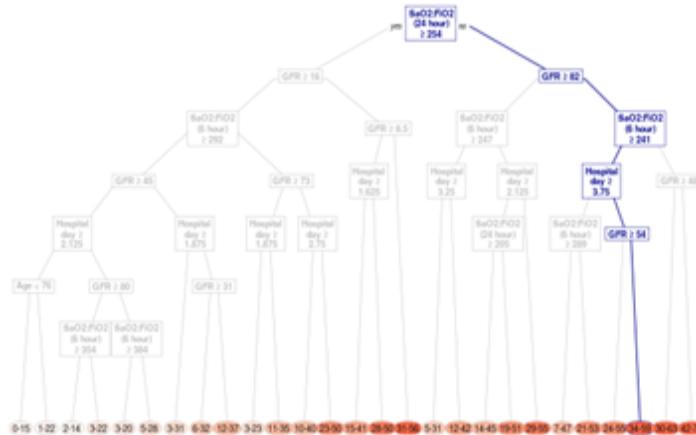
14%

in the next day



45%

in the next week

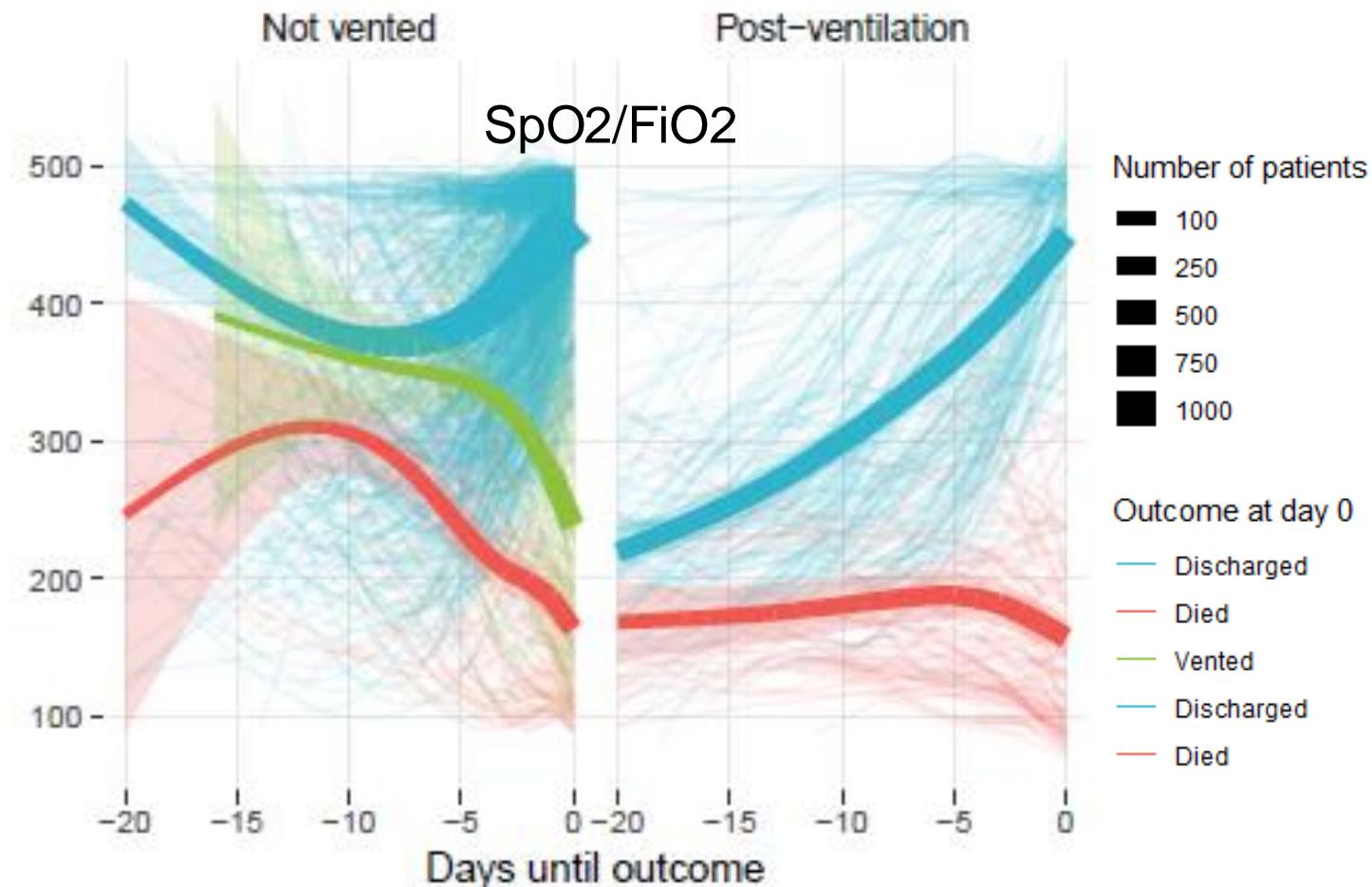


Question: Given my current state, what are my future risks of discharge, intubation or death; what are my expected biomarker trajectories for each outcome?

- Competing discrete hazards of (3) events on each future day (Project 1 approach)
- Retrospective longitudinal data analysis of multiple biomarkers *given event outcome* (define t=0 at event time) (Bowring, MG, Wang, Z et al, 2021)
- Bayes rule to calculate the probability of a future event given baseline and biomarker data until current time.



Longitudinal Data Analysis



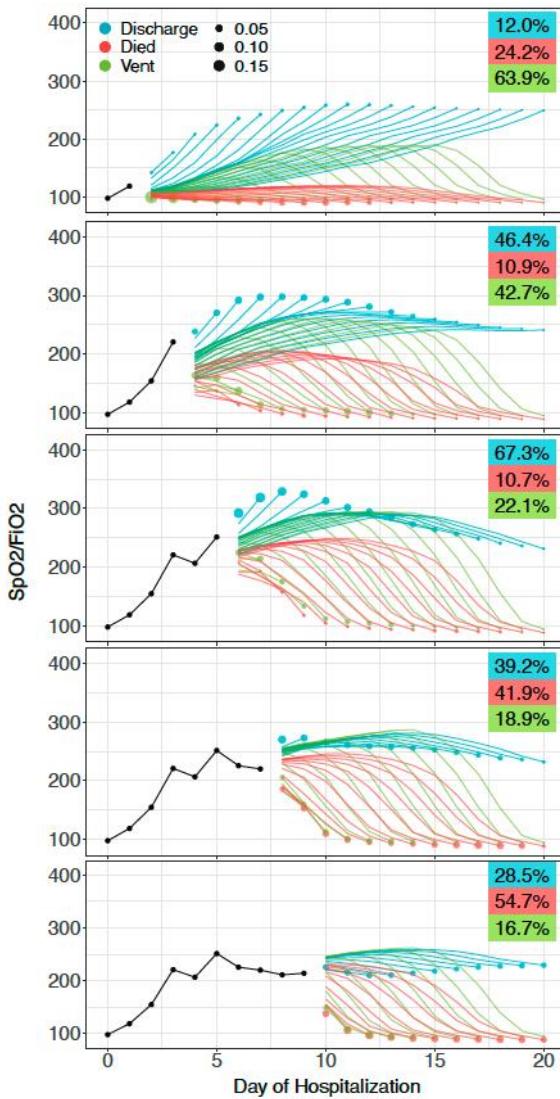


FIG 4. Prediction capabilities of the retrospective model based on data from the first 2, 4, 6, 8 and 10 days of hospitalization for a given patient. Black lines represent the observed $\text{SpO}_2/\text{FiO}_2$ for the patient during hospitalization. Blue, red, and green lines indicate the predicted $\text{SpO}_2/\text{FiO}_2$ trajectories leading to discharge, death, or ventilation events respectively. The marker size at the end of each predicted trajectory indicates the probability of that event on that day given the past biomarker history. The probabilities of each event happening eventually (before day 20) are shown in the colored boxes on the right, with colors indicating event types.

Statistical Challenges

1. Wrangling gigabytes of transactional EHR data into a longitudinal data set for thousands of patients; limitations of measurements
2. Brand new disease without a clinical evidence base
3. Observational, not experimental data
...=> Treatments(i,t-1) => Outcomes(i,t) => Treatments(i,t) =>...
4. Outcomes comprise many biomarkers, major events and treatment choices
5. Competing risks of three major events: discharge, intubation, death
6. Significant fraction of deaths complicate off-the-shelf longitudinal data analysis approaches
7. Predictors are numerous and dynamic; many are irrelevant; need to find the important ones
7. Potentially useful results are unhelpful until translated into terms clinicians can understand to improve their decisions

References

- Rosen, A. and **Zeger, S.L.**, 2019. Precision medicine: discovering clinically relevant and mechanistically anchored disease subgroups at scale. *The Journal of clinical investigation*, 129(3), pp.944-945.
- Wu Z**, Rosen L, Rosen A, **Zeger SL**. 2020. A Bayesian approach to Restricted Latent Class Models for scientifically-structured clustering of multivariate binary outcomes. *Biometrics*, 2020.
- Wu Z**, Casciola-Rosen L, Shah A, Rosen A, **Zeger SL**. 2019. Estimating autoantibody signatures to detect autoimmune disease patient subsets. *Biostatistics*. PMID: 29140482, 2017.
- Garibaldi BT, **Fiksel J**, et al. 2021. Patient trajectories among persons hospitalized for COVID-19: a cohort study. *Annals of internal medicine*.
- Wongvibulsin S**, Garibaldi BT et al, 2021. Development of Severe COVID-19 Adaptive Risk Predictor (SCARP), a Calculator to Predict Severe Disease or Death in Hospitalized Patients With COVID-19. *Annals of internal medicine*.
- Bower MG**, et al. Outcome stratified analysis of biomarker trajectories for patients with SARS-CoV-2 infection. to appear in *American Journal of Epidemiology*
- Wang Z**, **Bowring MG**, Rosen, A, Garibaldi B, **Zeger S**, Nishimura A. 2022. Learning and predicting from dynamic models for COVID-19 patient monitoring. *Statistical Science*, 37(2), pp.251-265.
- Kim JS**, Shah AA, Hummers LK, **Zeger SL**. Predicting clinical events using Bayesian multivariate linear mixed models with application to scleroderma. *BMC medical research methodology*. 2021 Dec;21(1):1-2.

Thank you



This Photo by Unknown Author is licensed under [CC BY-SA NC](#)

Autoimmune Disease Auto-antibody “Signature”

**Jisoo Kim, Livia Casciola-Rosen, Antony Rosen,
Ami Shah, Zhenke Wu, Scott Zeger**

Disease Mechanism (ξ_i) Causes Health Trajectory (η_{it})

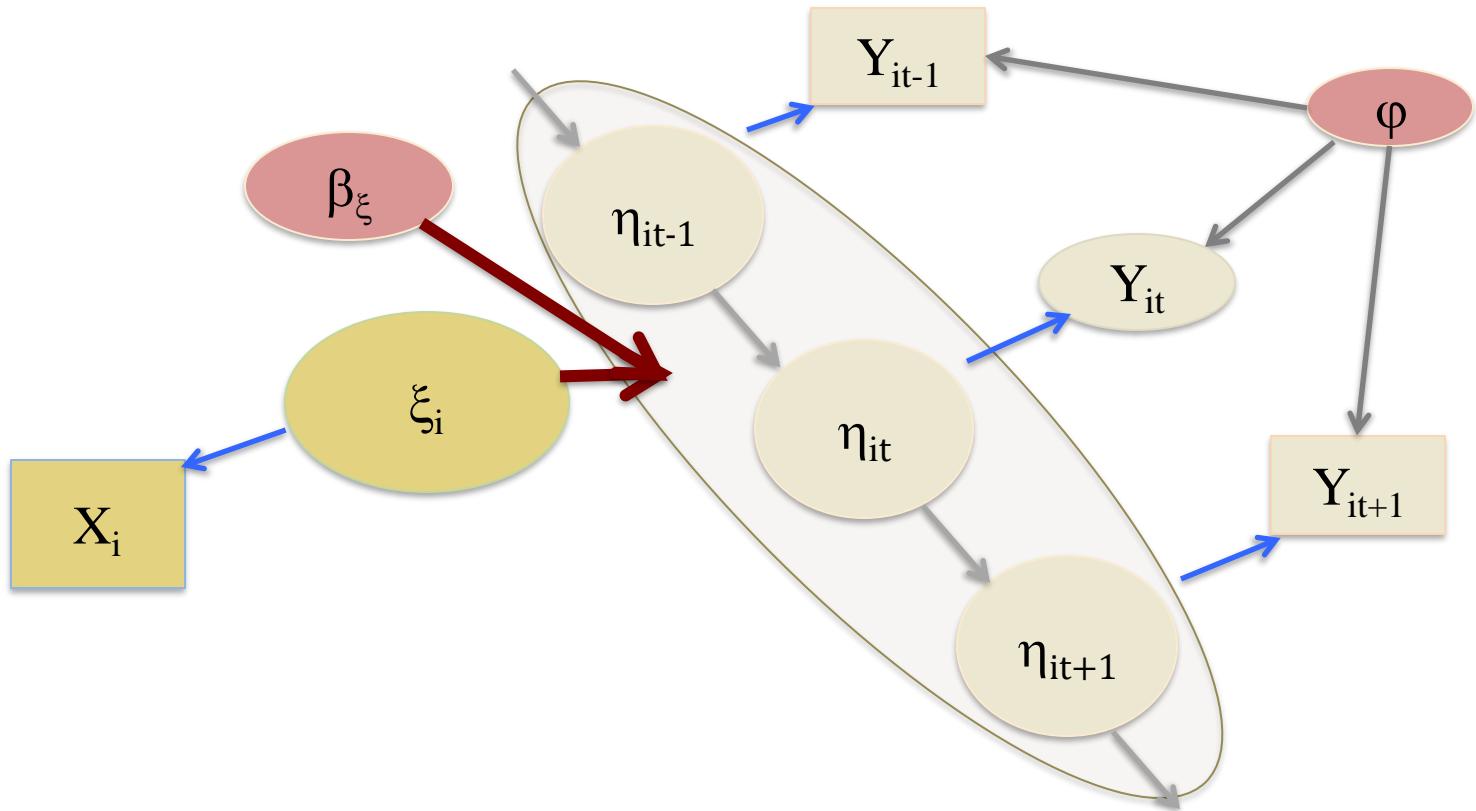
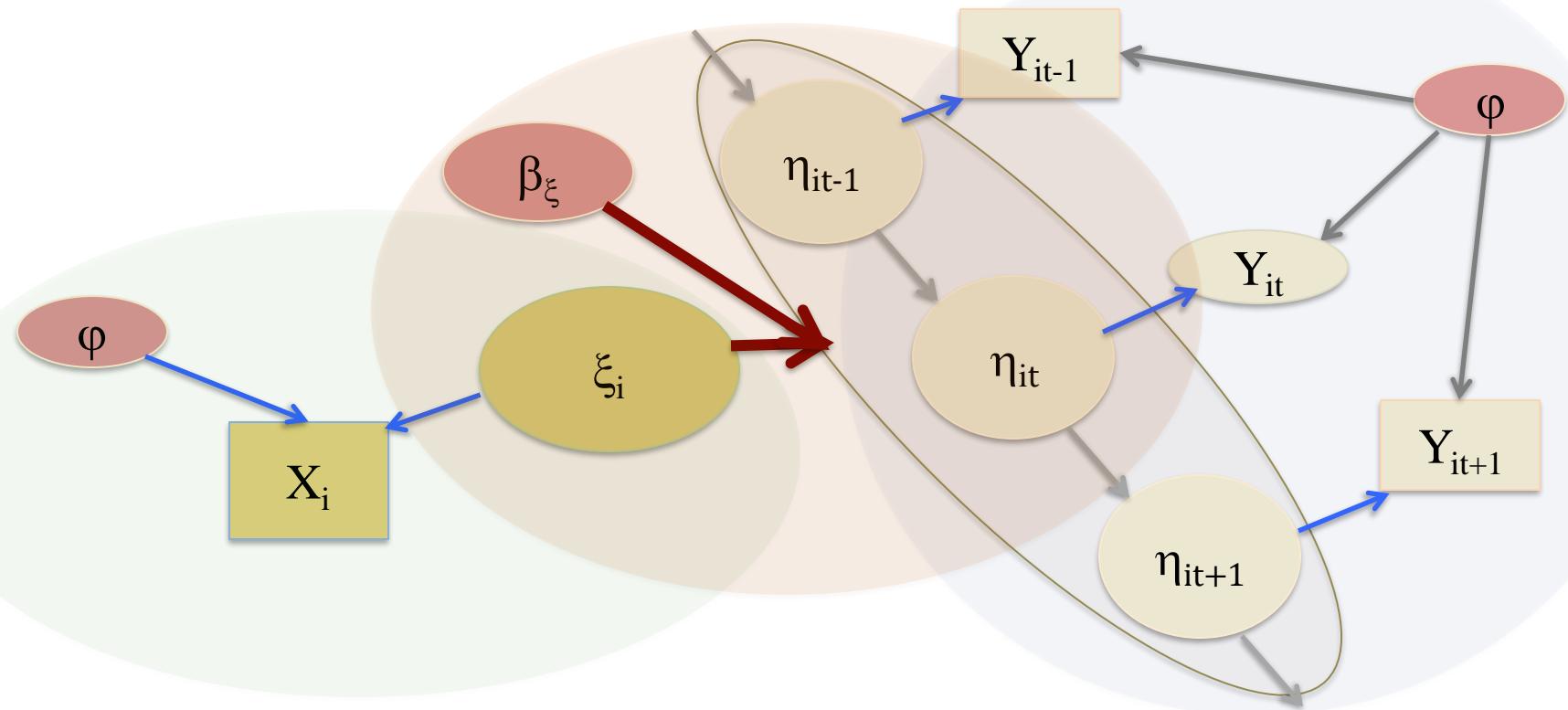
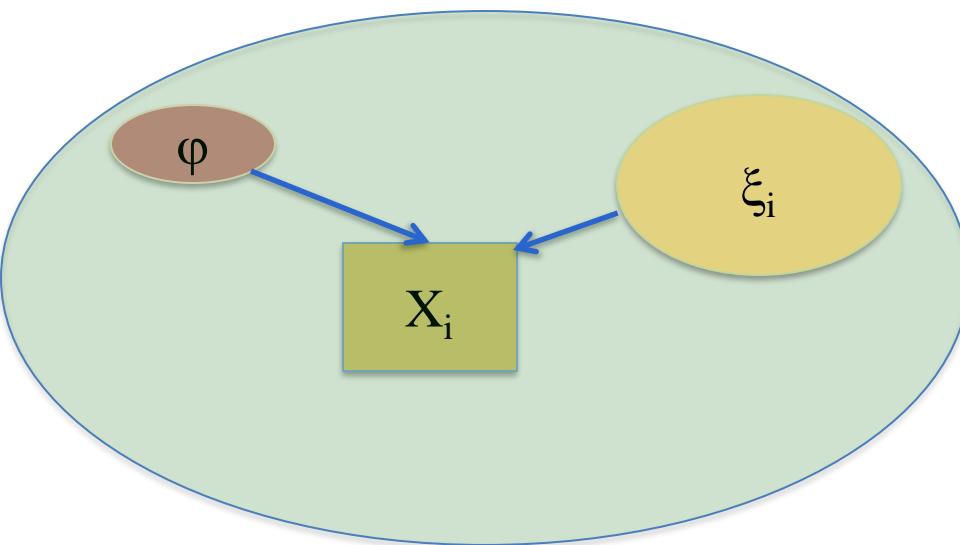


Figure 1. Disease mechanism ξ , measured by X , causes health trajectory η , measured by Y – *Structural Equation Model*



Disease Mechanism ξ , Measured by X, Causes Health Trajectory η , Measured by Y

Restricted latent class analysis => $[\xi|X]$ – Zhenke Wu



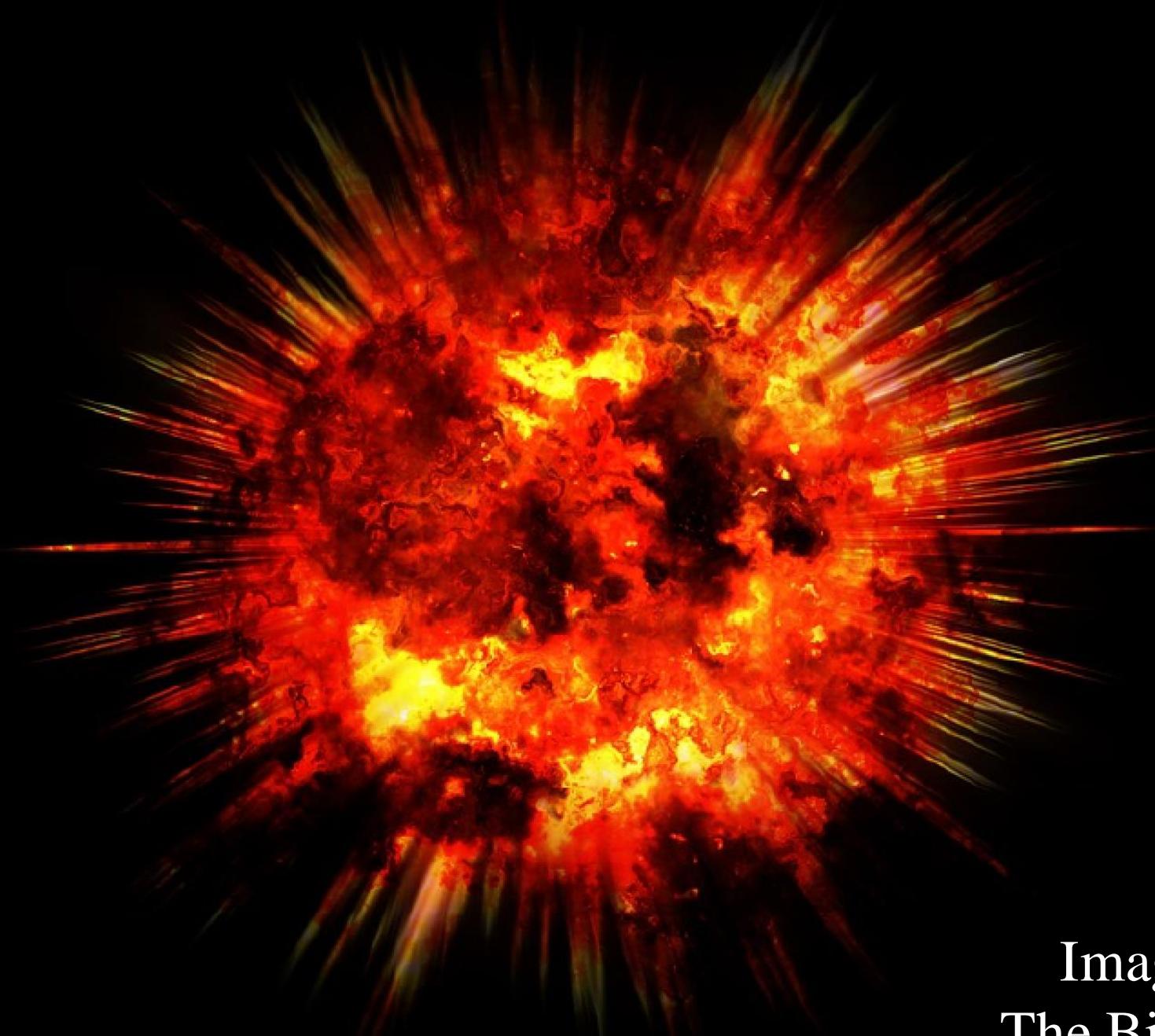
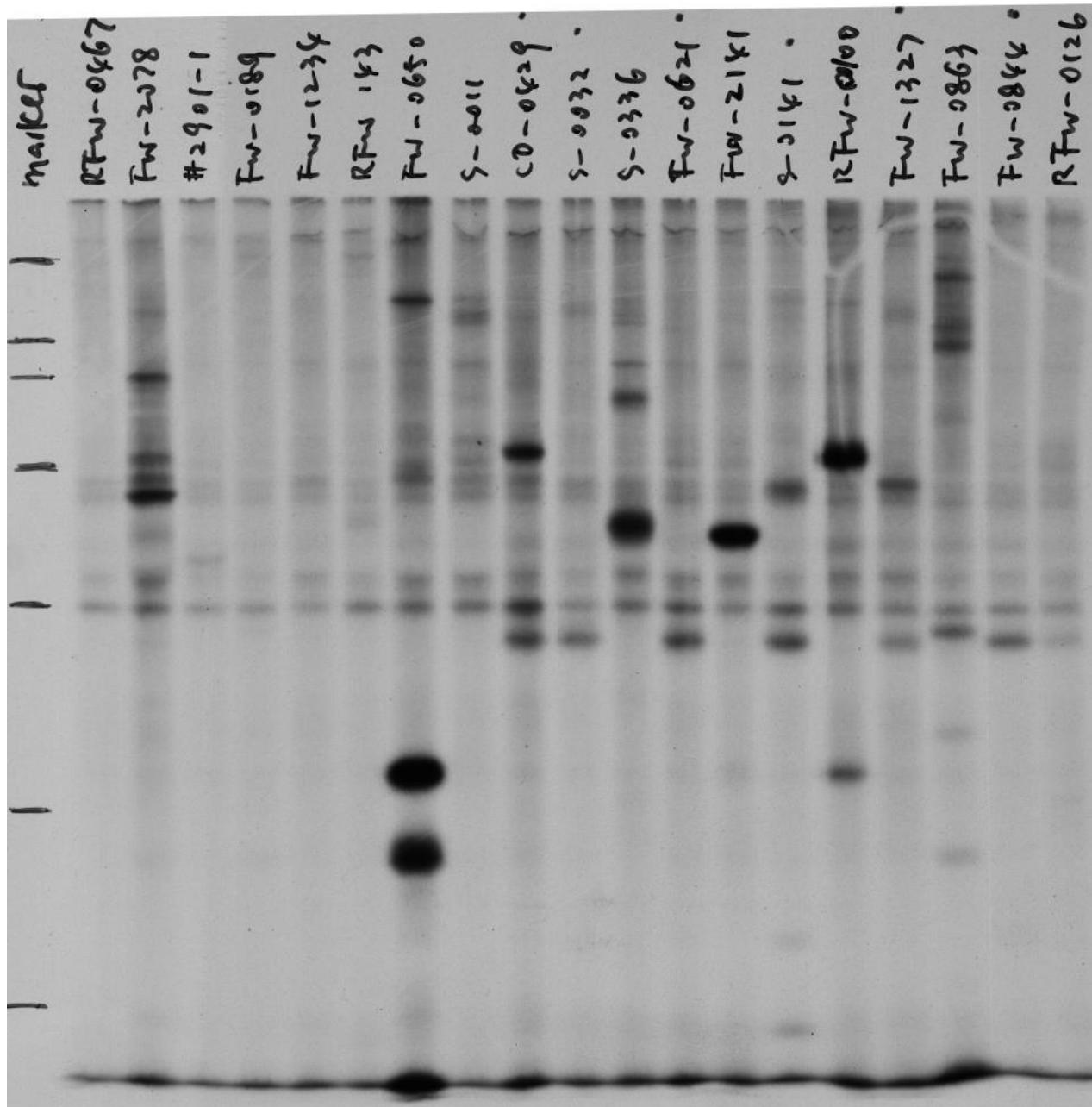
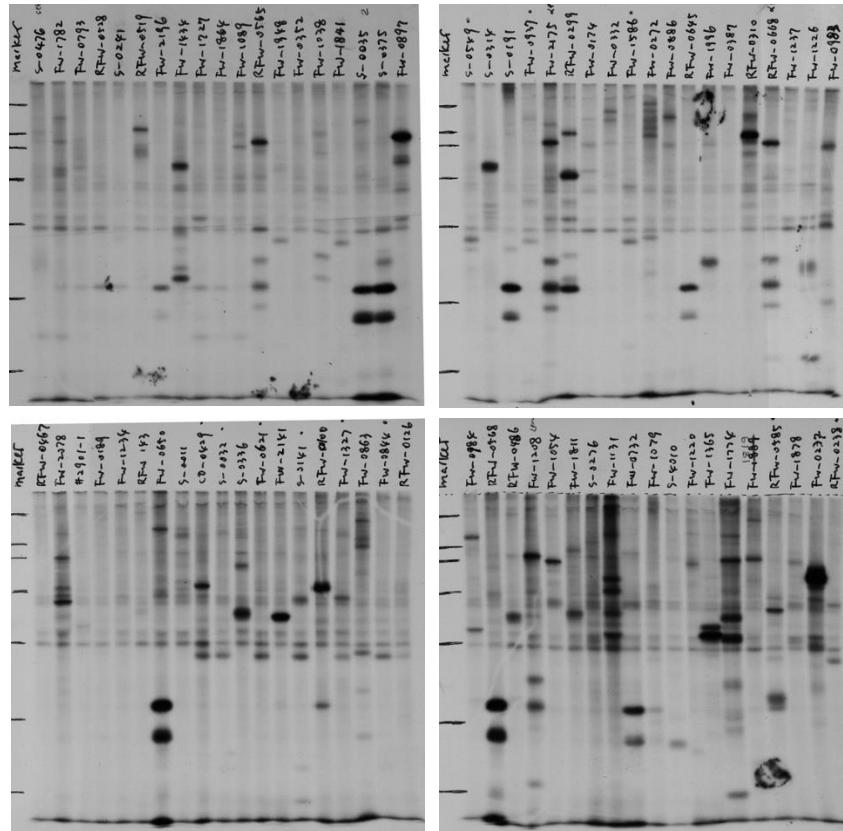


Image of
The Big Bang

Autoantibody profiles in patient sera – immunoprecipitation readout

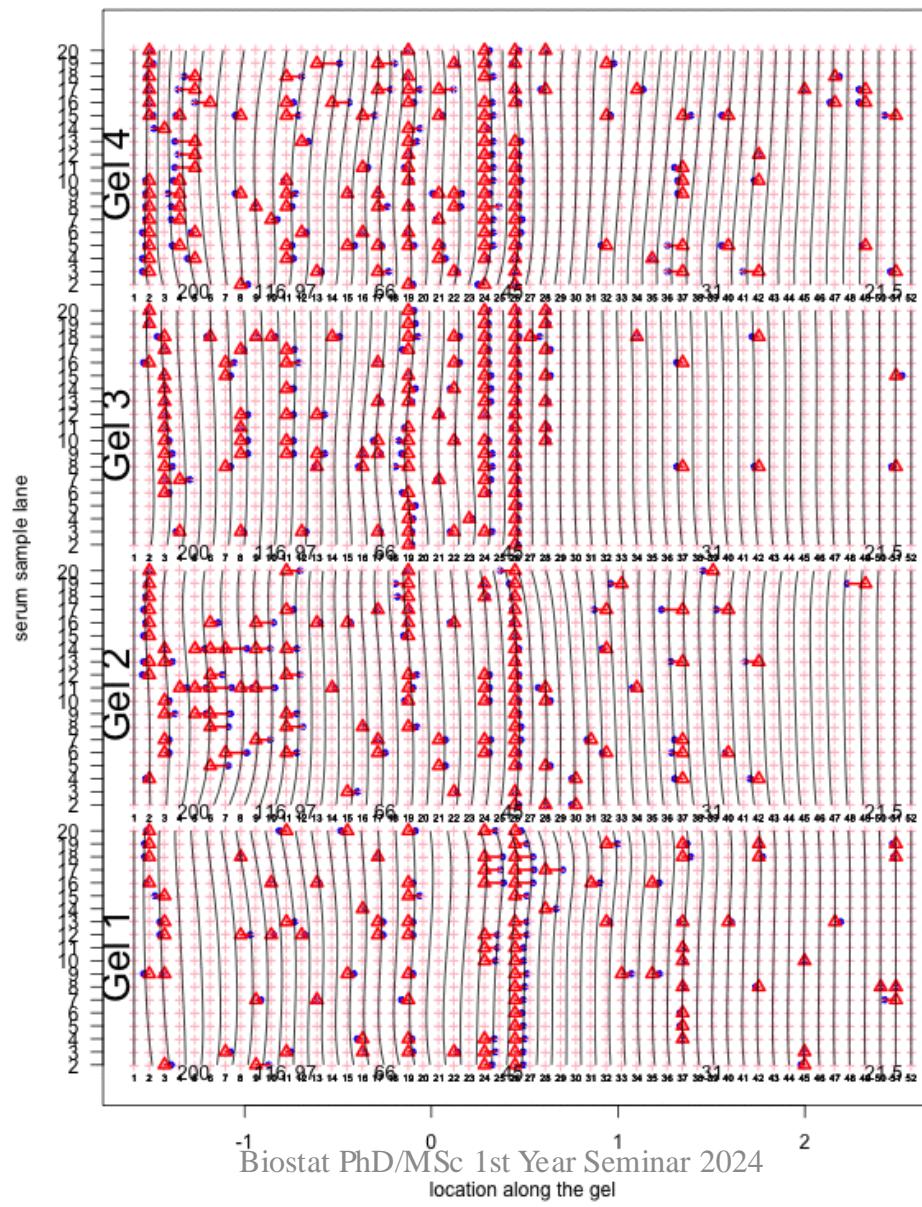


- 76 autoantibody patterns from patients with rheumatic disease & cancer
- all were negative for autoantibodies against prominent defined specificities

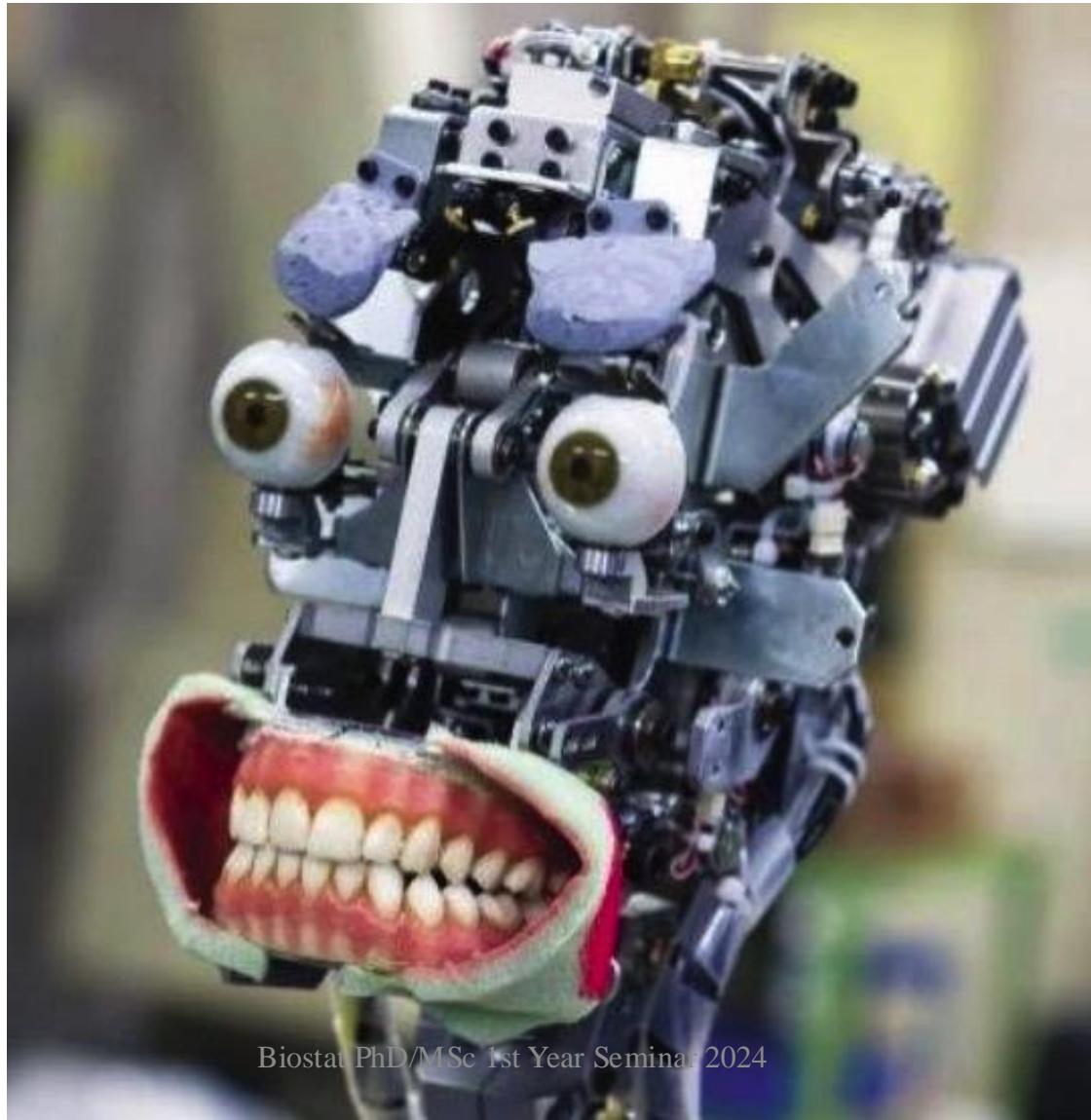


Bayes Rule to identify common autoantibody signatures?

Align the peaks

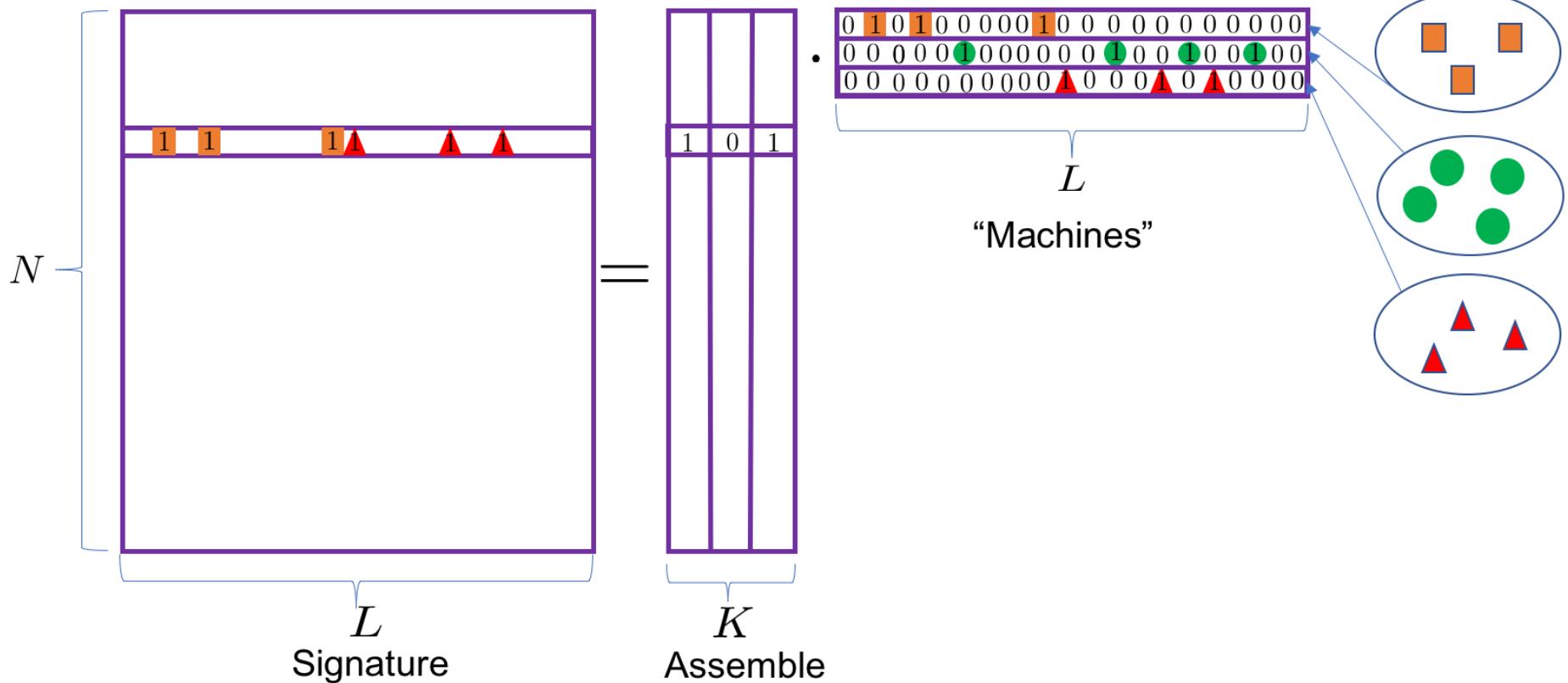


Cellular Machines Model – Machine is a set of proteins that do a cellular job; machines are targeted by the immune system



Clustering informed by the biology that antibodies target multiple components in complexes (intermolecular spreading)

$$\xi_i = \eta_i^T \cdot Q$$



Preliminary clustering results based on machine models

Data: CTP negative sera

Method: Bayesian machine-based algorithm

Figure: Three estimated clusters (top three panels) with distinct enrichment of three distinct estimated machines (bottom panel)

Colored labels: red, blue, green - for clusters obtained by eyeballing; the algorithm is agnostic to them.

