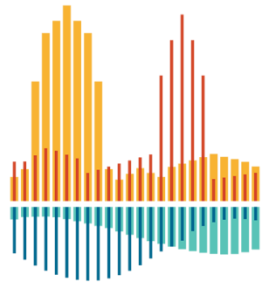


# LARGE-SCALE PROCESSING OF GENOMIC DATA COUPLED TO ELECTRONIC HEALTH RECORDS



INSTITUTE FOR  
COMPUTATIONAL  
BIOLOGY

William S. Bush, PhD, MS

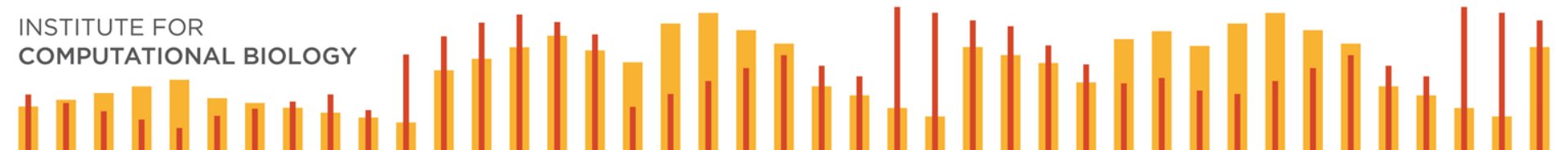
Assistant Professor

Department of Epidemiology and Biostatistics

Institute for Computational Biology

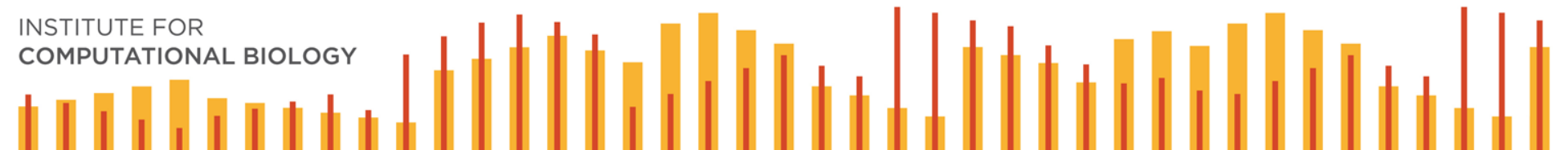
# DISCLOSURES

- I have no financial conflict of interests to report
- I liked data before it was big...



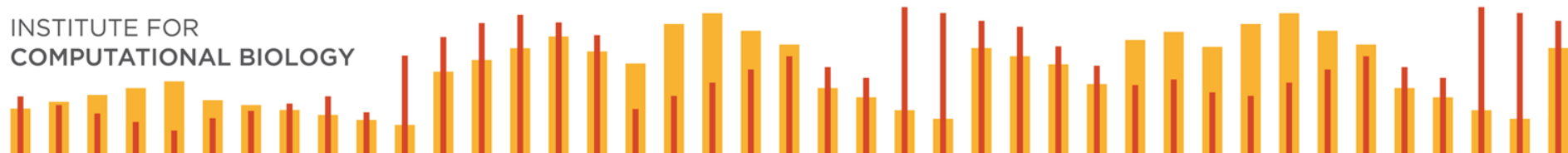
# OVERVIEW

- Use of Electronic Health Records for Research
- Applications of High Performance Computing to EHR Research
- Large-scale Genomic Analysis



# USE OF ELECTRONIC HEALTH RECORDS FOR RESEARCH

INSTITUTE FOR  
COMPUTATIONAL BIOLOGY



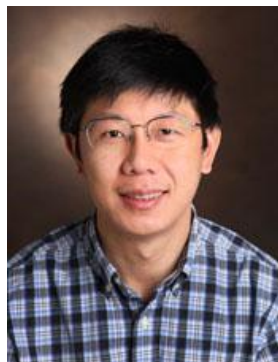
# THE INSTITUTE FOR COMPUTATIONAL BIOLOGY



Jonathan Haines



Ricky Chan



Chun Li



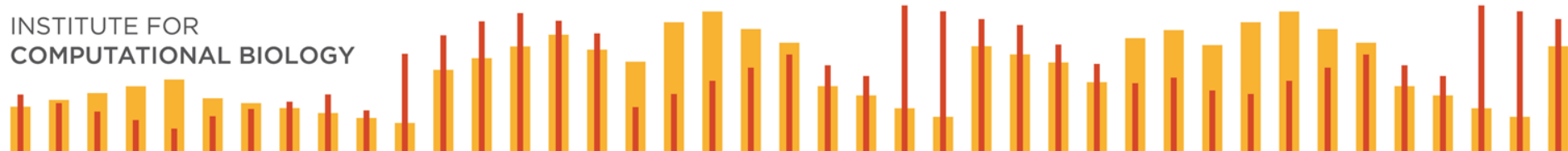
Dana Crawford



Will Bush

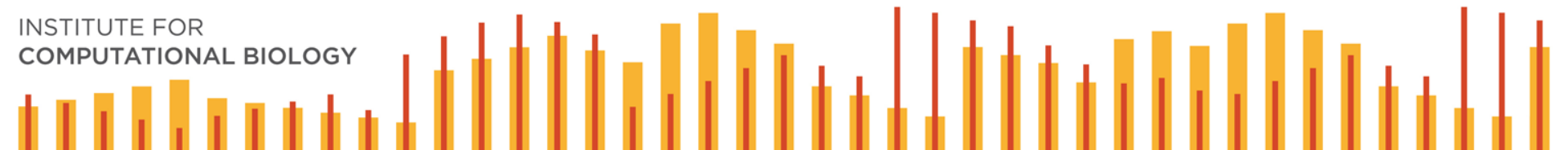


Jill Barnholtz-Sloan



# ELECTRONIC HEALTH RECORDS

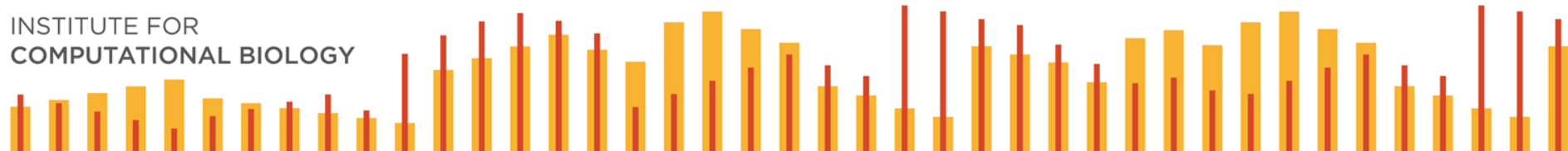
- Health Information Technology for Economic and Clinical Health (HITECH) Act of 2009
- Incentivizes the nationwide adoption of EHR systems in the US
- While there are standards and commonalities, EHR systems are very heterogeneous
- Epic Systems, Allscripts, Meditech, Cerner, IBM, McKesson, Siemens, GE Healthcare
- Epic is most common at large medical centers



# THE CLEVELAND METRO AREA



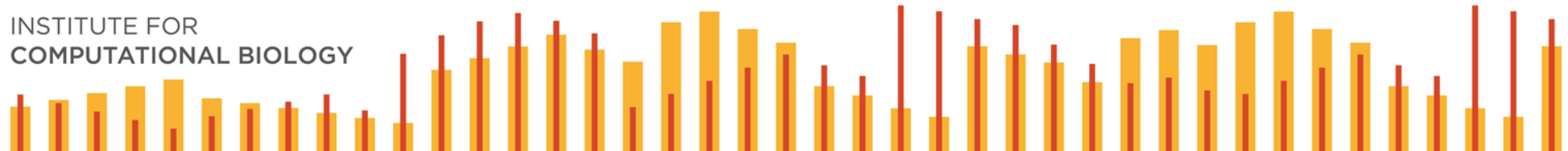
Approaching ~90% of all healthcare in the  
Cleveland Metropolitan Area



# ELECTRONIC HEALTH RECORD COMPONENTS

- Structured Elements
- Pseudo-Structured Text
- Unstructured Text

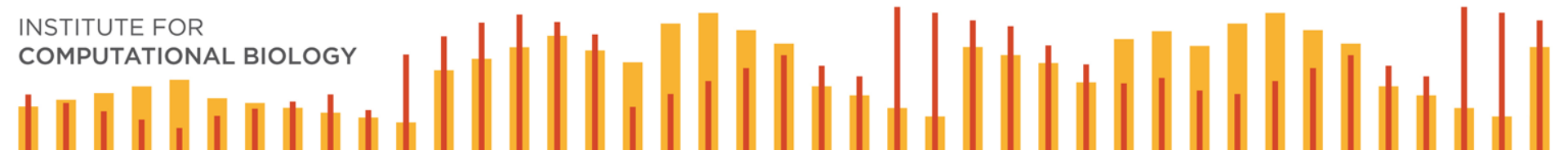
EHR elements can differ between systems and clinics





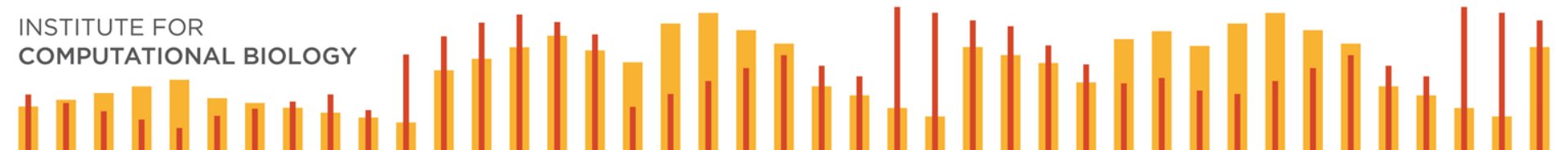
# STRUCTURED ELEMENTS

- Demographic information
  - Date of birth, race/ethnicity, gender
- Vital Signs
  - Heart rate, blood pressure, height, weight, body temperature
- Some Laboratory Values
  - WBC, insulin, glucose, glomerular filtration rate (GFR)
- Billing and Procedure codes
  - ICD9/10, CPT, ICD-O-3



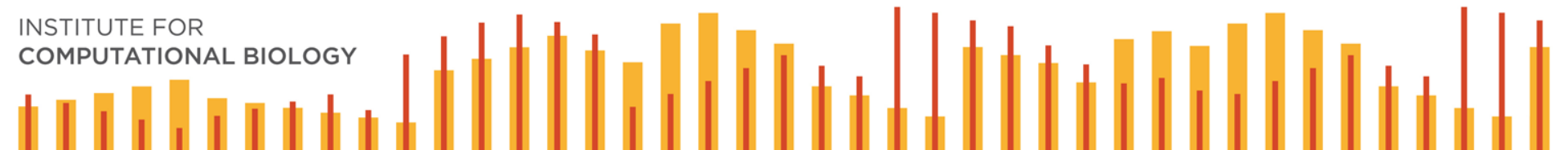
# PSEUDO-STRUCTURED TEXT

- Free text documents with a loosely standardized format
- History and Physical Examination (HPE)
- *“... approximately 10 months status post bilateral bunionectomy with metatarsal head resections 2 through 5 bilaterally. She denies any pain in her feet although she is a little upset that she has had a recurrence of her bilateral hallux valgus, left worse than the right.”*
- *“She does have some residual hallux valgus, worse on the left than on the right, but this is not uncommon following bunionectomy. We discussed the fact that her bunions were so bad to begin with, that her feet actually look pretty good.”*
- Problem List / Known significant medical conditions and procedures / Allergies



# UNSTRUCTURED TEXT

- Clinical communications (phone calls, prescription refills, etc)
- Discharge notes
- Clinic-specific notes
- Laboratory / Procedure reports (CT scan, colonoscopy, etc)
- “Medical flotsam”

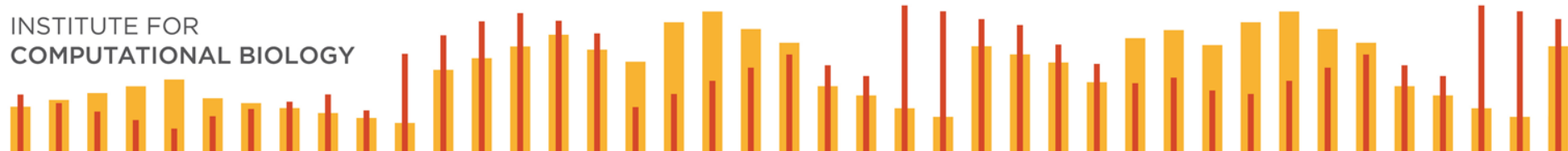


# MY EHR EXISTENTIAL CRISIS

At some point in everyone's life, they realize that...

- Their parents don't know everything.
- Their doctor doesn't know everything.
- Their medical record may be wrong.

EMRs are very useful, but they can be noisy and inaccurate.



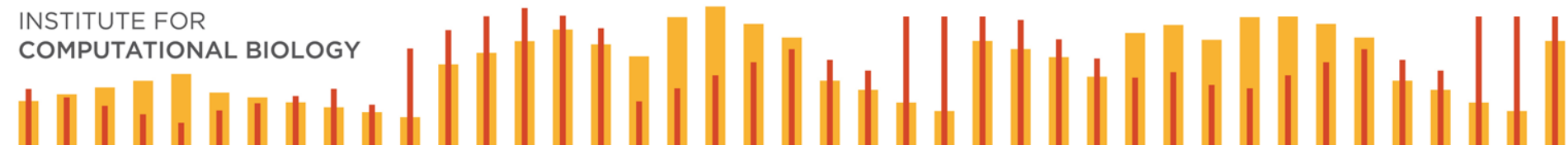
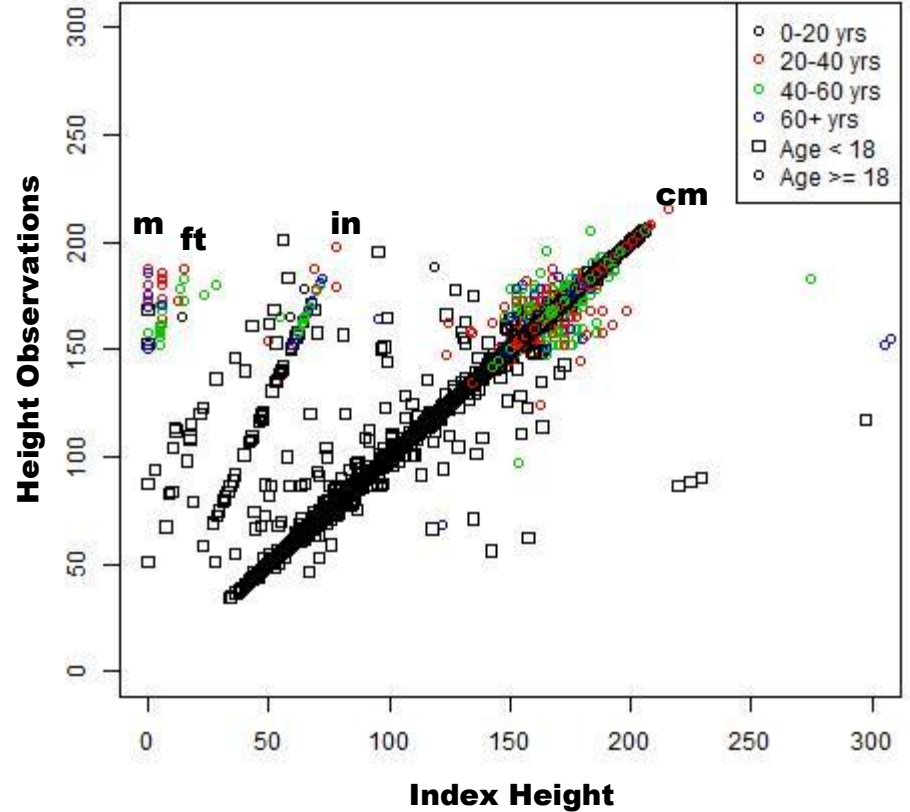
# EXAMPLE: BODY MASS INDEX

Height and Weight are the most ubiquitous measures reported in an EHR



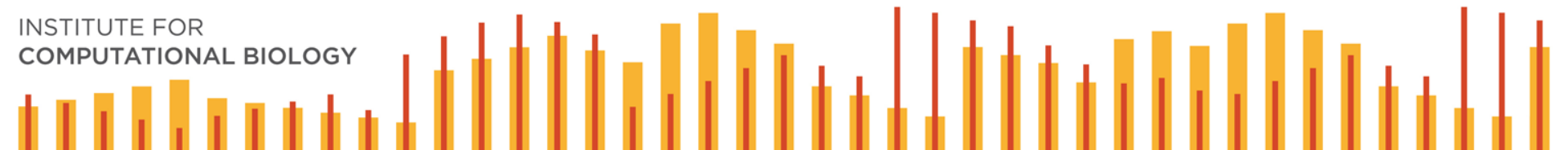
Robert Goodloe

Eagle-BioVU, Vanderbilt University Medical Center



# CLINICAL VERSUS RESEARCH USE

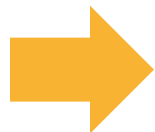
- In clinical practice, records and values are examined individually
  - Errors are easy to spot and ignore
  - Aggregate data is synthesized by experts in the context of a patient
- In research, records and values are examined in aggregate
  - Errors are not so easy to identify
  - Data points are examined outside the context of the patient



# STRUCTURING DATA FOR RESEARCH USE

record_id	icd_code	Date
24067	729.5	40909
24067	070.30	40910
24067	783.1	40910
24067	45.23	40911
24067	77	40911
24067	455.6	40912

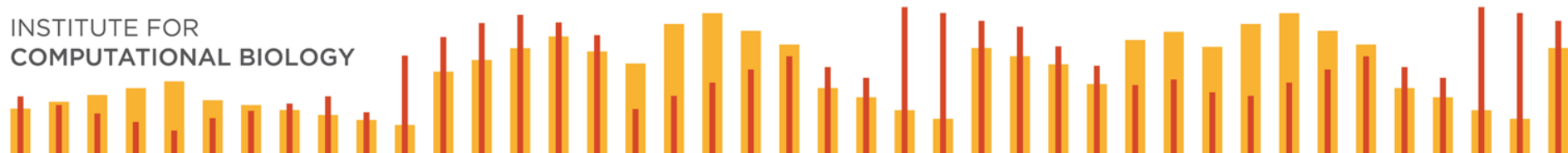
record_id	vital_sign	value	Date
24067	hr	72	40912
24067	sbp	160	40912
24067	dbp	92	40912
24067	height	150	40912



“patient started Simvastatin 800 mg”

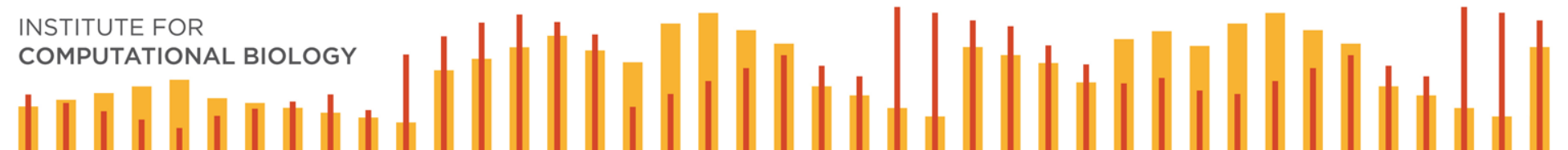
“42 year old male, overweight”

OfficeEMR by iSALUS Healthcare



# APPLICATIONS OF HIGH PERFORMANCE COMPUTING TO EHR RESEARCH

INSTITUTE FOR  
COMPUTATIONAL BIOLOGY





# EXAMPLE: CALCINEURIN-INHIBITOR TOXICITY

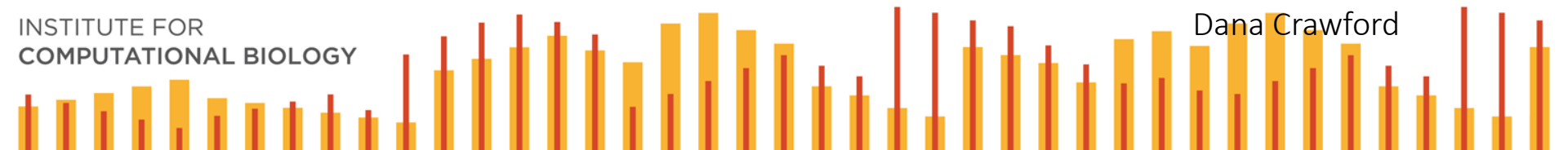
- Given post-transplant to prevent organ rejection
  - Tacrolimus and cyclosporine
- Narrow therapeutic window
- Nephrotoxicity is a serious and common complication
- Serum creatinine and glomerular filtration rates (GFR) are monitored post-transplant to assess kidney function



Matt Oetjens

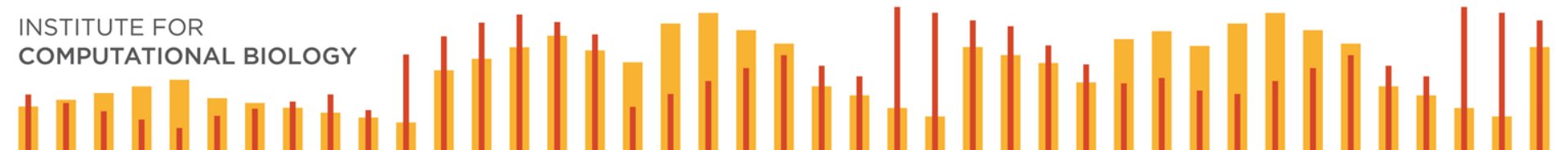


Dana Crawford



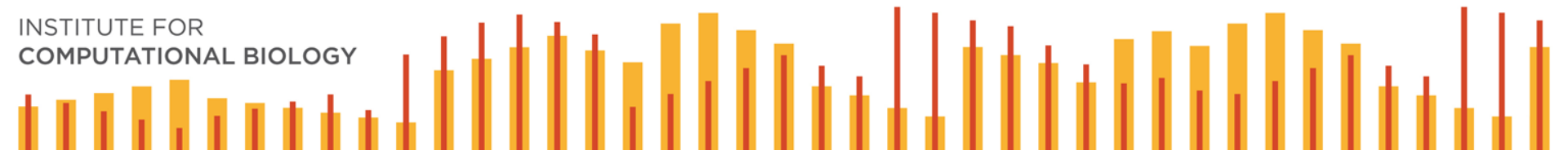
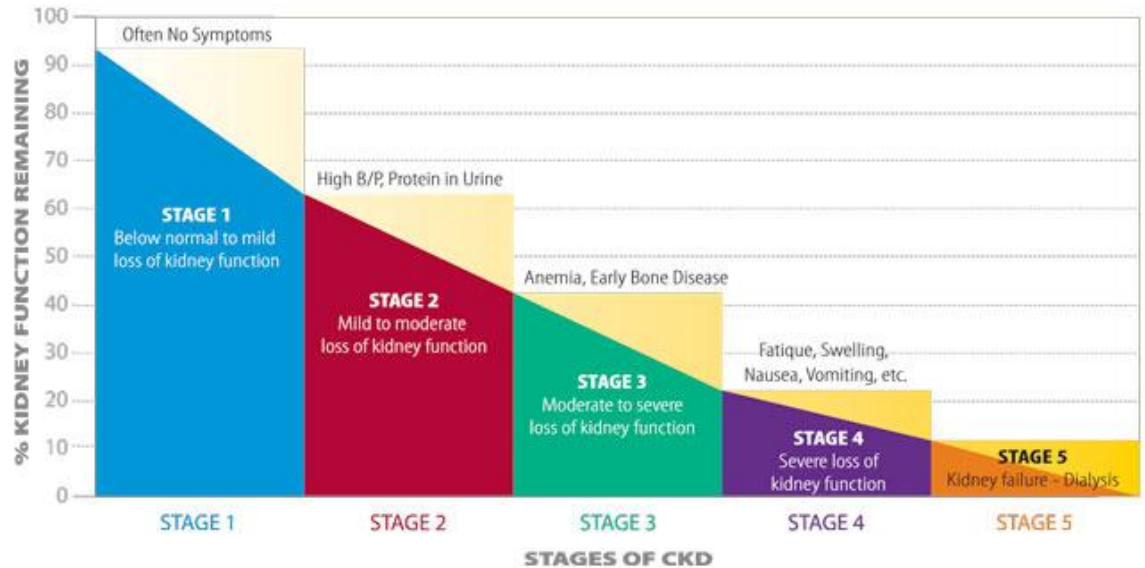
# STUDY DESIGN

- Patients identified having:
  - Heart transplant documented with  $\geq 3$  ICD9 Code V42.1 (heart replaced by transplant) and/or one CPT Code 33945 (cardiectomy with heart transplant)
  - One or more mention of an immunosuppressant
  - Age  $> 15$  at date of transplant
  - Available DNA



# CHRONIC KIDNEY DISEASE

Chronic kidney disease is classified in 5 stages of severity (determined by estimated GFR)

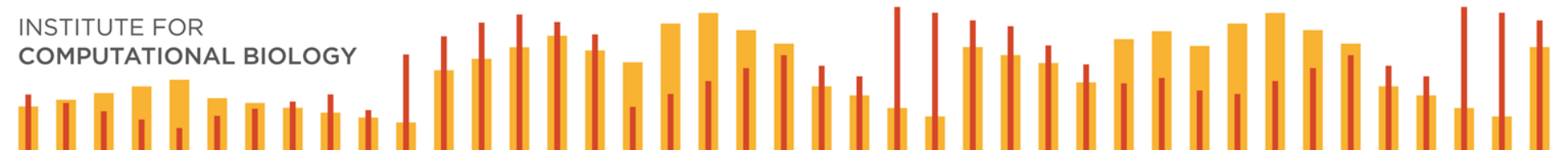


# OUTCOME FOR THIS STUDY

- eGFR is calculated

$$186 \times \text{Serum Creatinine}^{-1.154} \times \text{Age}^{-0.203} \times [1.212 \text{ if Black}] \times [0.742 \text{ if Female}]$$

- Severe Kidney Disease: Post-transplant eGFR < 30 mL/min/1.73m<sup>2</sup> for 3 consecutive months
- Time to development of severe nephrotoxicity clinically attributed to calcinurin inhibitor toxicity

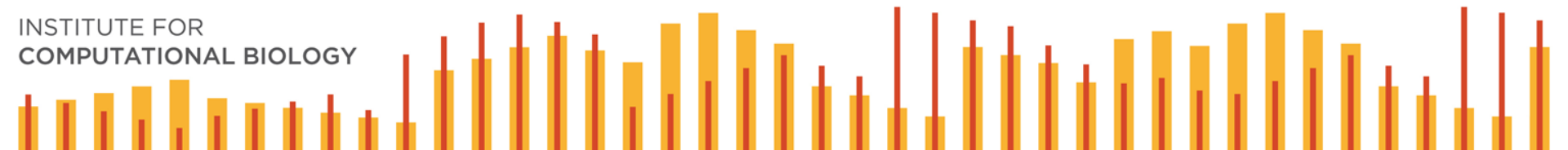
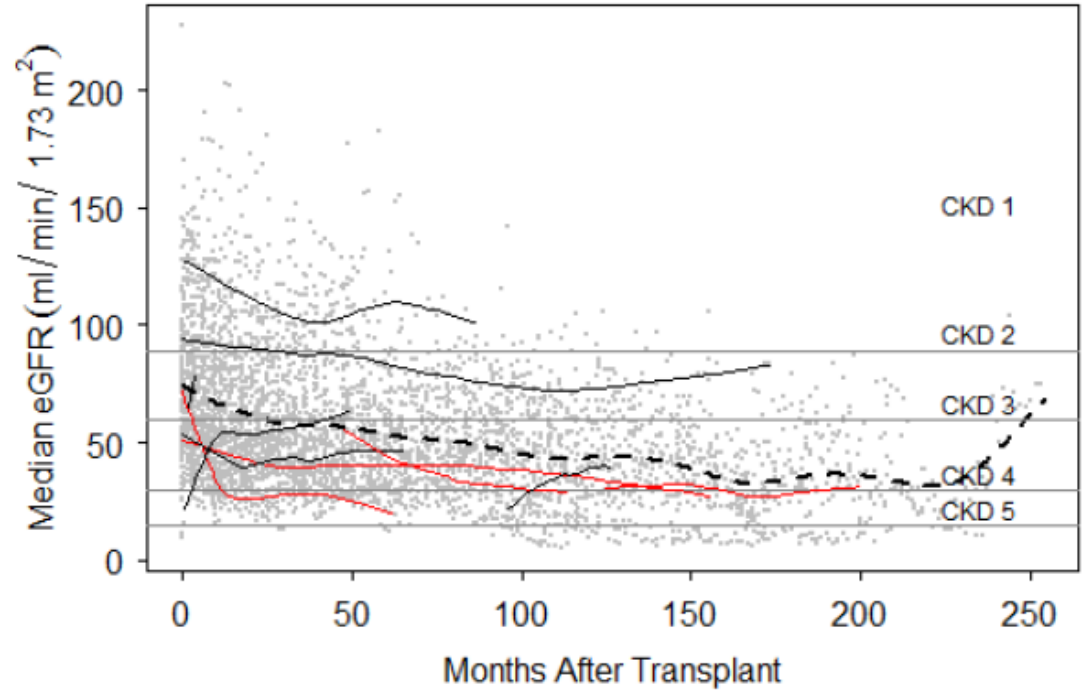


# CLINICAL TRAJECTORIES

A functional variant, rs1801265 in the gene *DPYD* associated to Cyclosporine induced nephrotoxicity (p = 0.03)

We can make new discoveries in EHR data, but it takes some work

Vanderbilt University Medical Center



# EXPLORING “PHENOMES”

Human Genome

Genome-Wide Association Study (GWAS)

## PHEWAS

- Can we explore the human genome in a new way?



Marylyn Ritchie  
Geisinger/PSU



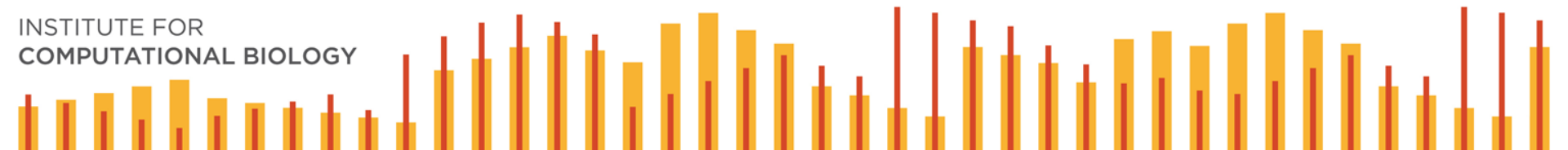
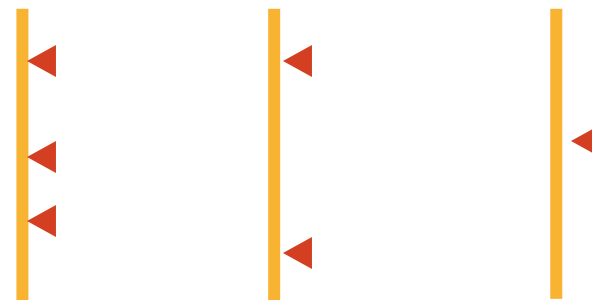
Sarah Pendergrass  
Geisinger



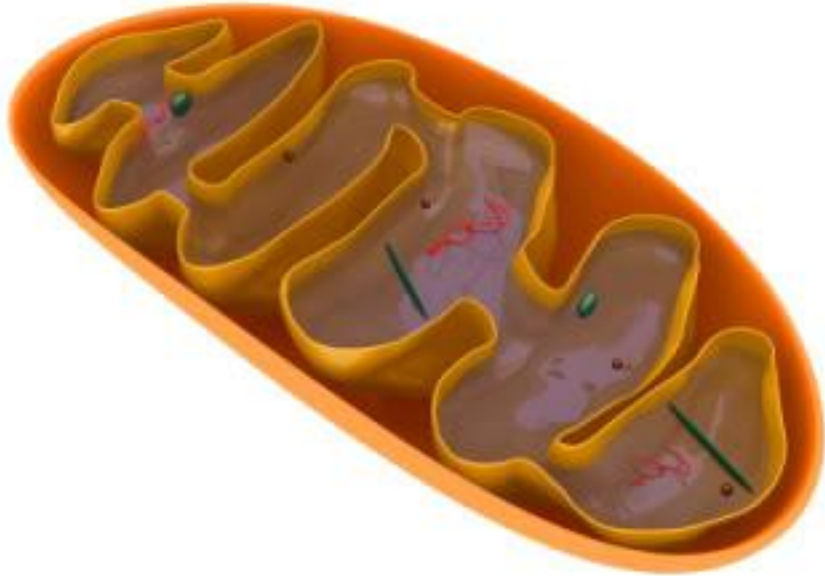
Dana Crawford  
CWRU

Human Phenome

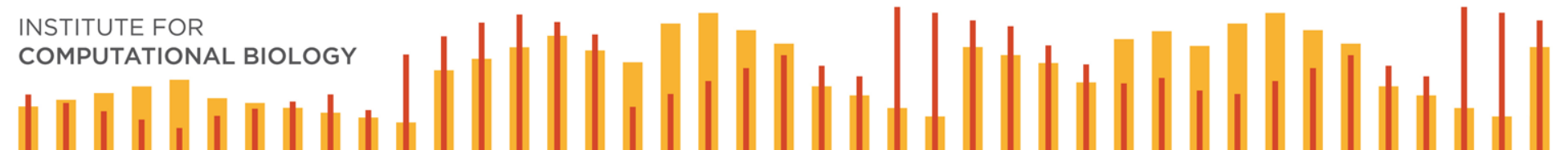
Diagnoses    Severity/Onset    Lab Measures



# THE MITOCHONDRIA



- *Power plants* of all living cells
- Have their own genome
- 26,000 base pairs versus 3 billion



# MITOCHONDRIAL PHEWAS

- Multiple quantitative traits for cardiovascular disease
- A series of mitochondrial genetic variants



Jake Hall



Sabrina Mitchell

Mitchell et al. *BioData Mining* 2014, 7:6  
<http://www.biodatamining.org/content/7/1/6>

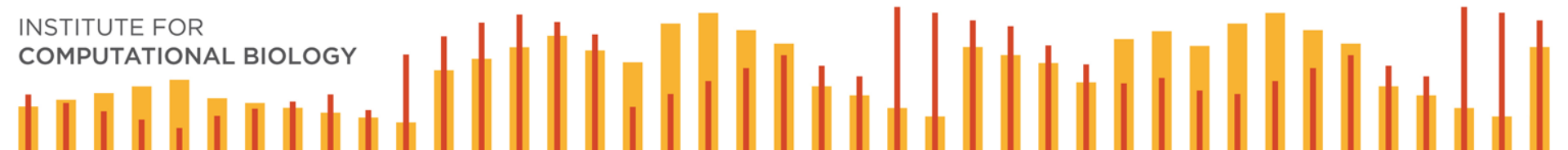


## RESEARCH

Open Access

Investigating the relationship between mitochondrial genetic variation and cardiovascular-related traits to develop a framework for mitochondrial phenome-wide association studies

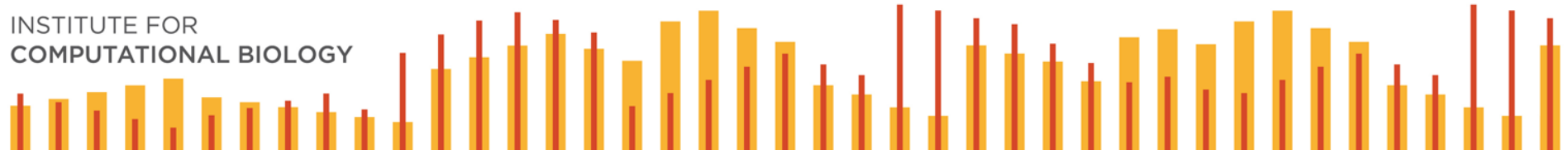
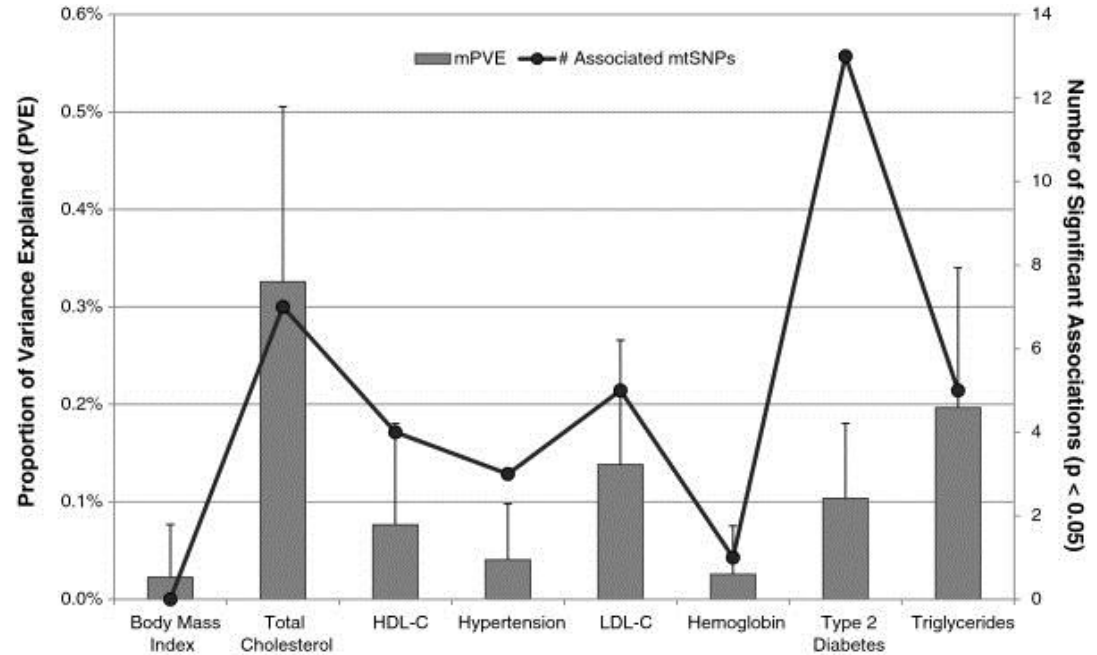
Sabrina L Mitchell<sup>1,2\*</sup>, Jacob B Hall<sup>1\*</sup>, Robert J Goodloe<sup>1</sup>, Jonathan Boston<sup>1</sup>, Eric Farber-Eger<sup>1</sup>, Sarah A Pendergrass<sup>3</sup>, William S Bush<sup>1,4</sup> and Dana C Crawford<sup>1,2\*</sup>





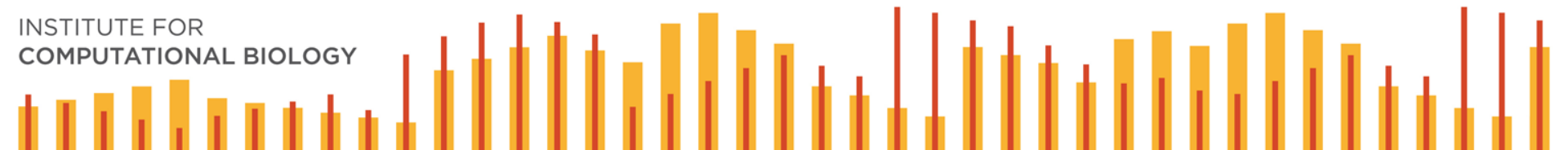
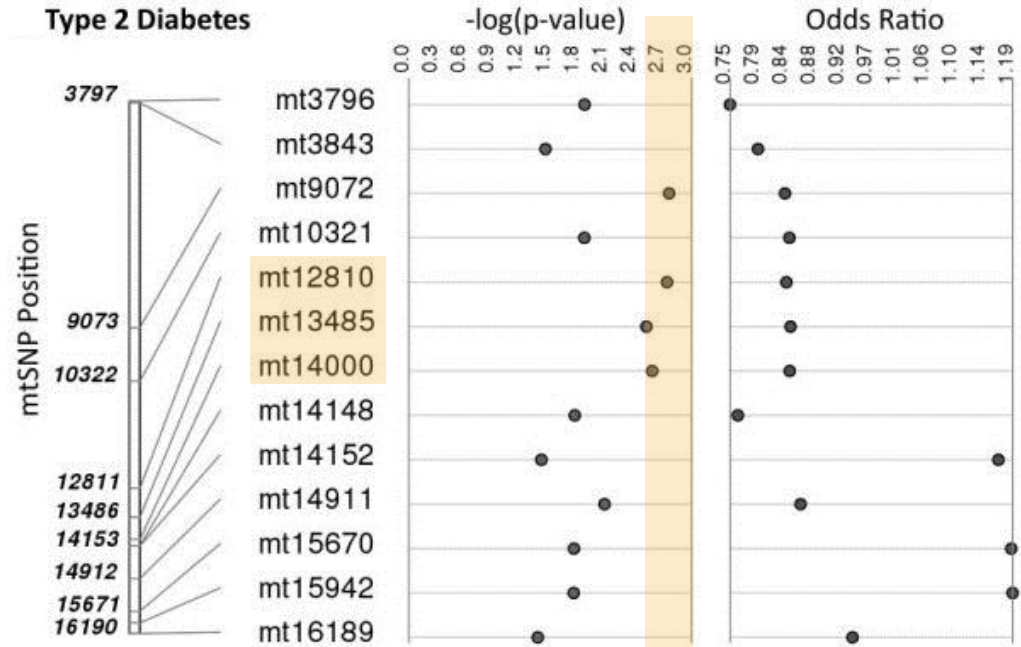
# MITOCHONDRIAL PHEWAS

- Total Cholesterol and Type II Diabetes show a significant amount of risk explained



# MITOCHONDRIAL PHEWAS

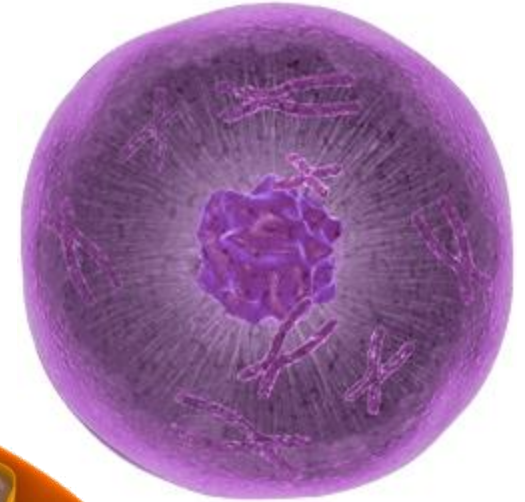
- A series of mitochondrial variants are protective for T2D risk



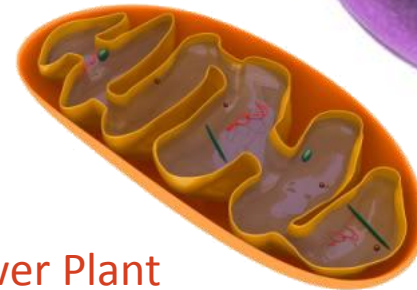
# MITOCHONDRIAL NUCLEAR INTERACTIONS

- There are complex chemical signals between the nucleus of the cell and mitochondria

Google Data Center

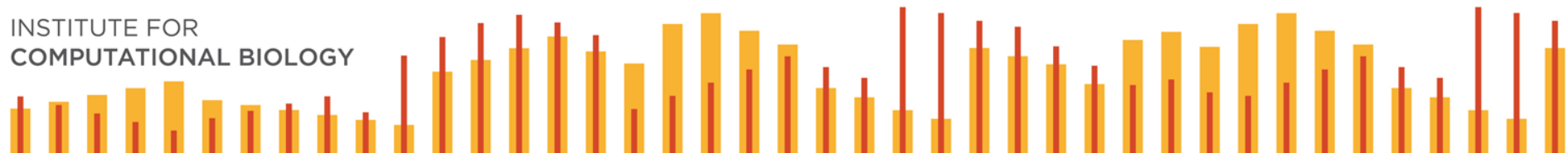


Mitochondria



Nucleus

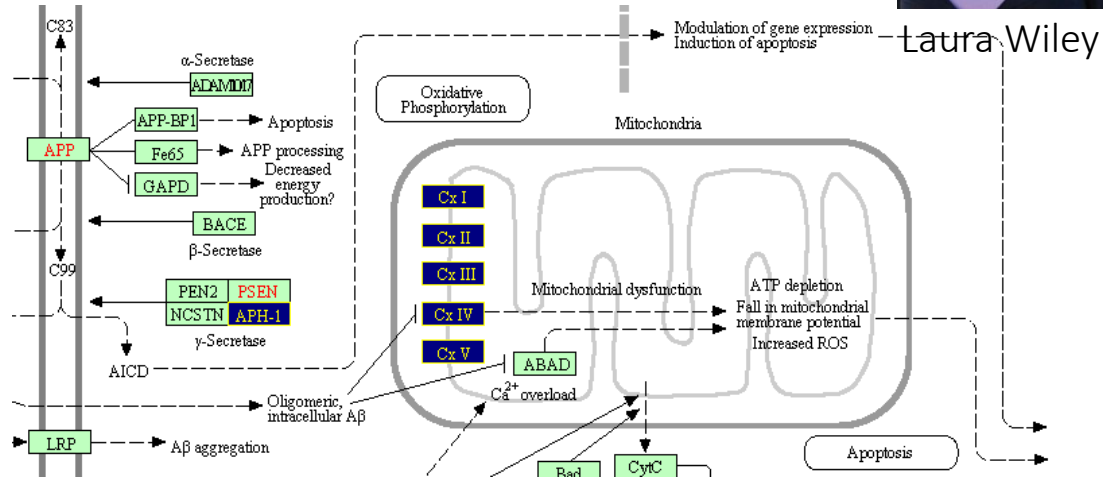
Power Plant



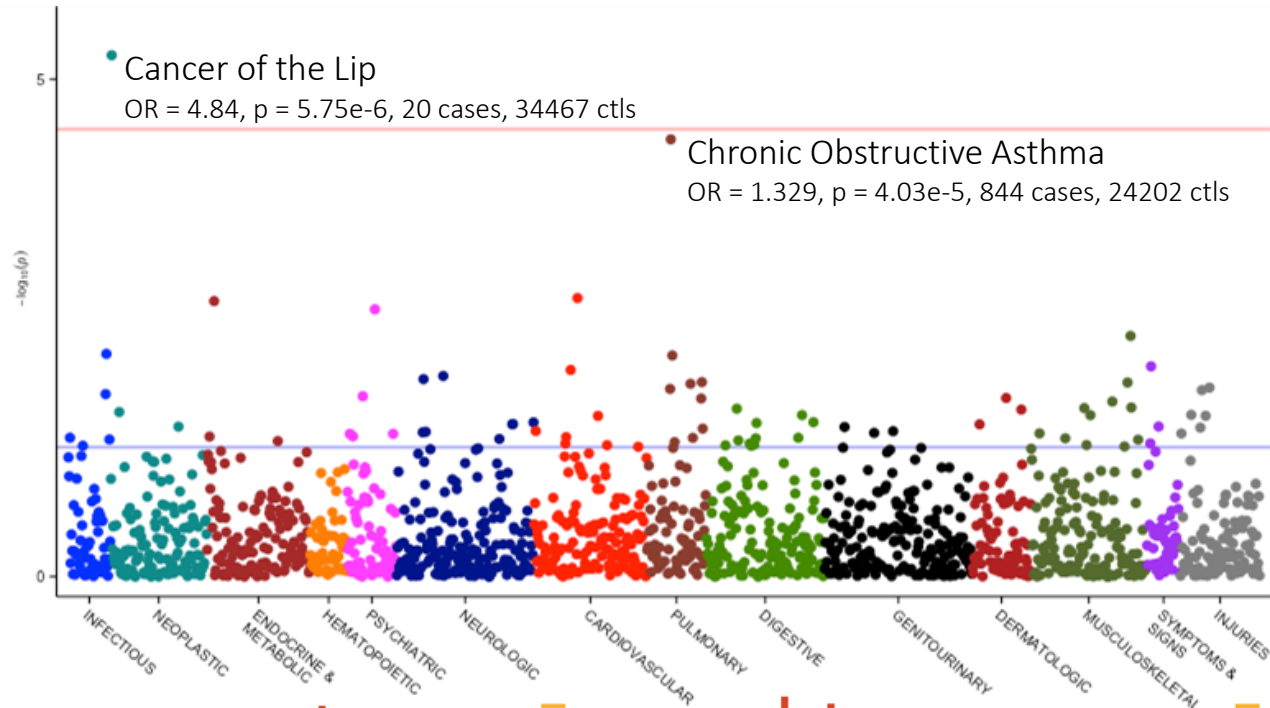
# MITOCHONDRIAL NUCLEAR INTERACTIONS



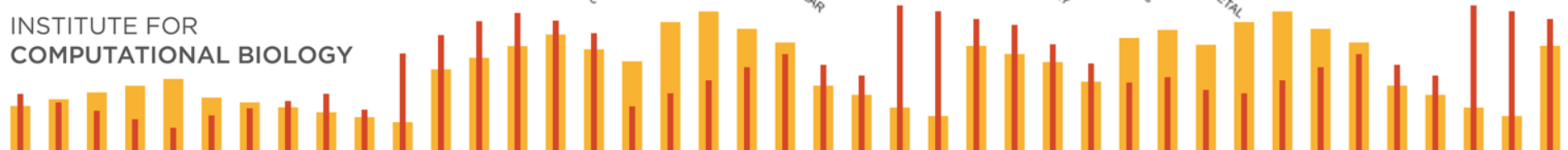
- We discovered a nuclear genetic change that influences expression of mitochondrial genes



# FUNCTIONAL SNP PHEWAS



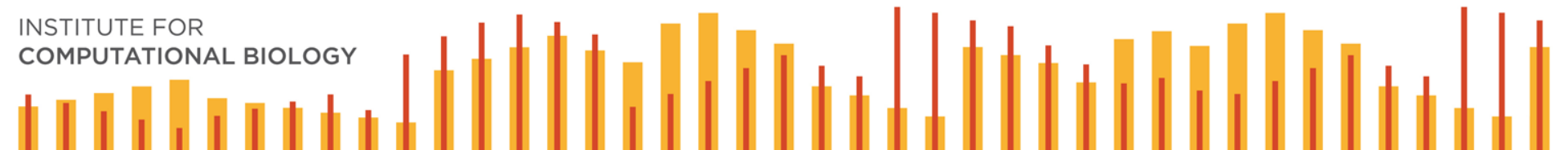
INSTITUTE FOR  
COMPUTATIONAL BIOLOGY



# LARGE SCALE GENOMIC ANALYSIS

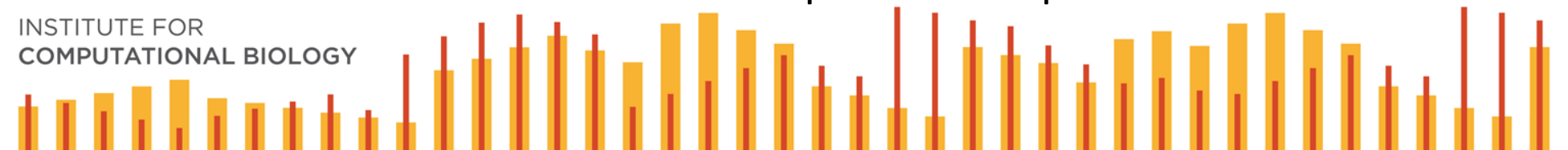


Jake Hall



# PATHWAY-BASED ANALYSES

- Individual genetic changes contribute very little to common diseases
- Examining changes in aggregate can inform biology as well (by mechanism)
- Combinations of genetic changes may contribute more risk than their independent parts



# MIXED-MODEL ANALYSES

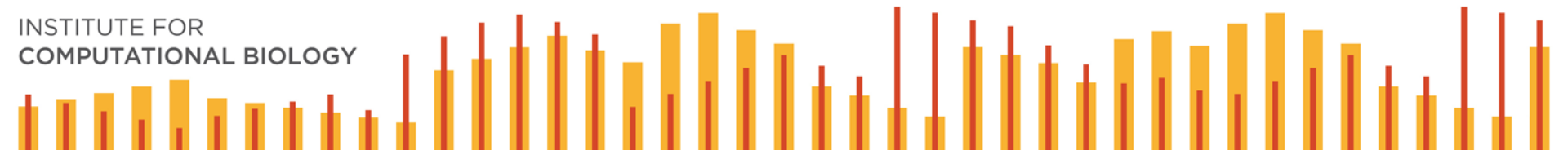
- Estimation of a genetic relationship matrix
- Produces an estimate of the risk explained by genetic sharing
- Assumes each genetic change contributes independently

$$\left[ A_{jk} = \frac{1}{N} \sum_{i=1}^N \frac{(\text{covar}.j)(\text{covar}.k)}{(\text{var})} \right]$$

46% of Type II Diabetes is explained



Nathan Morris

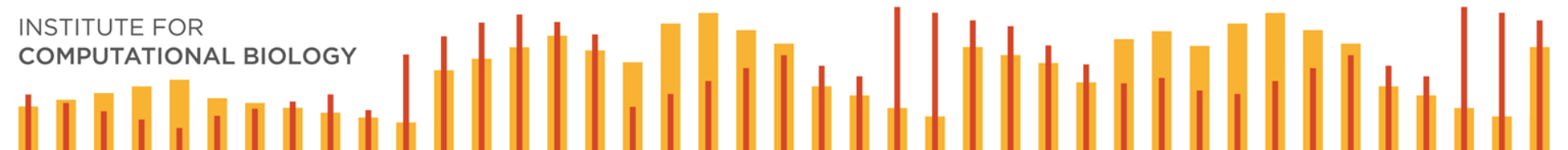




# EXPANDED GENETIC RELATIONSHIP MATRIX

- Can we model the effects of *interactions* among variants?

58% of Type II Diabetes is explained  
46% due to independent effects  
12% due to interaction effects



# EXPANDED GENETIC RELATIONSHIP MATRIX

## VARIANCE

Component	Equation
-----------	----------

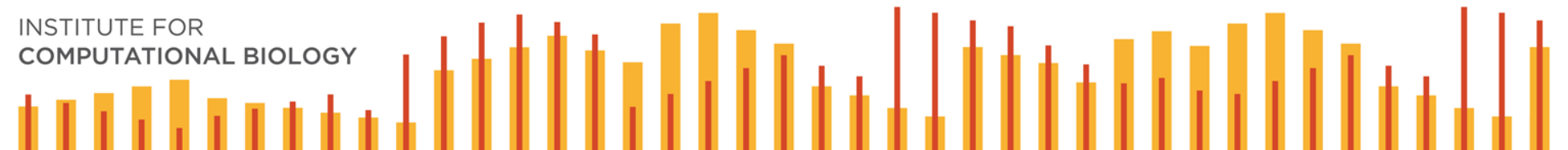
## COVARIANCE

Component	AABB	AABb	AAbb	AaBB	AaBb	Aabb	aaBB	aaBb	aabb
-----------	------	------	------	------	------	------	------	------	------

# Lots of Computation

A × D Interaction	$\frac{2xy}{(xy)^2}$
D × A Interaction	$\frac{2xy}{(uv)^2}$
D × D Interaction	$\frac{1}{(uvxy)^2}$

Interaction	$f_{\cdot 2}$	$f_{\cdot 1}$	$f_{\cdot 0}$	$f_{\cdot 2}$	$f_{\cdot 1}$	$f_{\cdot 0}$	$f_{\cdot 2}$	$f_{\cdot 1}$	$f_{\cdot 0}$
D × A Interaction	$\frac{2y}{f_{\cdot 2}}$	$\frac{(y-x)}{f_{\cdot 2}}$	$\frac{2x}{f_{\cdot 2}}$	$\frac{-4y}{f_{\cdot 1}}$	$\frac{-2(y-x)1}{f_{\cdot 1}}$	$\frac{4x}{f_{\cdot 1}}$	$\frac{2y}{f_{\cdot 0}}$	$\frac{(y-x)}{f_{\cdot 0}}$	$\frac{-2x}{f_{\cdot 0}}$
D × D Interaction	$\frac{1}{f_{22}}$	$\frac{-2}{f_{21}}$	$\frac{1}{f_{20}}$	$\frac{-2}{f_{12}}$	$\frac{4}{f_{11}}$	$\frac{-2}{f_{10}}$	$\frac{1}{f_{02}}$	$\frac{-2}{f_{01}}$	$\frac{1}{f_{00}}$



# ACKNOWLEDGMENTS

## The Bush Lab

- Jake Hall
- Mike Sivley
- Alexandra Fish
- Jeremy  
Fondran

## The Institute for Computational Biology

- Jonathan Haines
- Dana Crawford
- Chun Li
- Jill Barnholtz-Sloan
- Ricky Chan

## Grant Support

NIMH 095621

Asha Kallianpur/Todd Hulgan

NHGRI 004798

Dana Crawford

NIA 07133

Farrer, Schellenberg, Haines

