# Examples of how ENCODE facilitates biomedical research

ENCODE Users Meeting
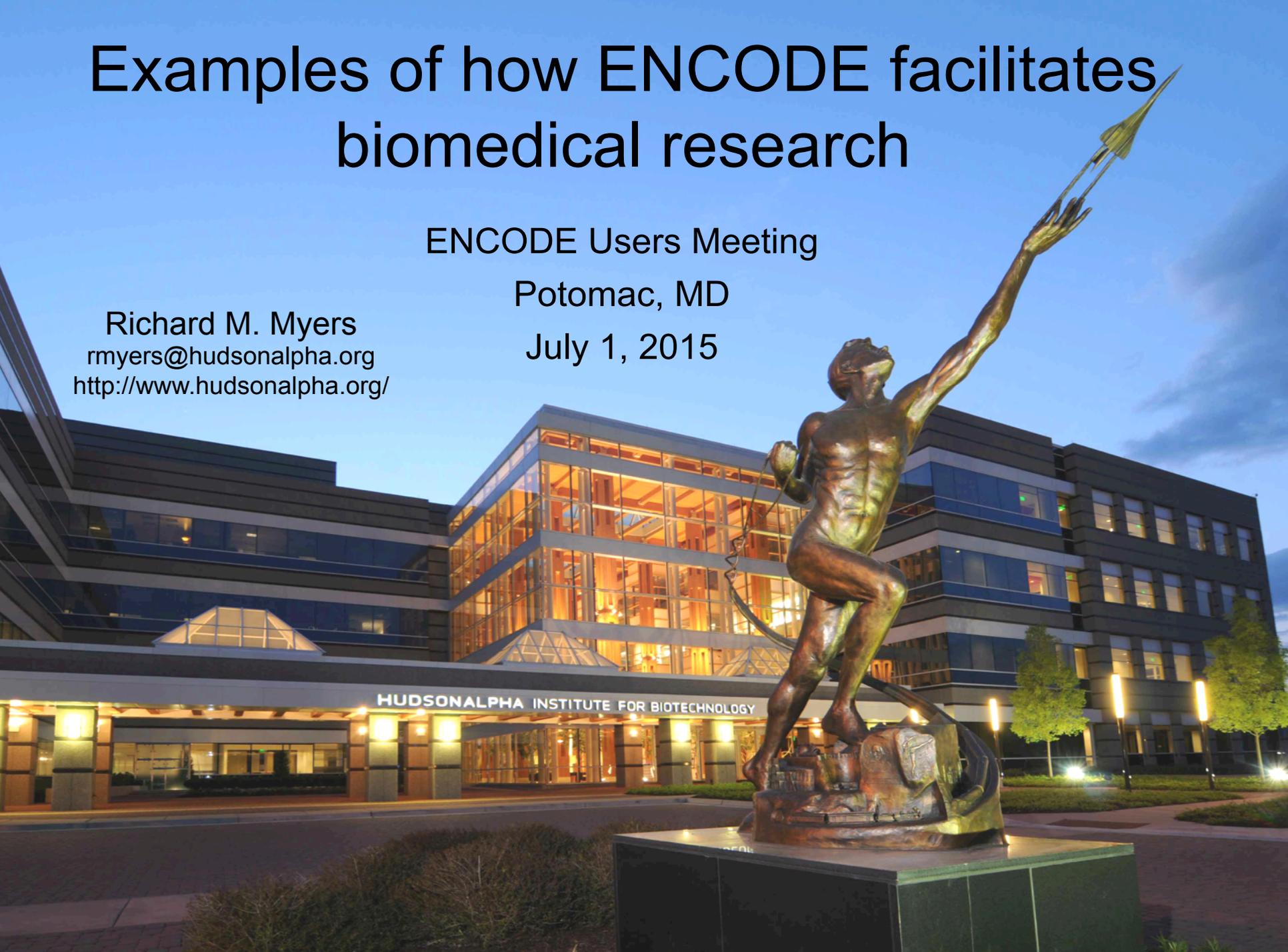
Potomac, MD
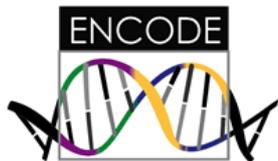
July 1, 2015

Richard M. Myers
rmyers@hudsonalpha.org
http://www.hudsonalpha.org/

# Disclosure

Our group has been part of the ENCODE Consortium since it began in 2003

**Rick Myers**

**Barbara Wold**

**Ross Hardison**

**Eric Mendenhall**

**Ali Mortazavi**

**Tim Reddy**

# Goals of ENCODE

Annotate the human genome

Disseminate data to researchers everywhere

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

# 5 examples of how we use ENCODE data to help in our research on human diseases

# 1.  Discovering the causes of undiagnosed genetic diseases

Richard M. Myers

# Childhood genetic disorders

1.5-3% of kids worldwide are born with 1 or more of:

- intellectual disability

- developmental delay

- heart defects

- craniofacial and skeletal abnormalities

- severe autism

- seizures

The vast majority of these problems have genetic causes

Richard M. Myers

# Diagnostic challenges for childhood genetic disorders

Inaccurate or undetermined causes (i.e., diagnoses) are a major hardship:

Years of expensive, invasive, and futile testing

Impossible to predict disease progression, symptoms

Treatment decisions are complicated

Slows research into developing new therapies

Impacts family planning

Results in feelings of parental guilt and lack of control

Thus, identifying the root genetic causes is essential

HUDSONALPHA
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

# HudsonAlpha Pediatric Genetics Project

Sequence whole genomes of 500 children with developmental/ intellectual delay of unknown etiology (and both parents' genomes too)

North Alabama
Children's Specialists

Univ. of Louisville
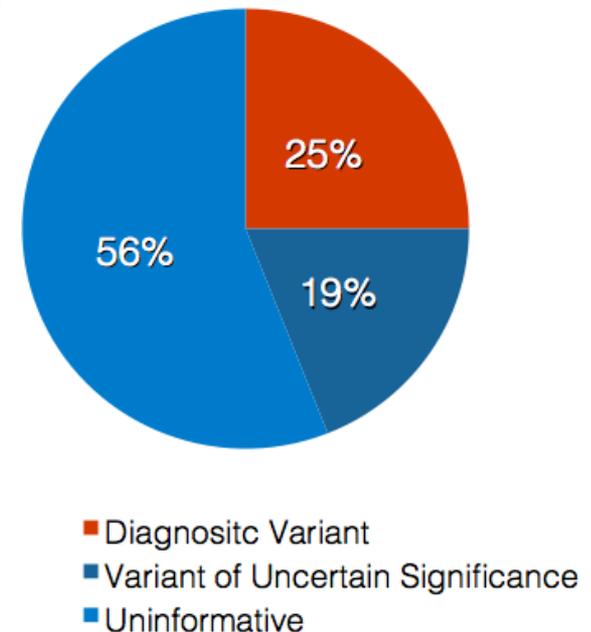
HudsonAlpha

Richard M. Myers

# Exome results so far

Exome sequencing completed for 171 families

Definitive genetic diagnosis in 25% of the children

- Pitt-Hopkins syndrome
- Dravet syndrome
- Rett syndrome
- Rubinstein-Taybi syndrome
- Noonan-like syndrome
- Many never-described causes



- 25%
- 56%
- 19%

■ Diagnositc Variant
■ Variant of Uncertain Significance
■ Uninformative

>20% of families receive uncertain genetic findings that will likely be definitively diagnostic in the future

Richard M. Myers

# Whole genome sequencing of trios

Illumina X Ten sequencers:  $ of 30X WGS = $ of exome

We have completed WGS of 30 trios in our Childhood Genetics Project

## Results:

Diagnostic rate is higher

Identified at least 3 cases where regulatory mutations were the causes

We relied heavily on ENCODE data to identify functional regulatory segments

**HudsonAlpha**
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

# Annotating genetic variants

## Problem:

HUGE number of sequence variants in each individual

Most are not important

How to find which variants have an effect on:
> The molecular/biochemical function of the gene
> The organism

Richard M. Myers

# CADD
## Combined Annotation Dependent Depletion
### Greg Cooper and Jay Shendure

## A general framework for estimating the relative pathogenicity of human genetic variants

Martin Kircher[1,5], Daniela M Witten[2,5], Preti Jain[3,4], Brian J O'Roak[1,4], Gregory M Cooper[3] & Jay Shendure[1]

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

# CADD integrates many features to give a single pathogenicity score



Legend:
- C score (91.64%)
- GerpS (84.79%)
- PhCons (83.09%)
- phyloP (83.03%)

Richard M. Myers

# Typical vs pathogenic CADD scores



Legend:
- HapMap Exome
- ClinVar Pathogenic

Promoter mutations that cause B-thalassemia
Enhancer mutations that cause pancreatic agenesis
Enhancer mutations that cause limb defects

CADD Score

Richard M. Myers

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

# Use the CADD webserver!

## http://cadd.gs.washington.edu

Richard M. Myers

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

# 2. Understanding renal cell carcinoma

ENCODE data were instrumental in helping us identify regions of the genome that are ~100% accurate diagnostic markers for kidney cancer

(and even for prognosis of different subtypes)

Richard M. Myers

# Genomic signatures of renal cell carcinoma

Brittany Lasseigne, Jim Brooks, Myers Lab

BMC Medicine

**RESEARCH ARTICLE**　　　　**Open Access**

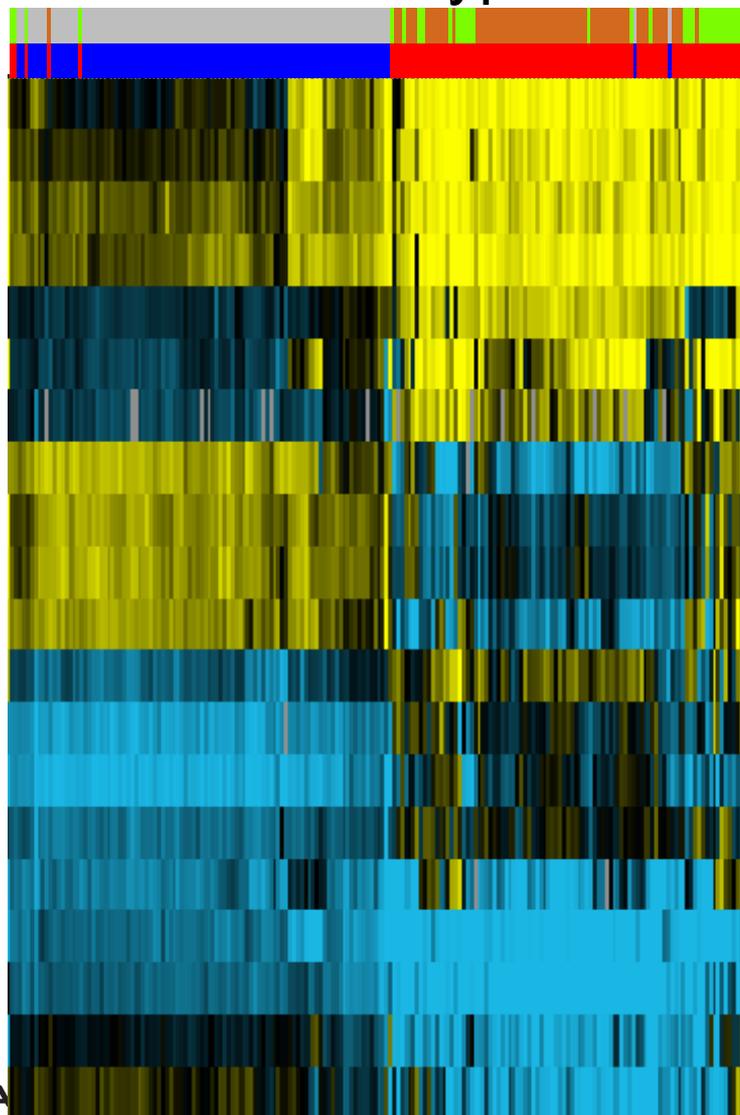## DNA methylation profiling reveals novel diagnostic biomarkers in renal cell carcinoma

Brittany N Lasseigne[1,2], Todd C Burwell[1], Mohini A Patil[3], Devin M Absher[1], James D Brooks[3] and Richard M Myers[1*]

We measured DNA methylation and copy number variants in 135 kidney tumors and matched non-tumor kidney tissues
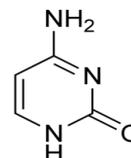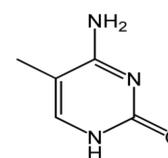
# Top 20 DNA methylation markers



All subtypes

20 CpGs

Kidney Tumor
Normal Tissue

Clear Cell
Other Subtypes
Normal Tissue

0%    50%    100%

Cytosine    5-Methyl Cytosine
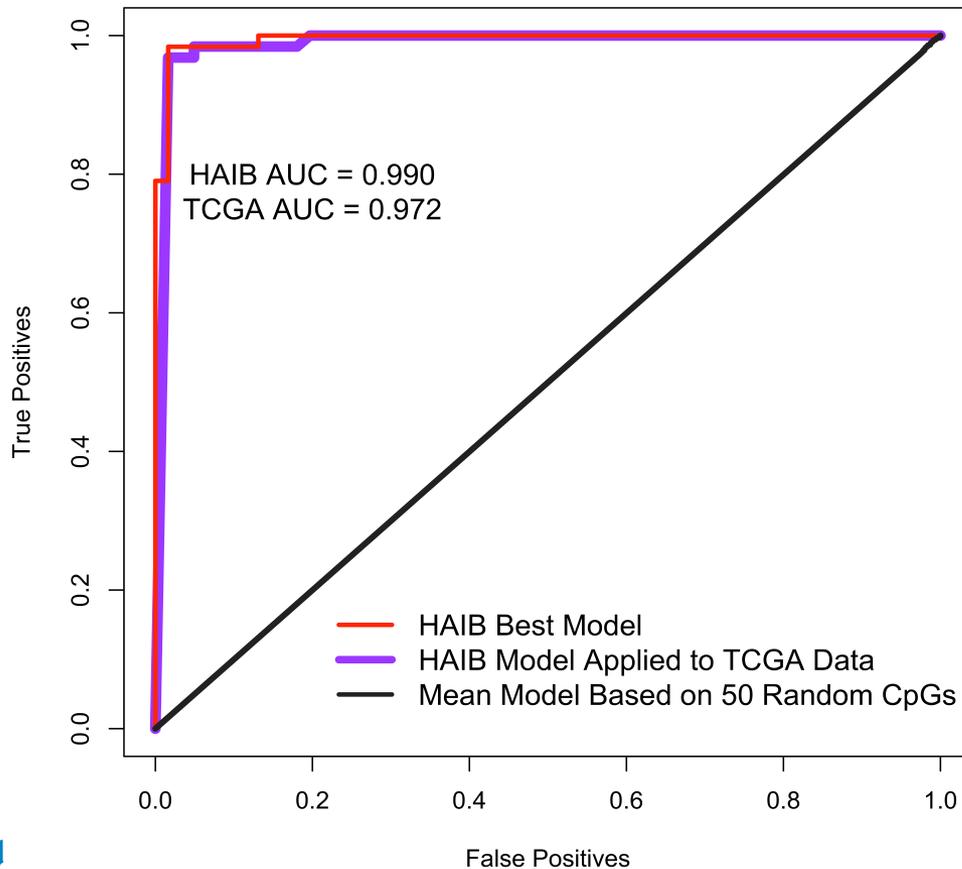
Richard M. Myers

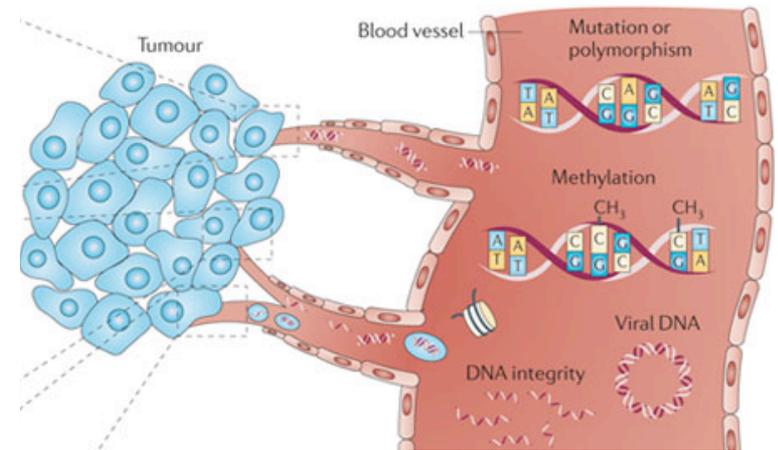HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

18

# DNA methylation patterns are highly accurate at predicting patients with renal cell carcinoma

ROC curves of DNA methylation results from 135 tumor and matched non-tumor samples from RCC patients



HAIB AUC = 0.990
TCGA AUC = 0.972

Legend:
- HAIB Best Model
- HAIB Model Applied to TCGA Data
- Mean Model Based on 50 Random CpGs

X-axis: False Positives
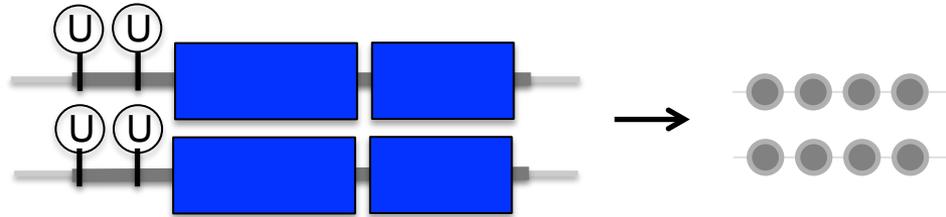Y-axis: True Positives

Apply these assays to urine or blood as a routine screening for early detection of kidney cancer



Schwarzenbach et al. Nature Reviews Cancer 11, 426-437
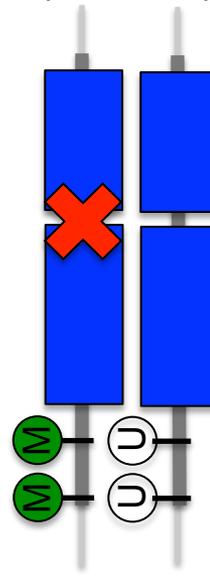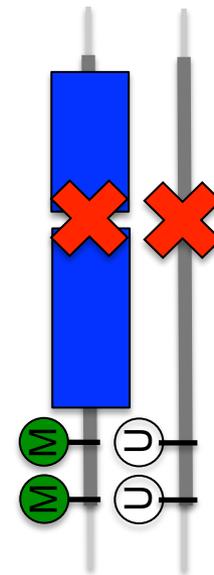
Richard M. Myers

# Integrating genomic signatures



HUDSONALPHA
INSTITUTE FOR BIOTECHNOLOGY
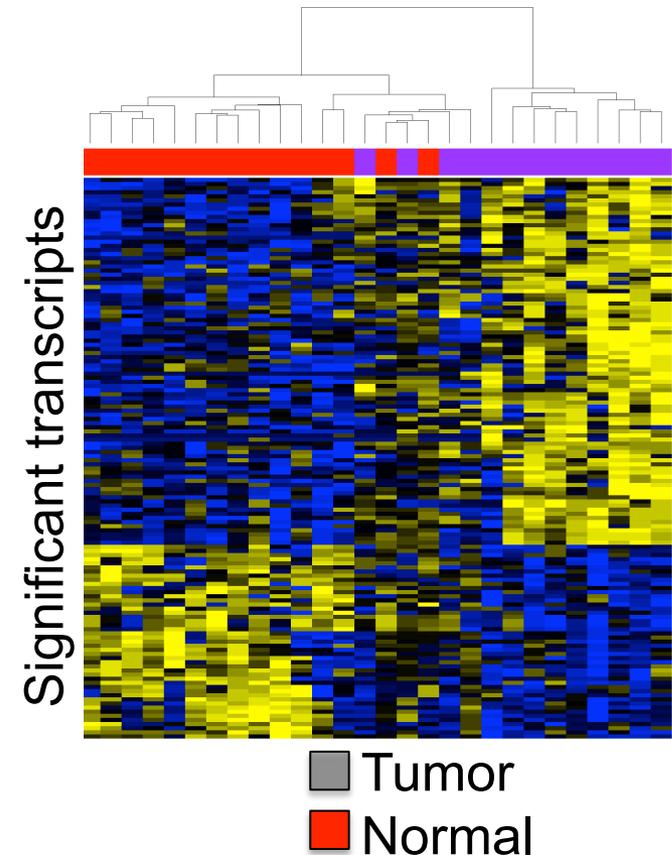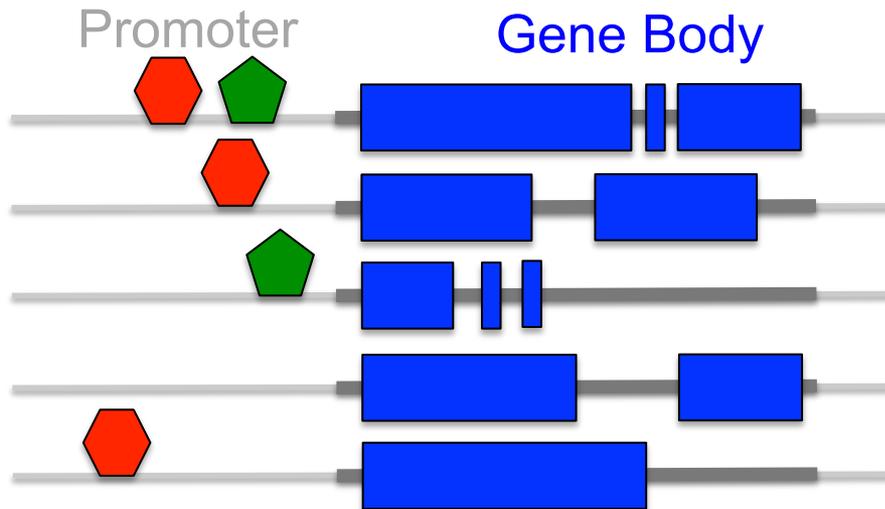
# Example: MSH4 gene

# 3. Using ENCODE TF data to prioritize cancer genetics and functional genomics data

Richard M. Myers

# Using ENCODE data to find cancer regulators

Genomic assays often reveal thousands of dysregulation events in cancer

These widespread genomic changes may be regulated by a few key transcription factors

# Differentially expressed genes in cancer are enriched for particular TFs

Genes differentially expressed in prostate cancer compared to normal prostate tissue are enriched for **EZH2**, **SUZ12**, and **CTPB2** binding sites (adjusted p-value < 0.05) and actual binding events (from ENCODE ChIP data)

# Intersect transcription factor binding sites from the ENCODE Project with genomic regions specifically unmethylated in basal breast cancer

25

# Master regulators (?) of different breast cancer subtypes

Intersect gene regulatory regions containing subtype-associated methylation with binding sites of 149 transcription factors in ENCODE datasets

Significantly enriched binding sites:

| Transcription Factor | Fold Enrichment |
|---|---|
| Estrogen Receptor | 6.9 |
| FOXA1 | 8.1 |
| GATA3 | 10.3 |

Richard M. Myers

# Master regulators (?) of different breast cancer subtypes

Intersect gene regulatory regions containing subtype-associated methylation with binding sites of 149 transcription factors in ENCODE datasets
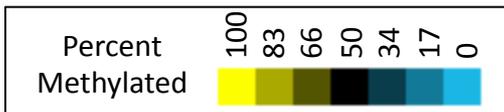
Significantly enriched binding sites:

| Transcription Factor | Fold Enrichment |
|---|---|
| Estrogen Receptor | 6.9 |
| FOXA1 | 8.1 |
| GATA3 | 10.3 |

| Transcription Factor | Fold Enrichment |
|---|---|
| STAT3 | 4.8 |
| GR (glucocorticoid receptor) | 4.2 |



Percent Methylated: 100 83 66 50 34 17 0

Richard M. Myers

# 4. Using RNA-seq to identify drug targets

Richard M. Myers

# Transcript fusions in cancer



Gene A    Gene B

PRECLINICAL STUDY

## Recurrent read-through fusion transcripts in breast cancer

Katherine E. Varley · Jason Gertz · Brian S. Roberts · Nicholas S. Davis ·
Kevin M. Bowling · Marie K. Kirby · Amy S. Nesmith · Patsy G. Oliver ·
William E. Grizzle · Andres Forero · Donald J. Buchsbaum · Albert F. LoBuglio ·
Richard M. Myers

K-T Varley with collaborators at UAB Comprehensive Cancer Center

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

# 3 fusion transcripts produce fusion proteins located in the cell membrane



IL17RC-CRELD1

SUM-149  MCF-7  ZR-75-1
+        +      −

220 / 120 kDa — Predicted Fusion Size = 102 kDa

50 / 40 kDa — CRELD1 Canonical Size = 46 kDa

SCNN1A-TNFRSF1A

HCC-1954  SUM-102  BT474
+         +        −

220 / 120 kDa — Predicted Fusion Size= 112 kDa

60 / 50 kDa — TNFRSF1A Canonical Size = 51 kDa

CTSD-IFITM10

MCF-7  BT-20  2-LMP
+      +      −

80 / 60 kDa — Predicted Fusion Size = 60 kDa

50 kDa — CTSD Canonical Size = 45 kDa

Potential therapeutic: Use drug-antibody complexes to direct a cellular toxin exclusively to cancer cells

1 ADC in plasma

2 ADC binds to receptor

3 ADC-receptor complex is internalized

4 Cytotoxic agent is released

5 Apoptosis (cell death)

CANCER CELL

Richard M. Myers

Seattle Genetics

# 5.  Which TF binding events are functionally important?

Richard M. Myers

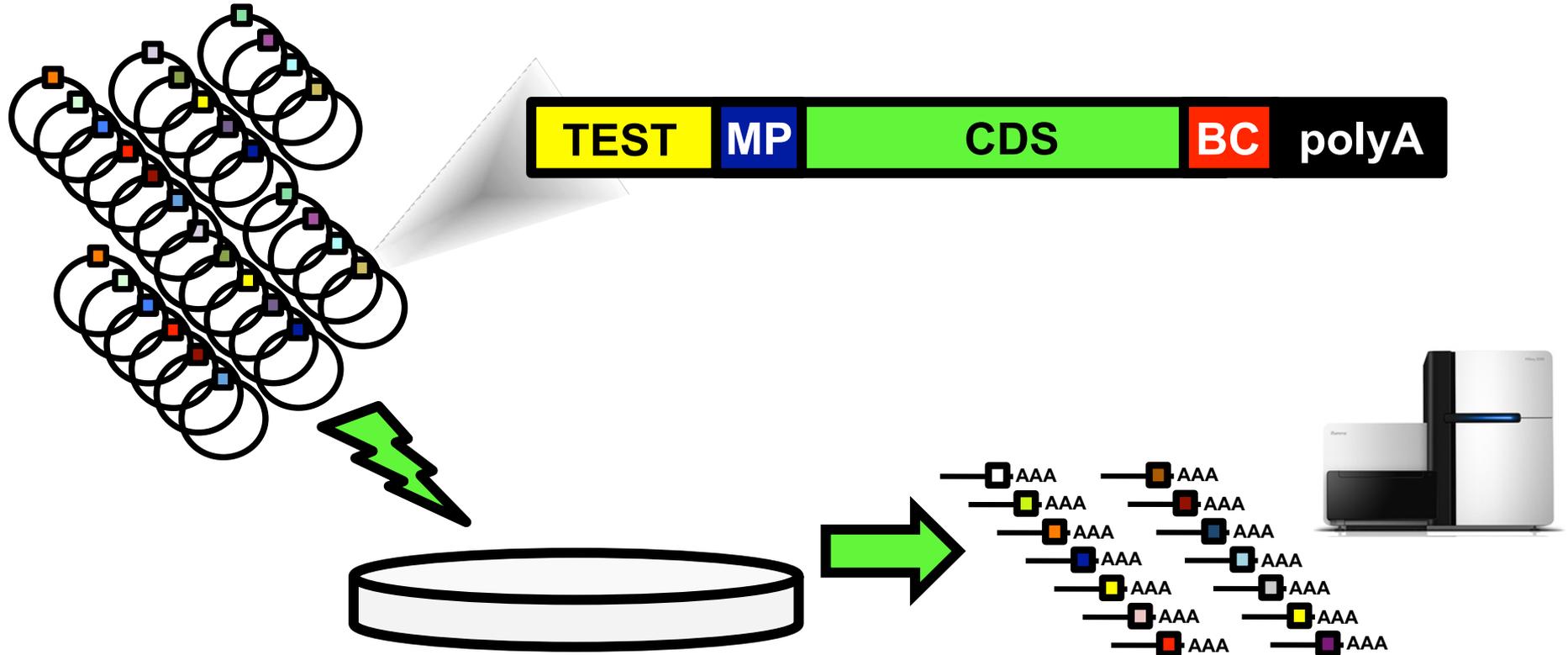# Using expression assays to identify functional transcription elements

## (especially long-distance ones)

Dan Savic, Brian Roberts, Chris Partridge, Barak Cohen, Greg CooperJay Gertz, Rick Myers

Test thousands of ENCODE-identified putative elements (based on TF binding, chromatin marks, etc.) in an ultra-high throughput reporter assay

**HUDSONALPHA**
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

# Massively parallel reporter assay
## Cis-Regulatory Element sequencing (CRE-seq)



Barcode abundance (sequence count)
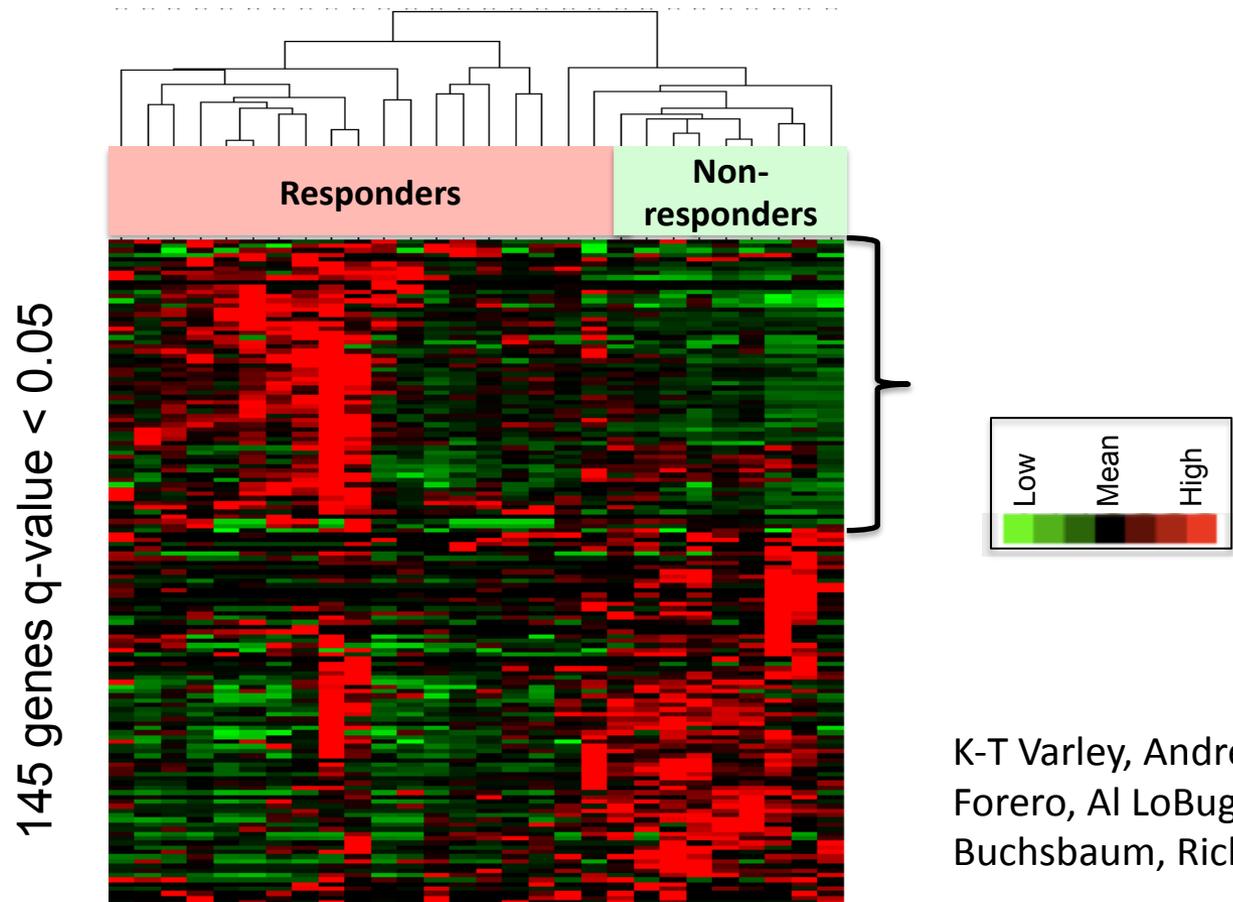is a proxy for test sequence activity

HUDSONALPHA
INSTITUTE FOR BIOTECHNOLOGY

Richard M. Myers

33

# Findings

RNAP2 at promoter-distal TF sites is a very strong mark of active regulatory elements

Richard M. Myers

Richard M. Myers

# 3. Using genomics to predict which patients will respond to various treatments

Richard M. Myers

HudsonAlpha
INSTITUTE FOR BIOTECHNOLOGY

# Clinical trial of a novel combination of drugs in ER+ breast cancer

Gene expression patterns in **responders** and **non-responders** during clinical trial of Letrozole (anti-estrogen) and Avastin (anti-angiogenesis)



K-T Varley, Andres Forero, Al LoBuglio, Don Buchsbaum, Rick Myers

Richard M. Myers