

# Integrating personalized gene expression profiles into predictive disease-associated gene pools

Emre Güney, PhD

*Institute for Research in Biomedicine (IRB)  
Barcelona, Spain*

ISMB/ECCB'17  
July 25<sup>th</sup>, 2017

ARTICLE

OPEN

# Integrating personalized gene expression profiles into predictive disease-associated gene pools

Jörg Menche<sup>1,2,3</sup>, Emre Guney<sup>1,4</sup>, Amitabh Sharma<sup>1,4,5</sup>, Patrick J. Branigan<sup>6</sup>, Matthew J. Loza<sup>6</sup>, Frédéric Baribaud<sup>6</sup>, Radu Dobrin<sup>7</sup> and Albert-László Barabási<sup>1,2,4,5</sup>

Gene expression data are routinely used to identify genes that *on average* exhibit different expression levels between a case and a control group. Yet, very few of such differentially expressed genes are detectably perturbed in individual patients. Here, we develop

ISMB/ECCB'17  
July 25<sup>th</sup>, 2017

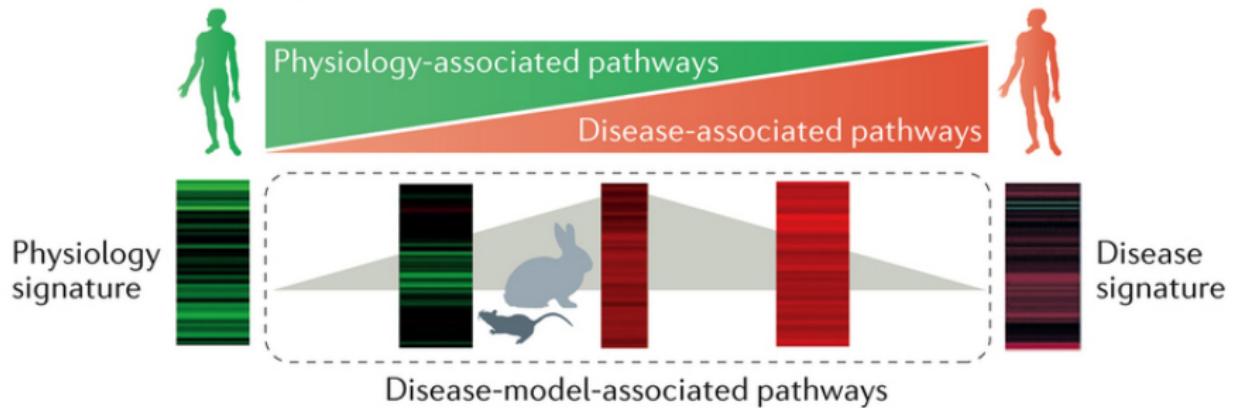


Image from *Moffat et al., 2017, Nat Rev Drug Discov*

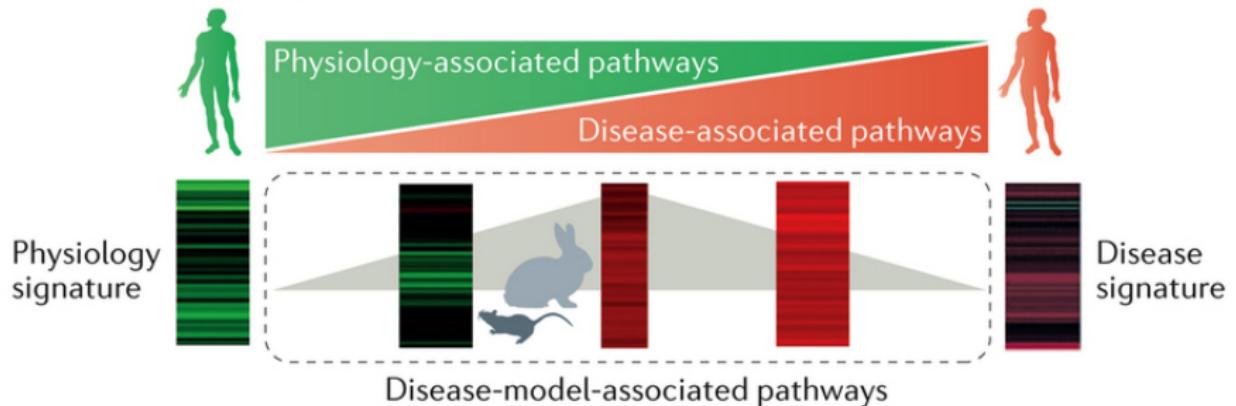


Image from *Moffat et al., 2017, Nat Rev Drug Discov*



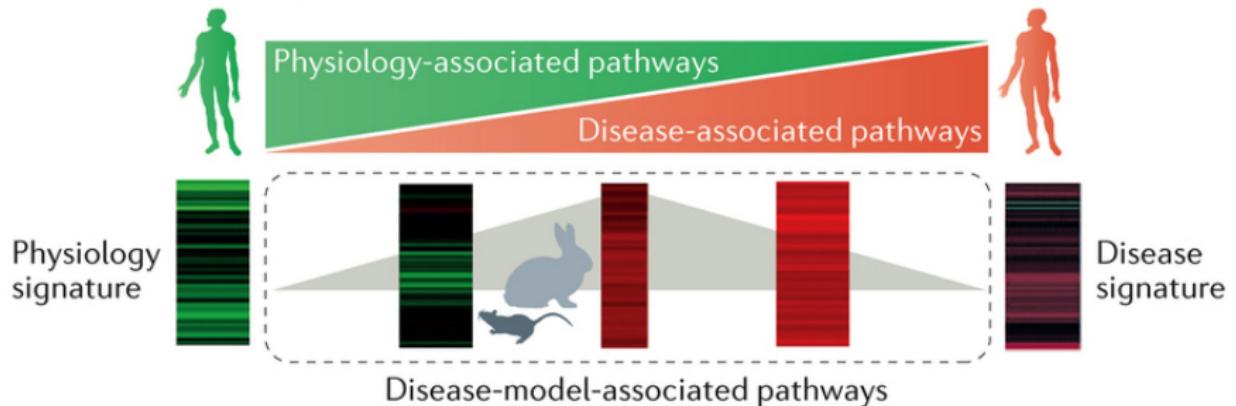


Image from *Moffat et al., 2017, Nat Rev Drug Discov*

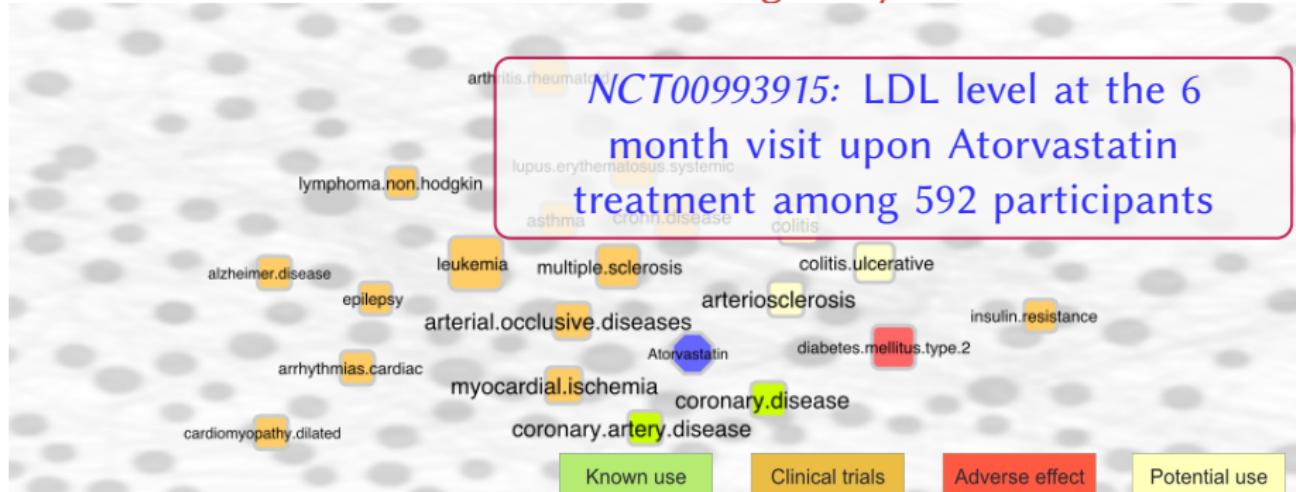


# Patient-level heterogeneity



Atorvastatin: one of the most successful cholesterol-lowering drug

# Patient-level heterogeneity



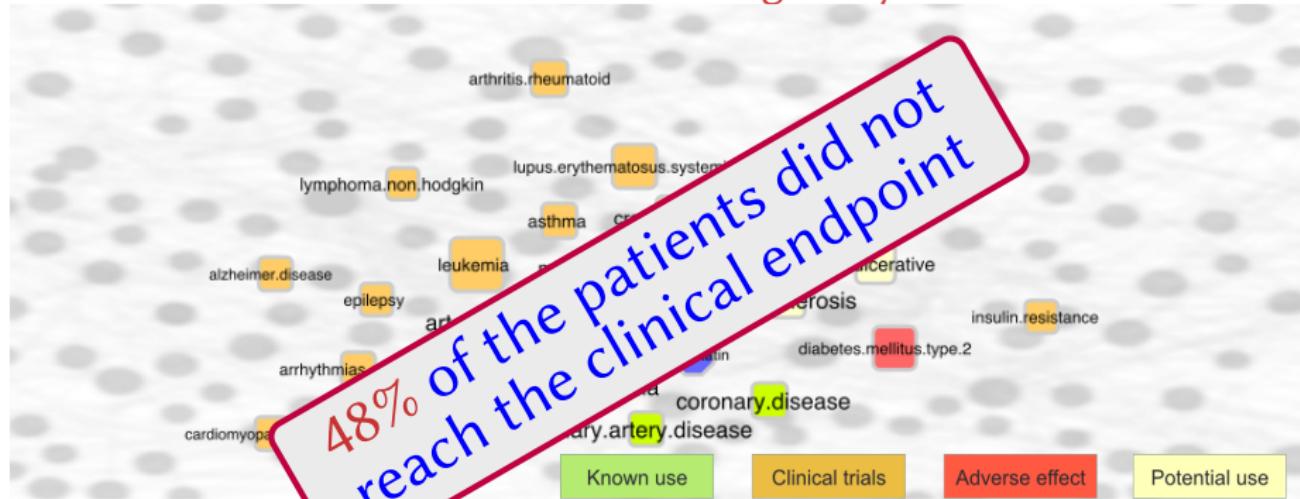
**Atorvastatin:** one of the most successful cholesterol-lowering drug

**Clinical endpoint:**  $\text{LDL} \leq 100 \text{ mg/dL}$  after treatment

[*Pre-treatment*] The average LDL level 179 mg/dL

[*Post-treatment*] 52% of the participants had  $\text{LDL} \leq 100 \text{ mg/dL}$

# Patient-level heterogeneity



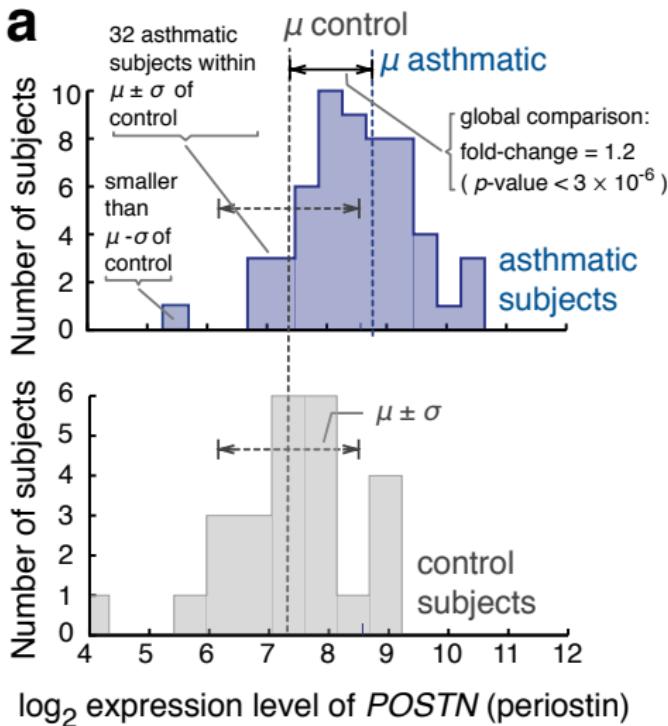
Atorvastatin: one of the most successful cholesterol-lowering drug

Clinical endpoint:  $LDL \leq 100 \text{ mg/dL}$  after treatment

[*Pre-treatment*] The average LDL level 179 mg/dL

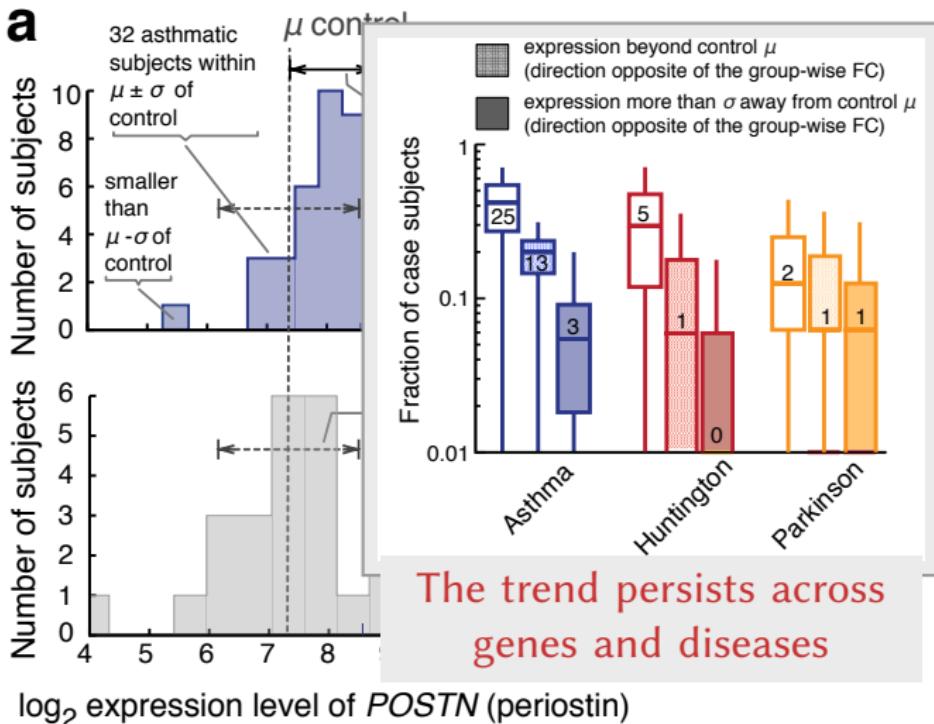
[*Post-treatment*] 52% of the participants had  $LDL \leq 100 \text{ mg/dL}$

## Gene expression of POSTN in asthmatic and control individuals



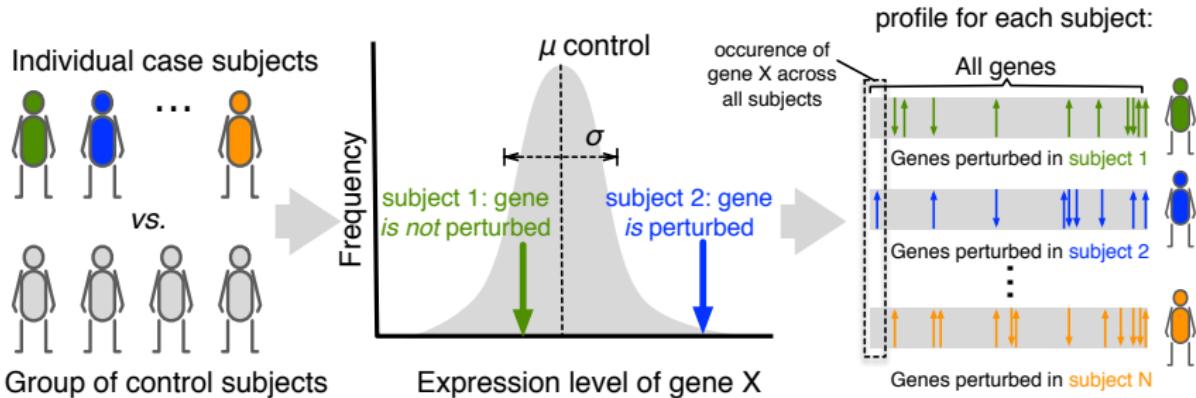
Periostin (*POSTN*) functions in the adhesion and migration of epithelial cells, whose up-regulation is a biomarker for Asthma.

# Gene expression of POSTN in asthmatic and control individuals



Periostin (*POSTN*) functions in the adhesion and migration of epithelial cells, whose up-regulation is a biomarker for Asthma.

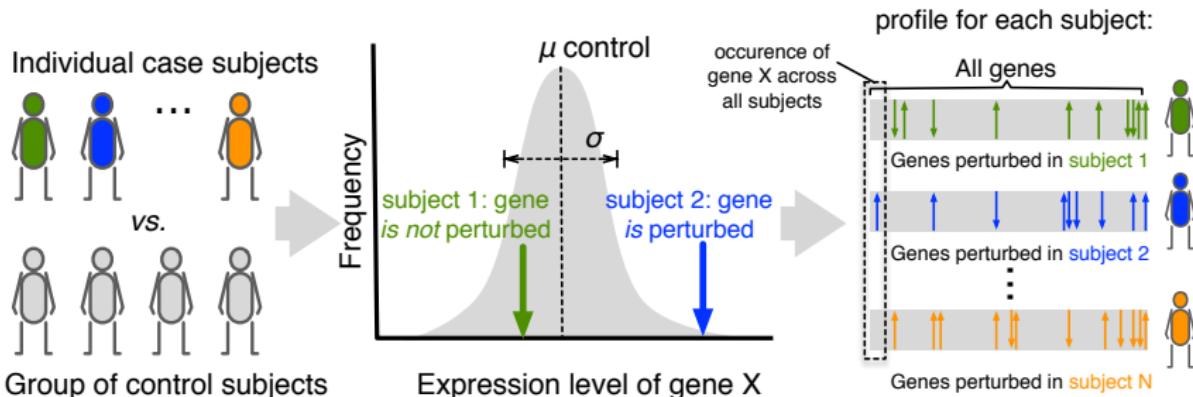
# PEEP: PErsonalized Expression Profile



	Samples			
Case	x	x	...	x
Control	c	c	c	c

$$z(\text{gene in } \textcolor{blue}{x}) = \frac{\text{expression}_{\textcolor{blue}{x}}(\text{gene}) - \mu_c(\text{gene})}{\sigma_c(\text{gene})}$$

# PEEP: PErsonalized Expression Profile



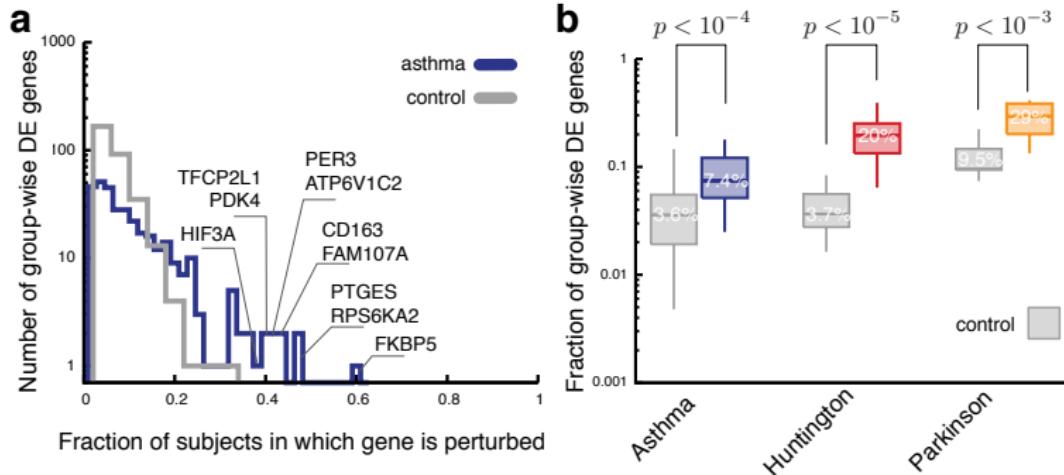
	Samples			
Case	x	x	...	x
Control	c	c	c	c

$$z(\text{gene in } \textcolor{blue}{x}) = \frac{\text{expression}_{\textcolor{blue}{x}}(\text{gene}) - \mu_c(\text{gene})}{\sigma_c(\text{gene})}$$

$$\text{PEEP}(\textcolor{blue}{x}) : \forall \text{gene } |z(\text{gene in } \textcolor{blue}{x})| > z_{threshold}$$

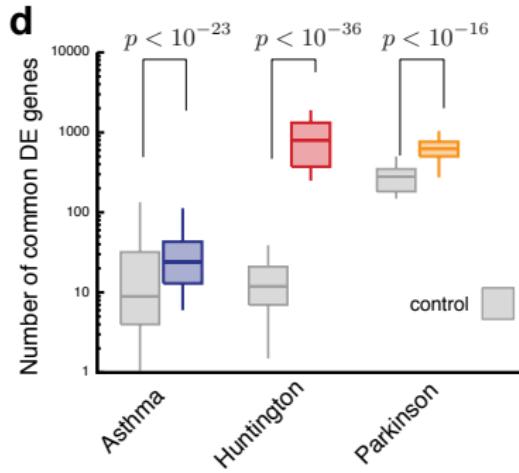
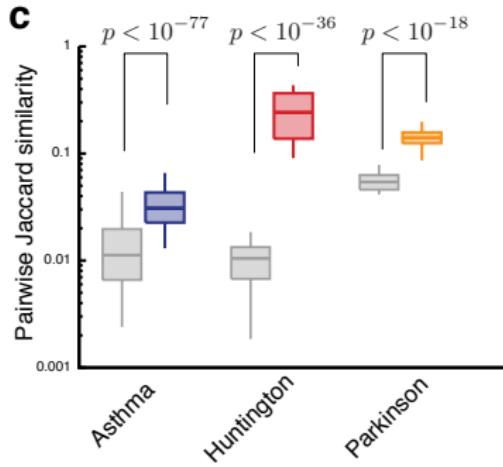
[ Genes that are significantly perturbed in each individual ]

## Quantifying the heterogeneity using PEEPs



Only a small fraction of group-wise DE genes appears in PEEP of each patient

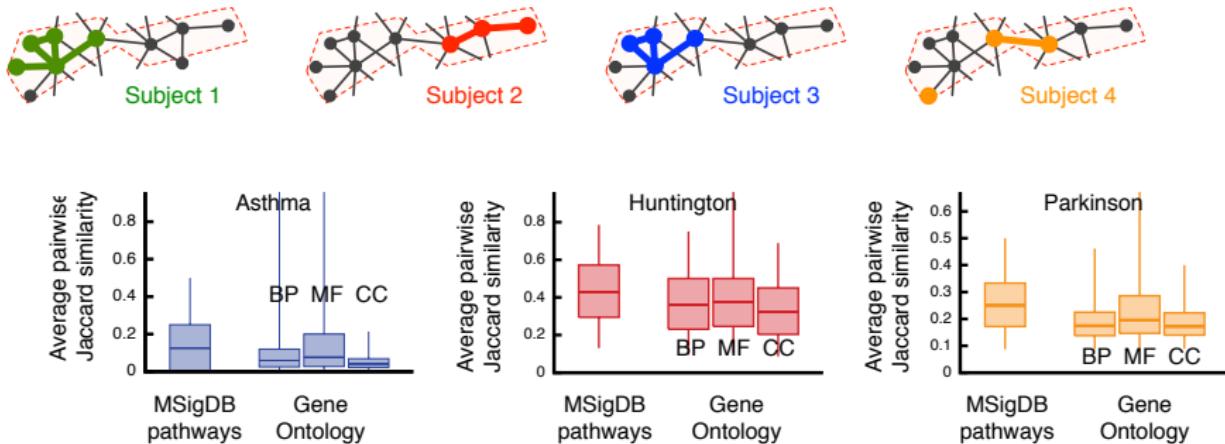
# Quantifying the heterogeneity using PEEPs



The overlap between PEEPs of two individuals with the same disease

- is low (< 30%), suggesting high heterogeneity at the transcription level
- is higher than the overlap between the PEEPs of healthy subjects

# From personalized gene-level signatures to pathway-level signatures

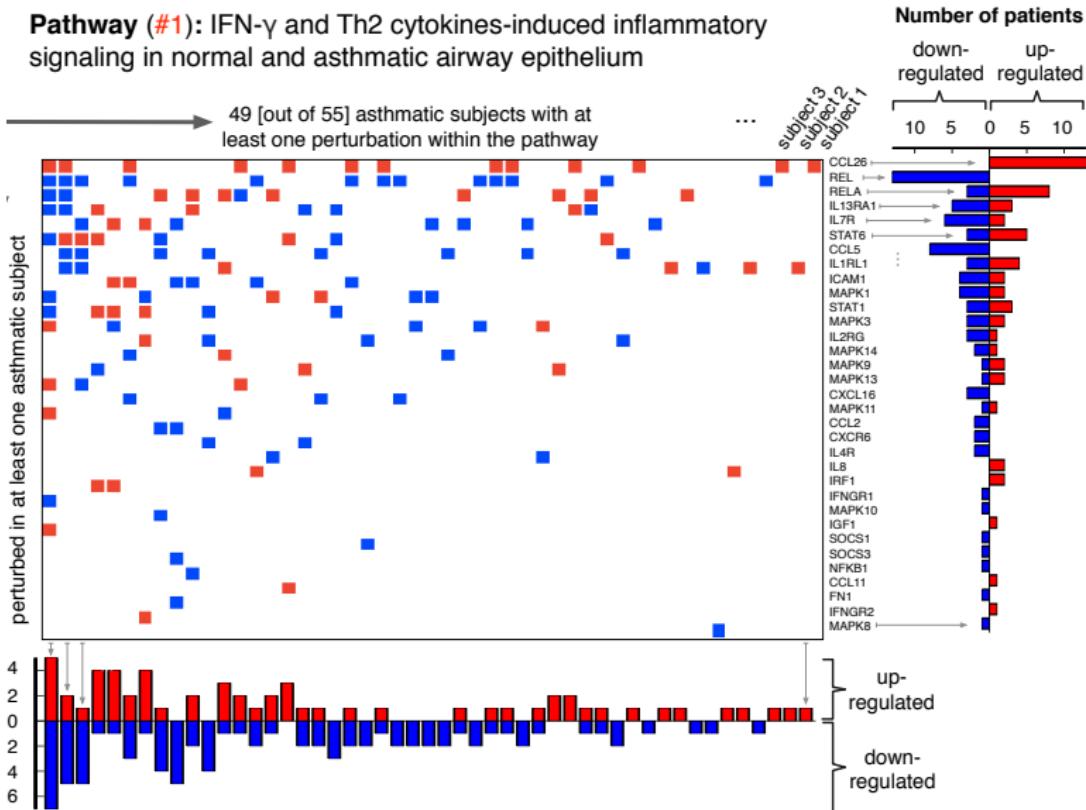


Enrichment of disease-specific pathways in PEEPs (assessed by Fisher's test followed by Bonferroni correction) reveals that

- almost all the individuals show significant perturbations in disease-specific pathways
- the specific perturbations differ greatly across subjects

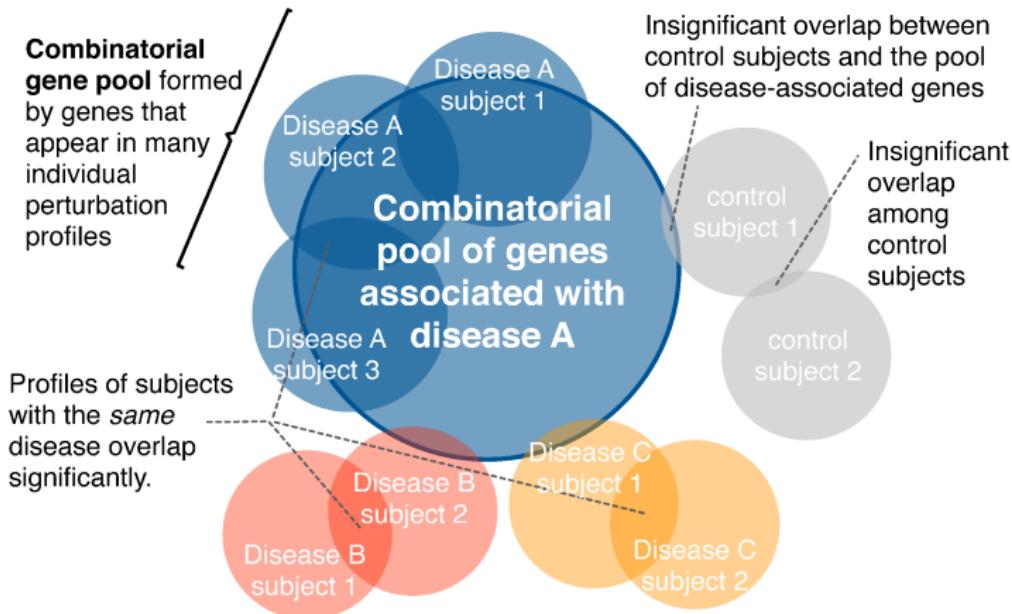
# Personalized pathway-level signatures

**Pathway (#1): IFN- $\gamma$  and Th2 cytokines-induced inflammatory signaling in normal and asthmatic airway epithelium**



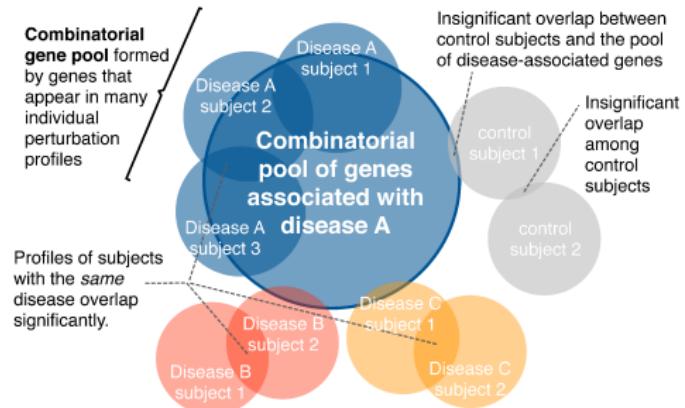
Pathways are perturbed in different ways in different patients

## Disease module hypothesis



- Different perturbations within certain molecular pathways lead to the same disease
- Genes that are perturbed in a significant fraction of the case subjects define the **disease module**

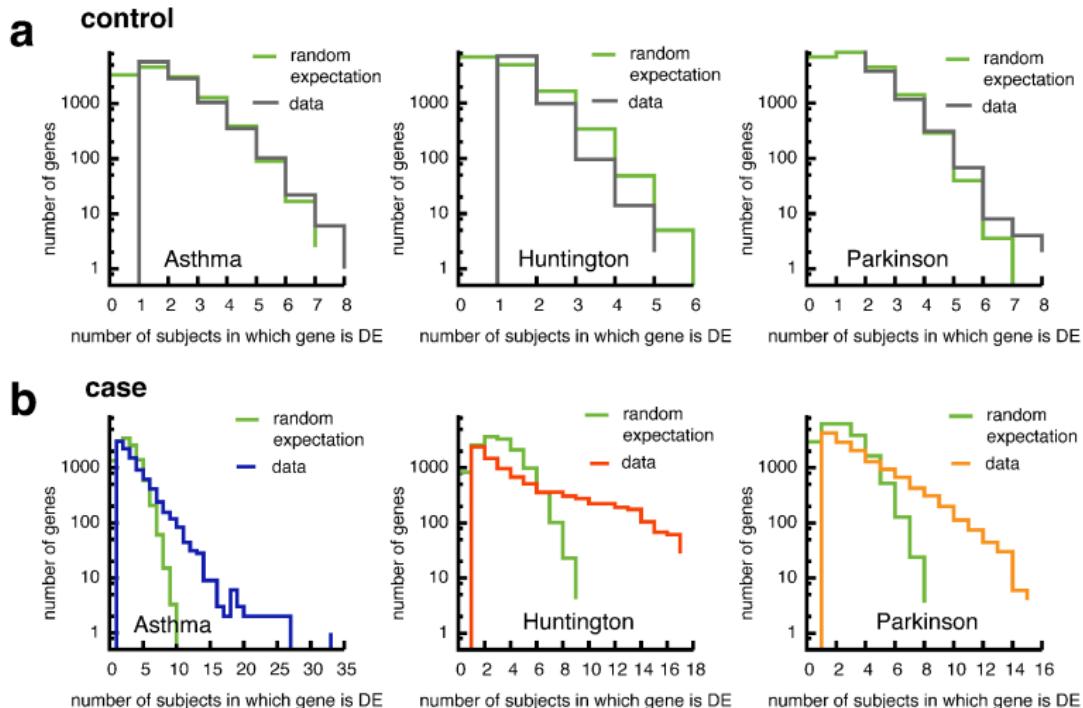
# Can PEEPs predict disease status?



Need to make sure that

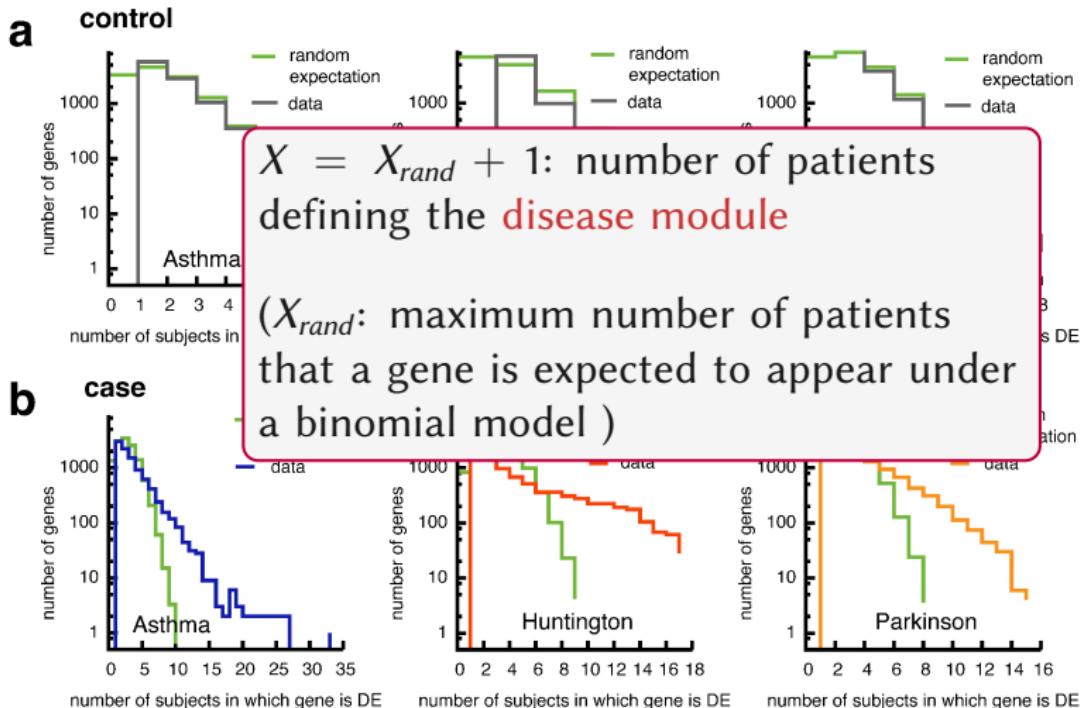
- PEEPs are not an artifact of technical or biological noise in transcriptomics data
- Disease module captures disease heterogeneity observed across patients

# Can PEEPs predict disease status?



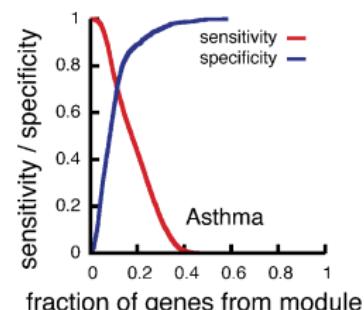
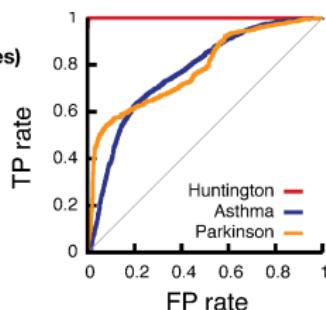
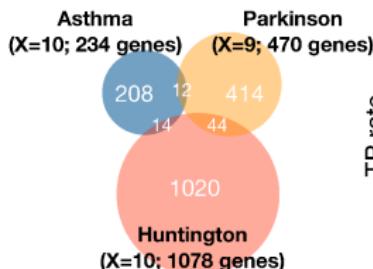
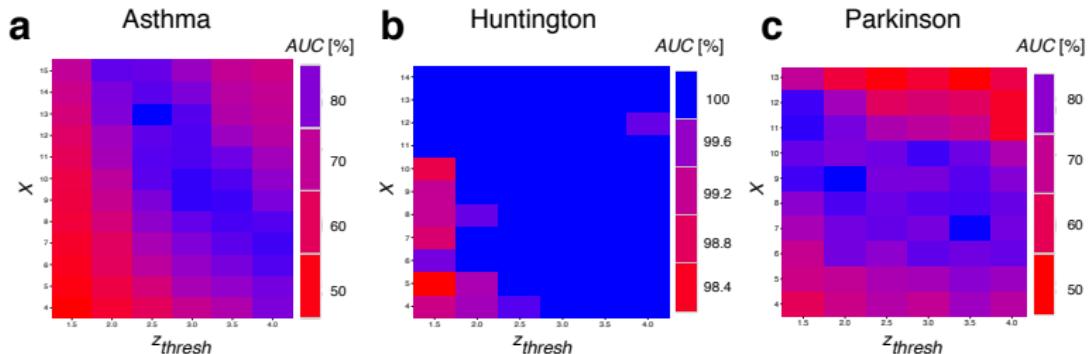
The number of shared genes between case subjects significantly exceeds the random expectation

# Can PEEPs predict disease status?



The number of shared genes between case subjects significantly exceeds the random expectation

# Predicting diseases from PEEPs



## Conclusions

- High heterogeneity across patients at the transcription level
- PErsonalized Expression Profiles (PEEPs)
  - quantify heterogeneity
  - define a pool of perturbed genes across patients
  - predict disease status

R package for generating PEEPs is available at  
[github.com/emreg00/pepper](https://github.com/emreg00/pepper)

## Acknowledgements



VIENNA SCIENCE  
AND TECHNOLOGY FUND



P50-HG004233, U01-HG001715, UO1-HG007690 from NHGRI

VRG15-005 from WWTF

PO1-HL083069, R37-HL061795, RC2-HL101543, U01-HL108630 from NHLBI

Beatriu de Pinós Fellowship from AGAUR

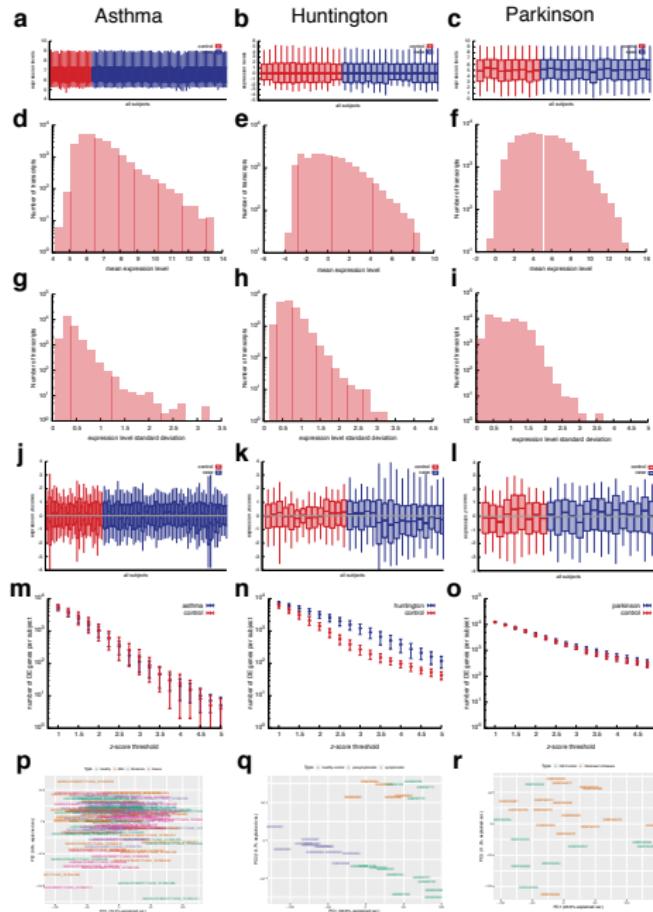
## Integrating personalized gene expression profiles into predictive disease-associated gene pools

Jörg Menche<sup>1,2,3</sup>, Emre Guney<sup>1,4</sup>, Amitabh Sharma<sup>1,4,5</sup>, Patrick J. Branigan<sup>6</sup>, Matthew J. Loza<sup>6</sup>, Frédéric Baribaud<sup>6</sup>, Radu Dobrin<sup>6</sup> and Albert-László Barabási<sup>1,2,4,5</sup>

Marc Santolini & Alan Karma

R package for generating PEEPs is available at  
[github.com/emreg00/pepper](https://github.com/emreg00/pepper)

# Preprocessing and normalization of transcriptomics data



## Defining disease modules

$X$ , the number of patients defining the disease module is given by

$$X = X_{rand} + 1$$

where  $X_{rand}$  given by

$$\operatorname{argmax}_{X_{rand}} \left( \sum_{k=X_{rand}}^n G * f(k; n, p) < 1 \right)$$

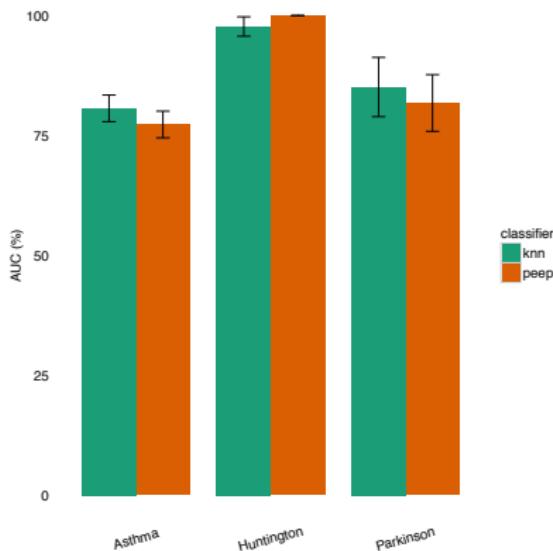
and

$$f(k; n, p) = Pr(x = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

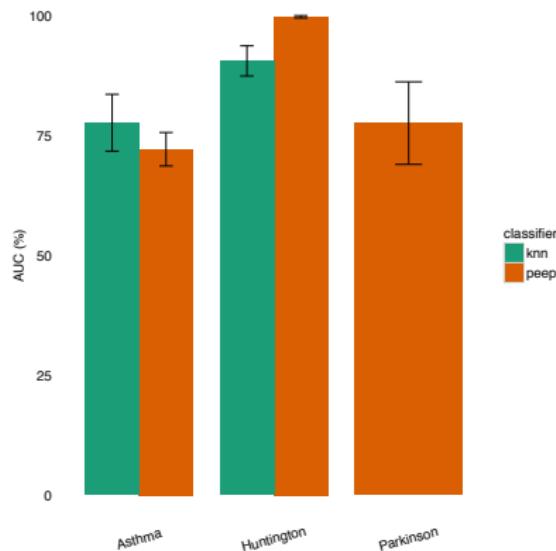
corresponding to the probability for a gene to be perturbed in exactly  $k$  out of  $n$  subjects given that each subject has  $g$  perturbed genes on average (out of all possible  $G$  genes in the expression profiling platform) and thus,  $p = g/G$

# PEEP vs K-nearest-neighbor based predictions

**a** Five-fold cross-validation

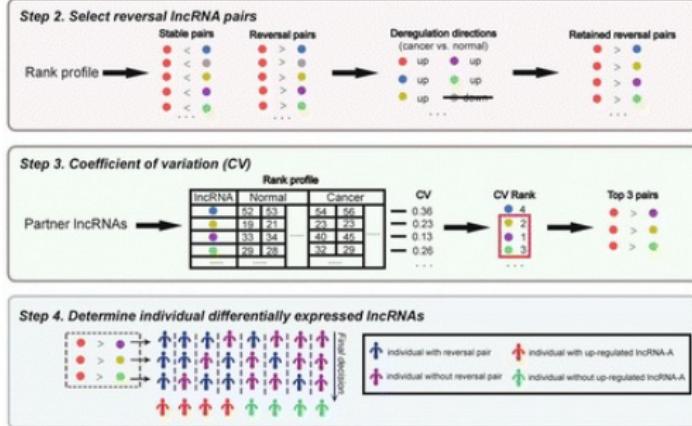
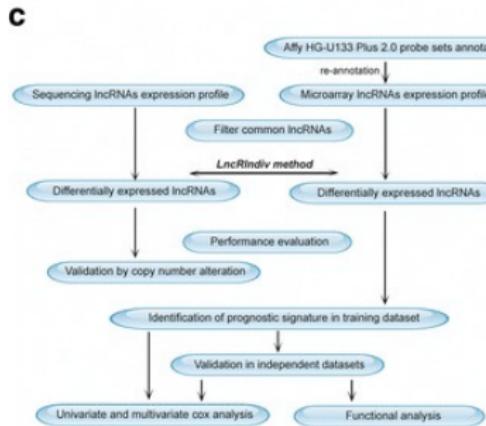
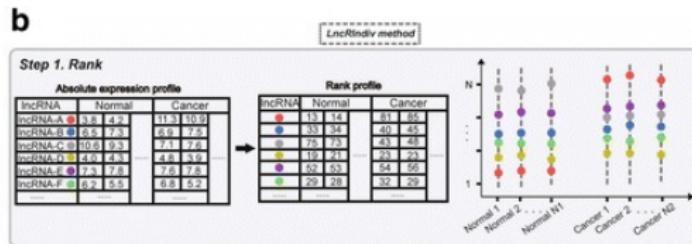
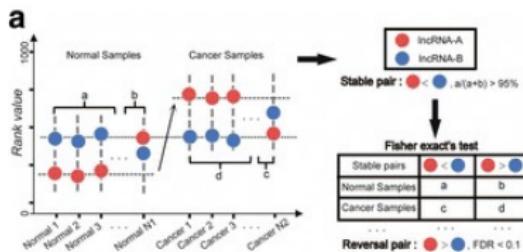


**b** Two-fold cross-validation



- Results with  $K = 15$  (maximizes accuracy among 3, 5, 10, 15, 20)
- AUCs are calculated using a repeated cross validation approach (100 repetitions)

# Rank-based approaches for calculating individual perturbation profiles



RankComp Wang et al., 2015, Bioinformatics

LncRIndiv Peng et al., 2017, Mol Cancer