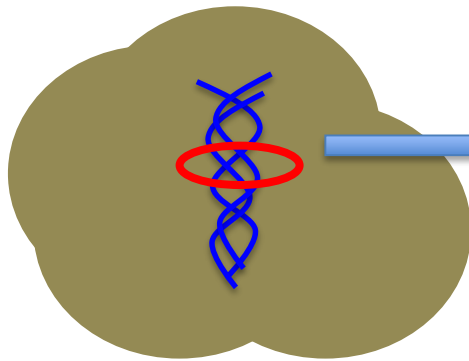


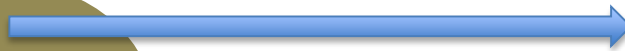
# Deciphering Signatures of Mutational Processes Operative in Human Cancer

# Tumor Cells Carry Somatic Mutations

**Tumor**



**Sequence**



gcttcgctagcgcccccttttaatcgatcccgatcg  
cccacgatcgatagctagatcgactgttttaatt  
agccacatcactatctcccttttgggagacgatc  
atgccccggtttcgaatgctaaaatgctaaagtt  
cccacgatcggatagctagatcgactgttttaatt  
cagctactgatcgttttgccggccccccgggagat  
atgccccggtttcgaatgctaaaatgctaaagtt



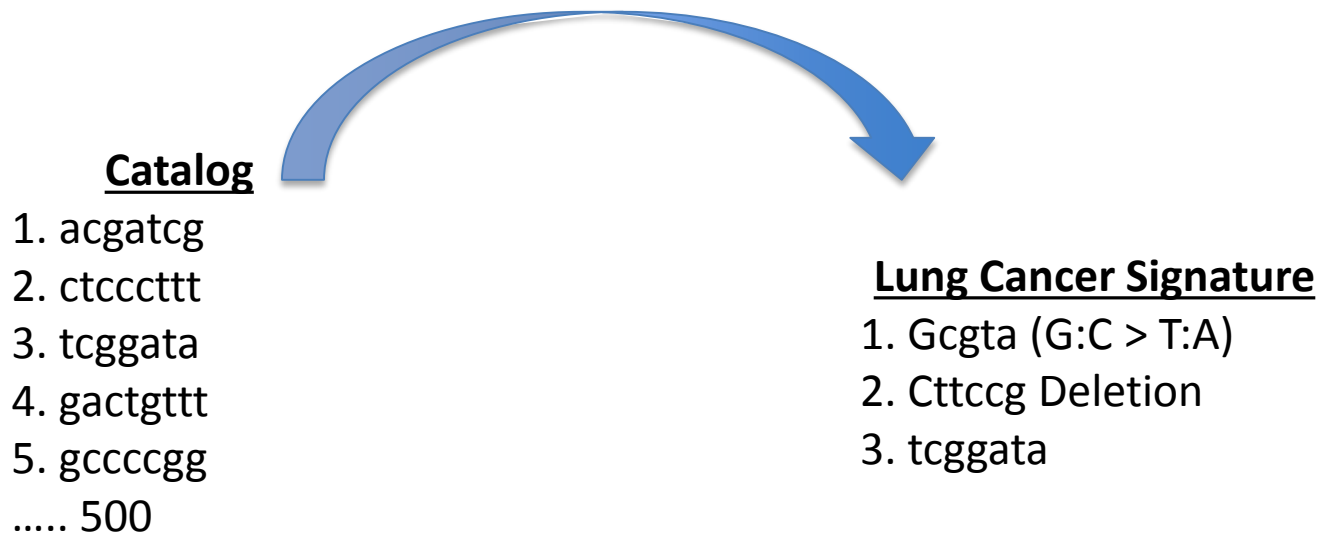
**Catalog**

1. acgatcg
2. ctcccttt
3. tcggata
4. gactgttt
5. gccccgg
- ..... 500

# Motivation

- Catalogs have heterogeneity
  - Different mutation types: Substitution, missense, nonsense, indels
  - DNA Repair mechanisms
  - Passenger mutations
- Many different cancer signatures

Aim to create computational framework to bridge the gap between the catalogs and signatures



# Feature of Signatures

$$P_1 = [p_1^1, p_1^2, \dots, p_1^K]^T$$

P = Mutational Signature

$p_{1\dots k}$  = probability P causes a certain mutation

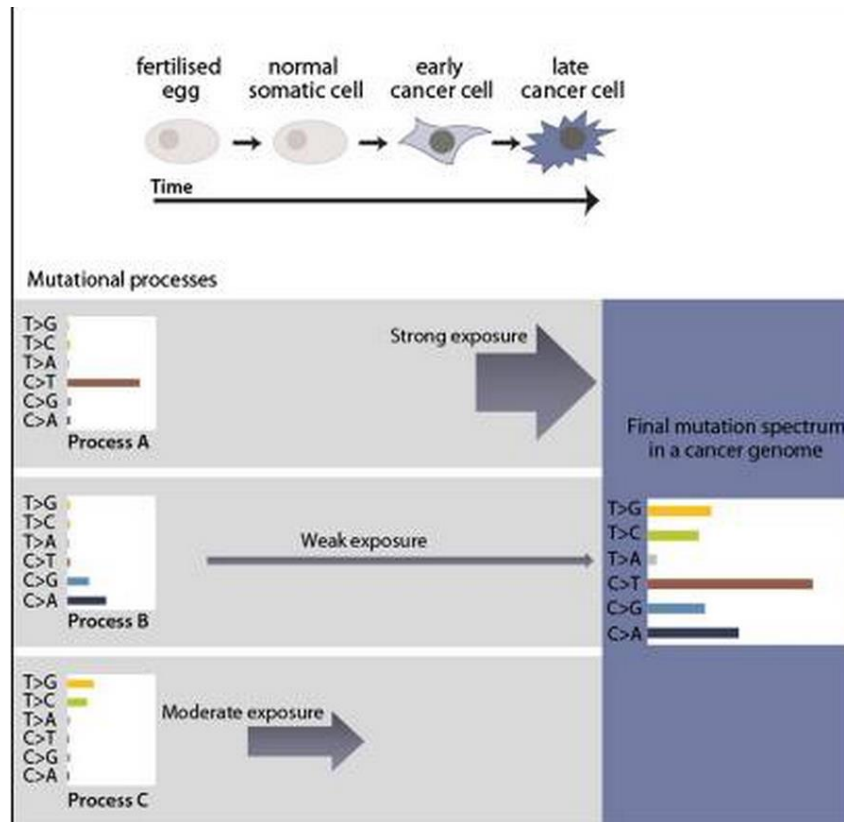
K = 96 (6 types of substitutions \* 4 types of 5' bases \* 4 types of 3' bases)

# Mapping of a Genome

$$m_g^i \approx \sum_{n=1}^N p_n^i e_g^n.$$

P = process/mutation

e = exposure/weight



# What we end up with

$$\begin{aligned}
 P &= \begin{bmatrix} p_1^1 & p_2^1 & \cdots & p_{N-1}^1 & p_N^1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ p_1^K & p_2^K & \cdots & p_{N-1}^K & p_N^K \end{bmatrix} \\
 &\quad \times \\
 E &= \begin{bmatrix} e_1^1 & e_2^1 & \cdots & e_{G-1}^1 & e_G^1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ e_1^N & e_2^N & \cdots & e_{G-1}^N & e_G^N \end{bmatrix} \\
 &= M = \begin{bmatrix} m_1^1 & m_2^1 & \cdots & m_{G-1}^1 & m_G^1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ m_1^K & m_2^K & \cdots & m_{G-1}^K & m_G^K \end{bmatrix}
 \end{aligned}$$

# Non-Negative Matrix Factorization

- Want to extract “P” and “e” from M

## **Step 1 and 2**

Reduce Matrix Dimensions

$$\sum_{r \in R} \sum_{g=1}^G m_g^r \leq 0.01 \times \sum_{k=1}^K \sum_{g=1}^G m_g^k,$$

Use bootstrap resampling



## ***Step 3&4: Non Negative Matrix Factorization***

- All inputs must be non-negative
- Aims to recreate P and e from M

Iterate until convergence

$$e_G^N \leftarrow e_G^N \frac{[P^T \widetilde{M}]_{N,G}}{[P^T P E]_{N,G}}$$

$$p_N^K \leftarrow p_N^K \frac{[\widetilde{M} E^T]_{K,N}}{[P E E^T]_{K,N}}$$

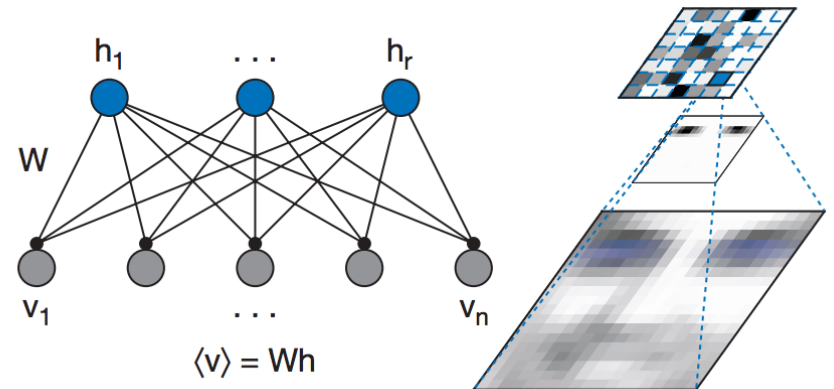
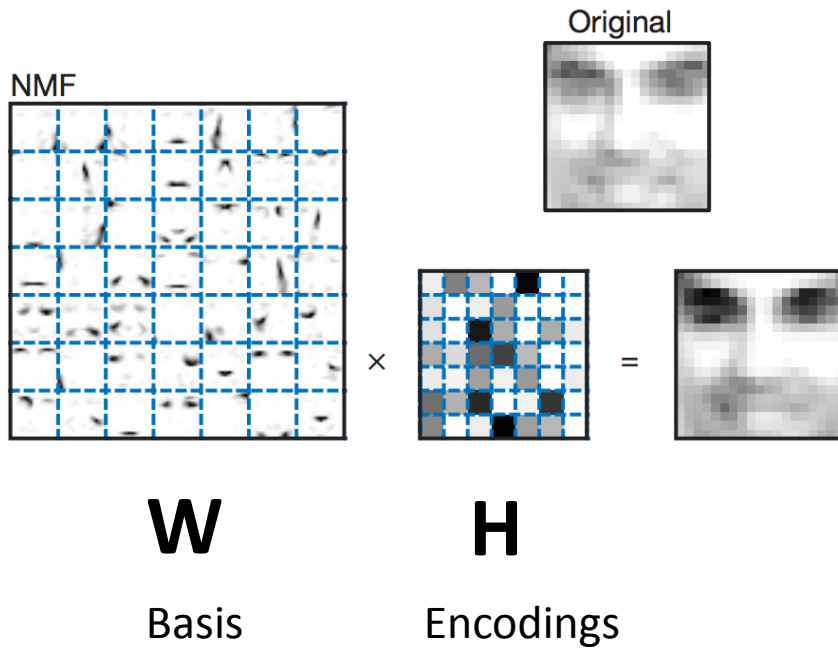
Minimize

**Cost Function**

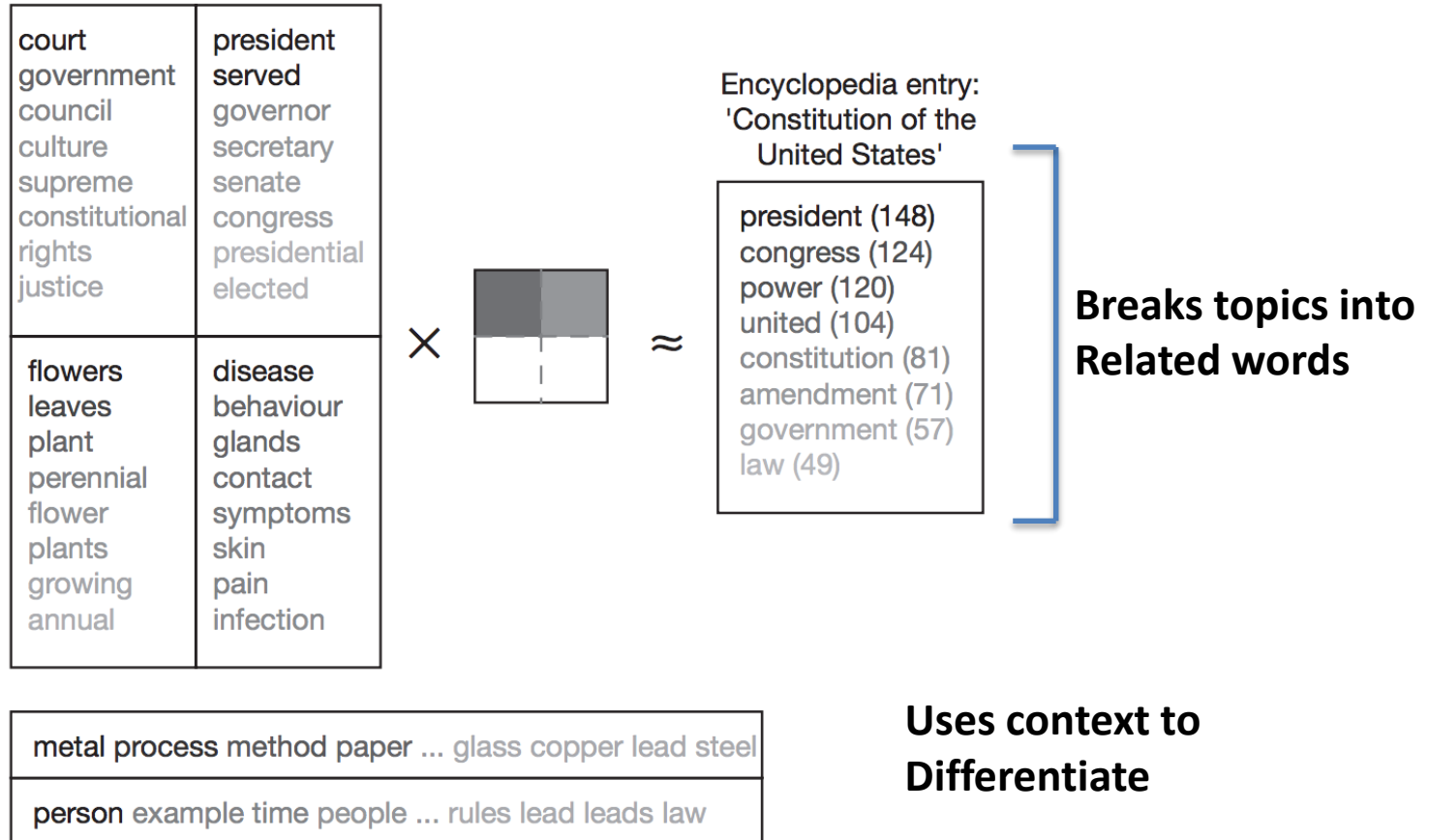
$$\|\widetilde{M} - P \times E\|_F^2$$

Equivalent to (K,N)<sup>th</sup> element of matrix

# NMF: Faces



# NMF: Encyclopedia



## ***Step 5: Clustering***

- Partition-clustering algorithm was applied to cluster data into N clusters

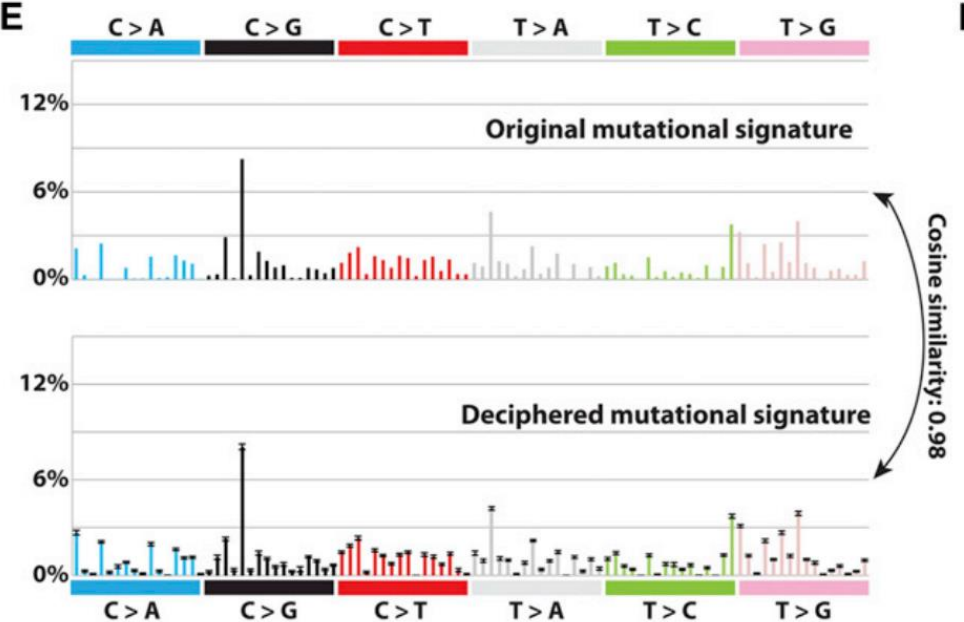
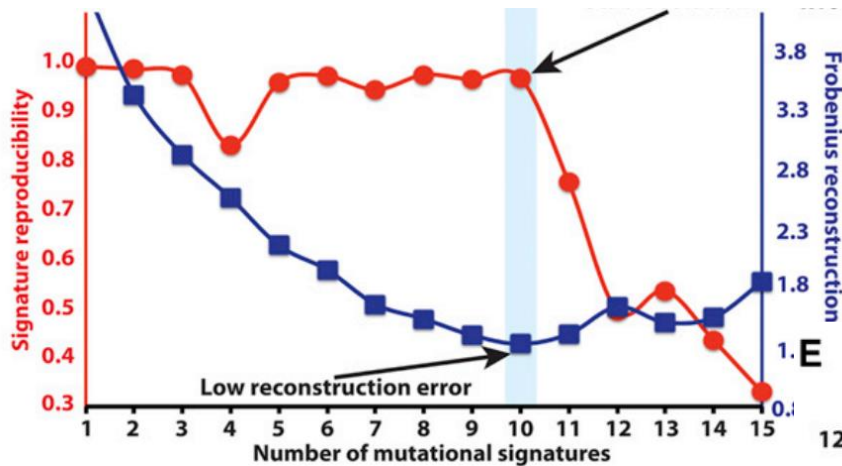
## ***Step 6: Evaluate***

- Look at Frobenius reconstruction error to evaluate for accuracy
- Compare mutational signatures:

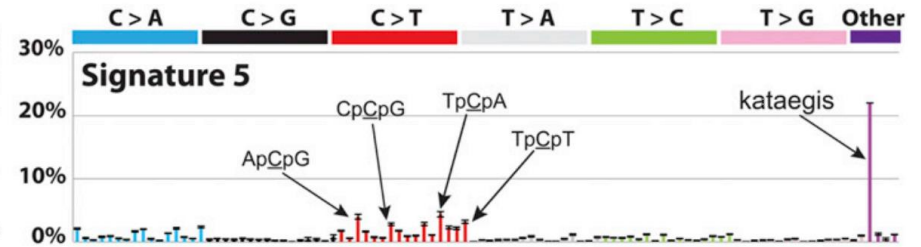
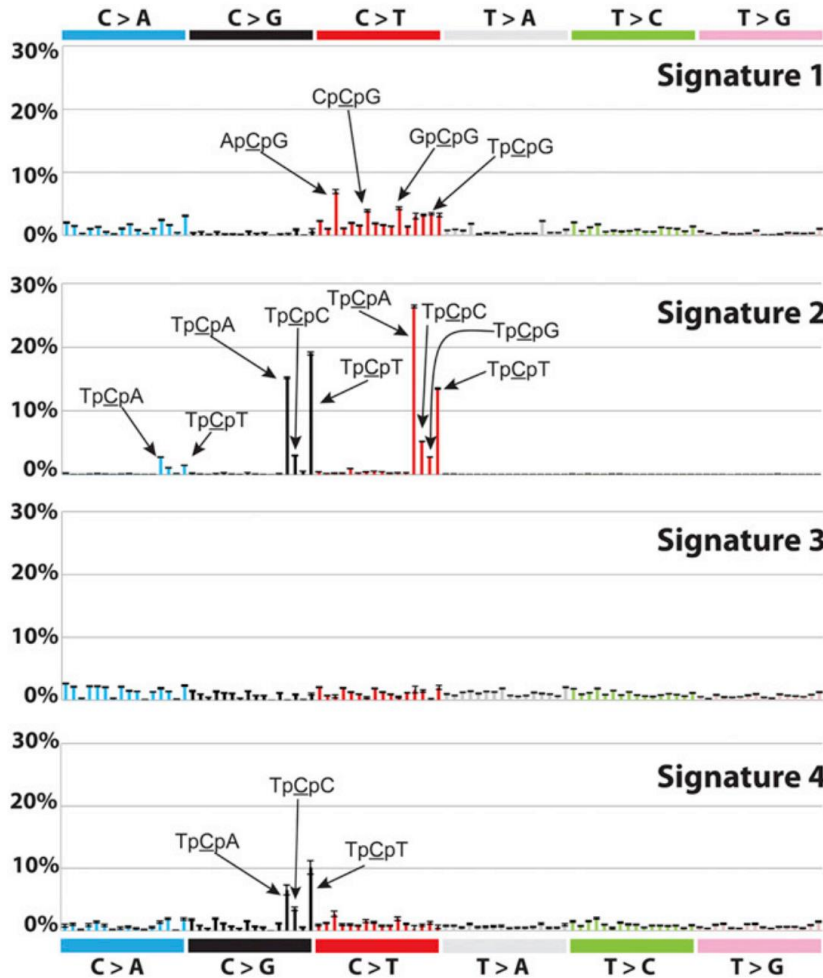
$$\text{sim}(A, B) = \frac{\sum_{k=1}^K A_k B_k}{\sqrt{\sum_{k=1}^K (A_k)^2} \sqrt{\sum_{k=1}^K (B_k)^2}}.$$

$\text{Sim}(A, B) = 1$  means same signature

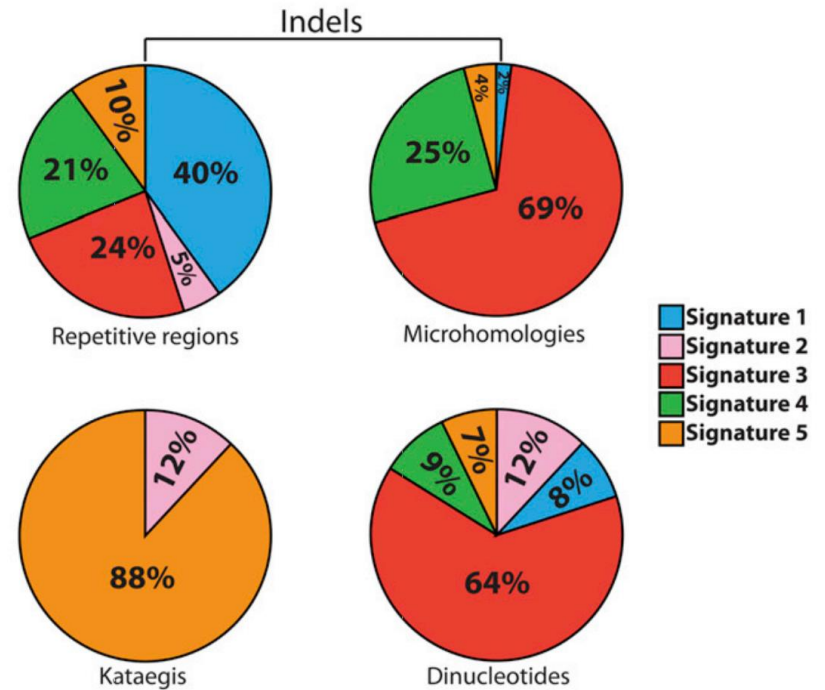
# Does it work?



# Breast Cancer Example



C



# Impact

- Ability to generate cancer signatures from comprehensive 'omic data
- Opens the door for further work. Eg. Sparsity constraint to use a minimum number of signatures