

Predicting the effects of mutations with deep generative models

Adam Riesselman

Predicting the effects of mutations with deep generative models

Adam Riesselman



John Ingraham



Debbie Marks

Proteins are the workhorses of biology

DNA

ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTCGCCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACTCGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATACCAAACGACGAG...

RNA

Protein

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERD TTMPAAM
ATTLRKL LTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

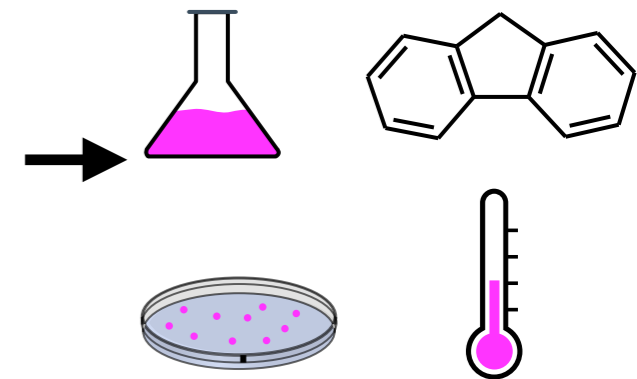
20 different amino acids

4 different bases

Structure



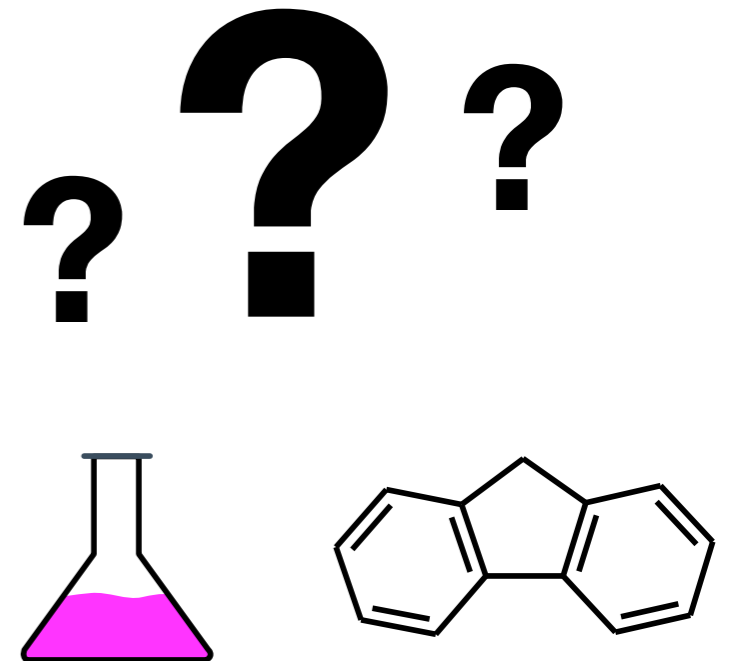
Function



Mutations impact protein function

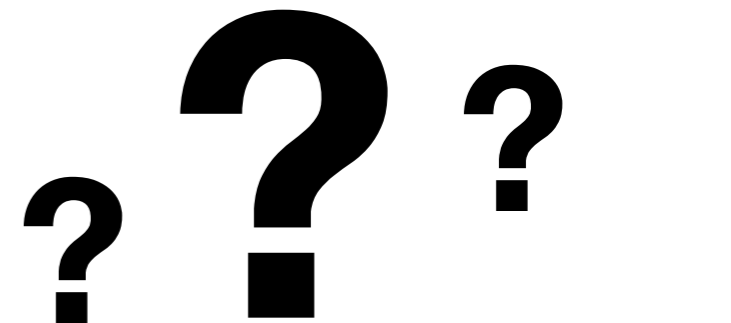
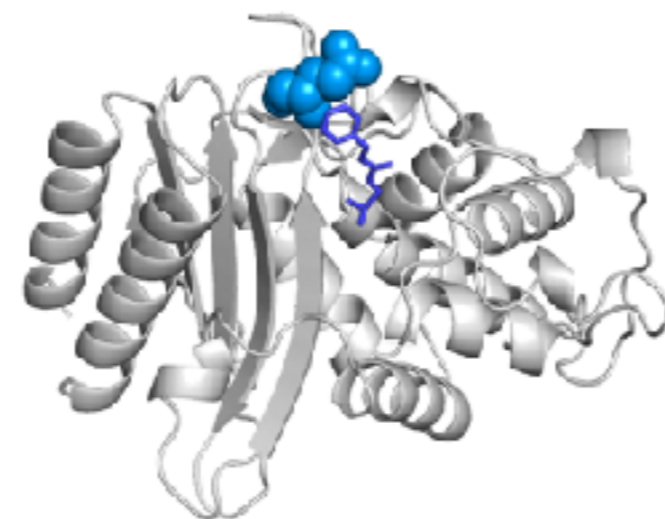
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → I



Sequence

Structure

Function

Mutation effect prediction is important

Understanding
disease

Biomedicine

Bioengineering

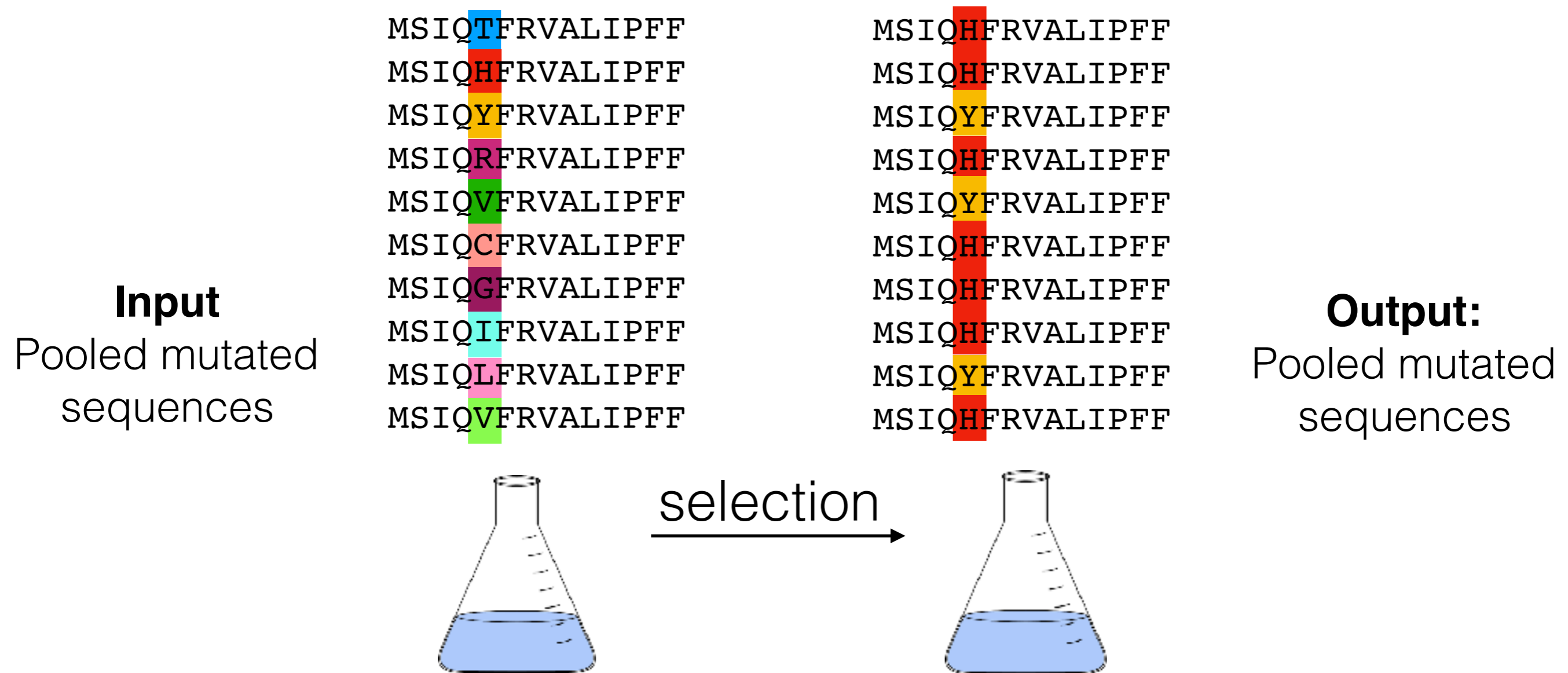


“Does this mutation
cause cancer?”

“Is this antibody
stable in a patient?”

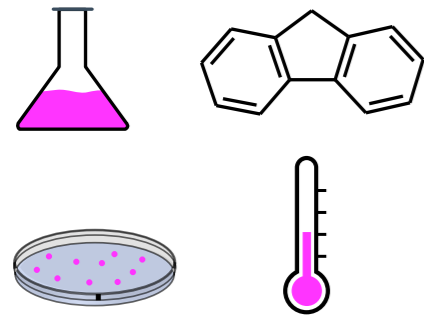
“Will this protein digest
wood better for
biofuels?”

State of art methods for measuring mutation effects

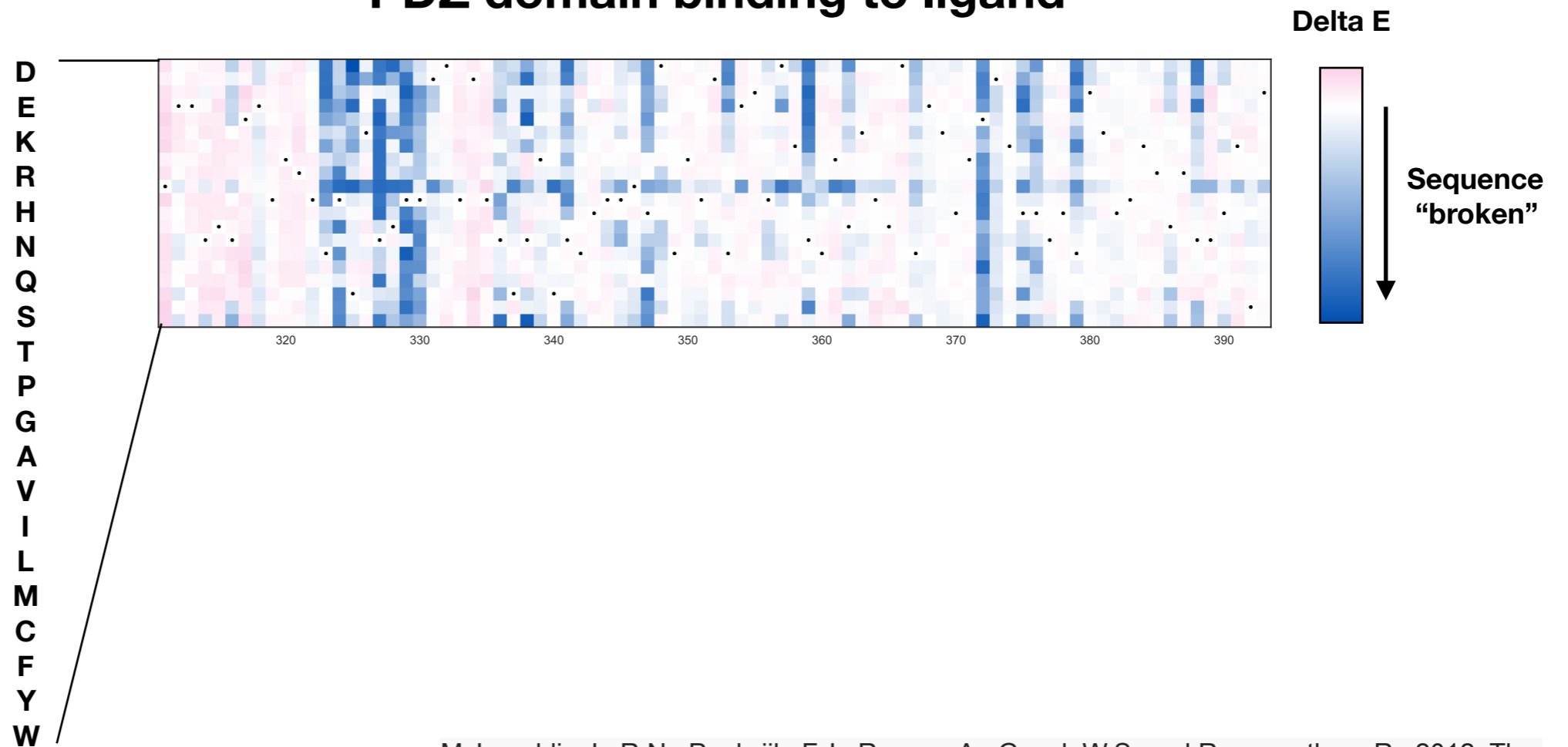


Compare ratio of sequences before and after selection

State of art methods for measuring mutation effects

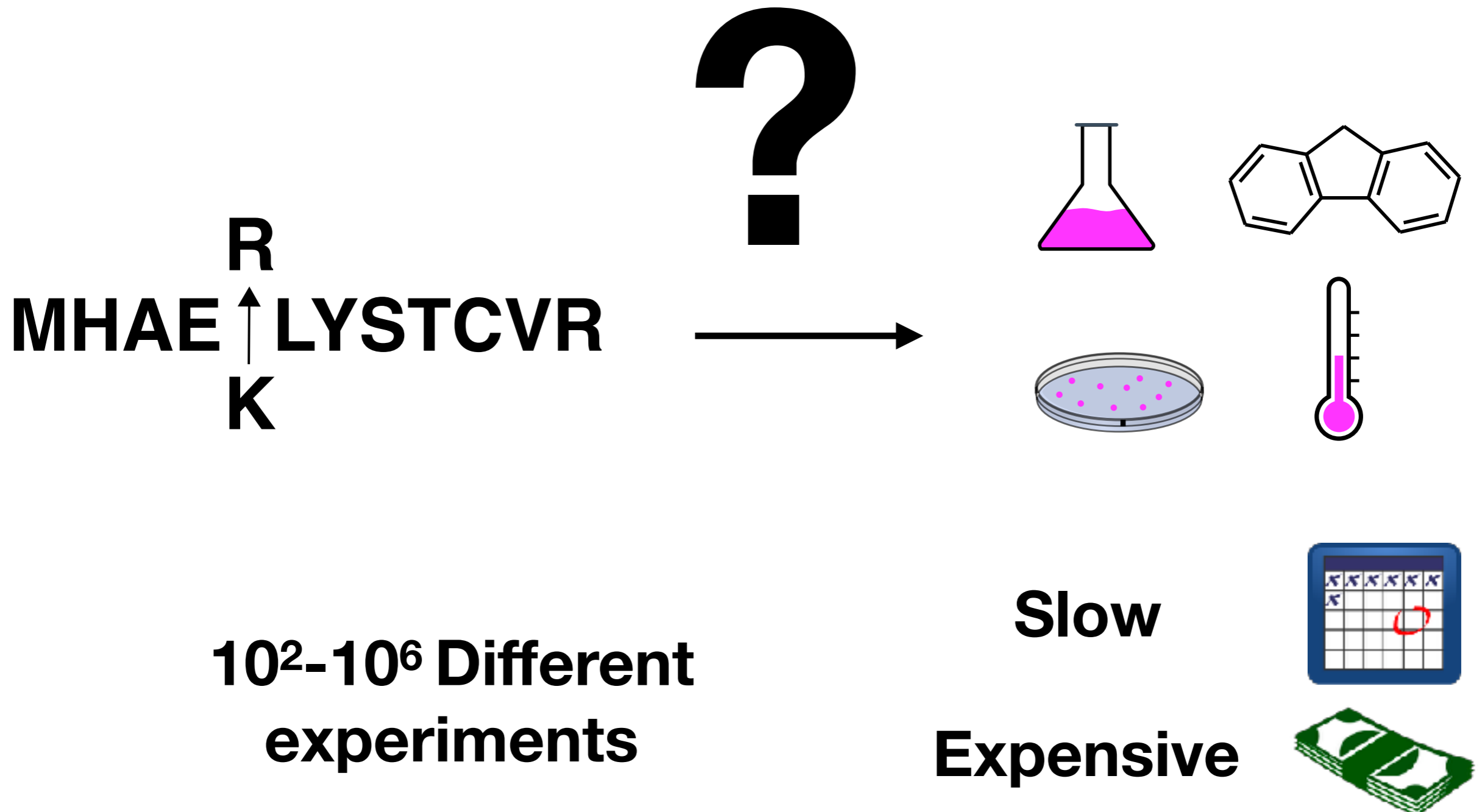


PDZ domain binding to ligand



McLaughlin Jr, R.N., Poelwijk, F.J., Raman, A., Gosal, W.S. and Ranganathan, R., 2012. The spatial architecture of protein function and adaptation. *Nature*, 491(7422), pp.138-142.

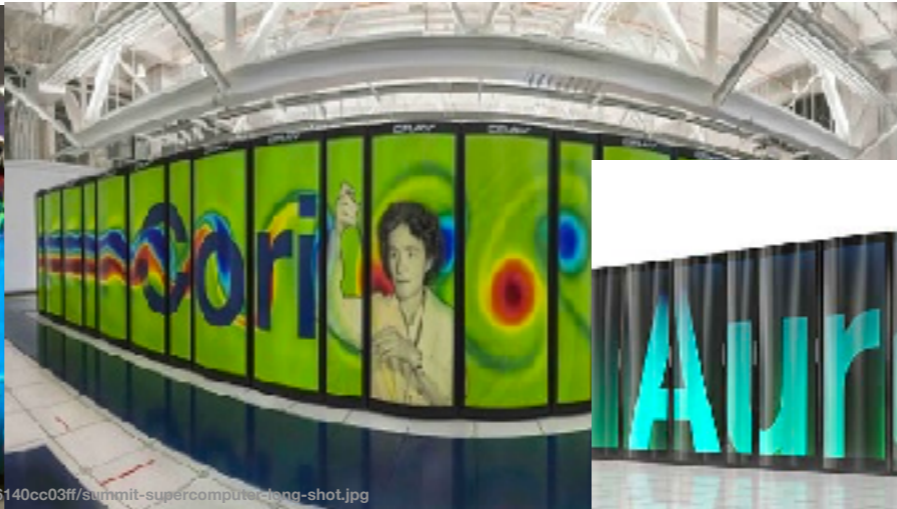
Understanding the effects of mutations is important





There has to be a better way...

with computers!



<https://www.extremetech.com/extreme/99413-titan-jaguar-computer-38400-processor-20-petaflop-successor-to-jaguar>
<http://ichef.bbci.co.uk/news/6/0/1/1/76/p01076fr.jpg>
<https://6111539m39y3h4e1qun0z2f1-wpeengine.netdna-ssl.com/wp-content/uploads/2017/04/Cori-NERSC-405x228.jpg>
<https://6111539m39y3h4e1qun0z2f1-wpeengine.netdna-ssl.com/wp-content/uploads/2016/12/aurora-675x380.jpg>
<https://cnet3.cbsistatic.com/img/7Af4Ub5p3mWfc1Lb0XK-QE8w8-/610x503/2018/06/08/06b0304d-1fc5-428a-9399-e76140cc03ff/summit-supercomputer-long-shot.jpg>
<https://3c1703fe8d.site-internapcdn.net/newman/csz/news/800/2016/summitting-down.jpg>
<https://www.extremetech.com/wp-content/uploads/2013/01/Hopper-1-348x196.jpg>
https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcRg9avPn1s3pNCCpRbisdm1cnROZqUFOT_KF7AgKShwBJvjf7Q
https://c1.staticflickr.com/6/5599/31681202785_24374e416b_b.jpg

Lots of other **examples** of **Protein X** are available



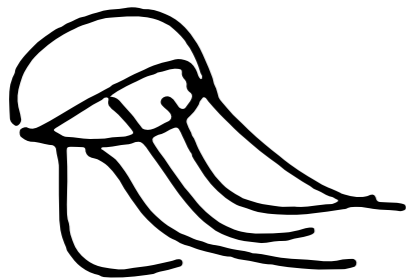
ADRLYMTKIHHEFEGD



ADRLYMTKIH HQFDGD

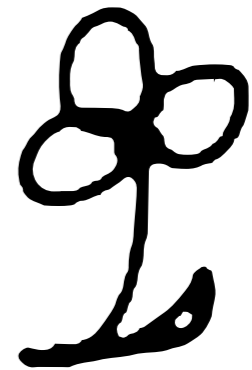
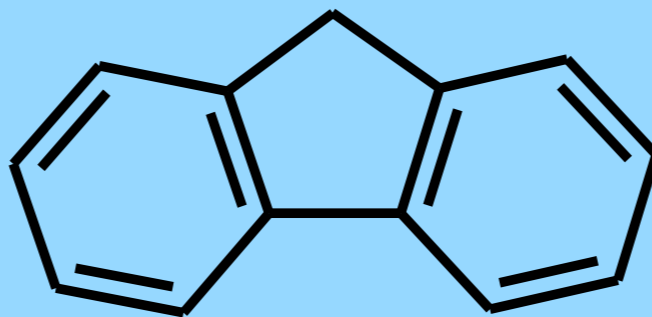


ADRLYMTKIHHEFEGD



ADRLYL TQIRNKFKGD

Protein X



TSKMYITKIGQEFEGD

All are **functional, homologous examples** of Protein X

Lots of other **examples** of **Protein X** are **available**

Sequences are found in public **genome databases.**



Natural **evolution** is an **experiment**, in **parallel.**

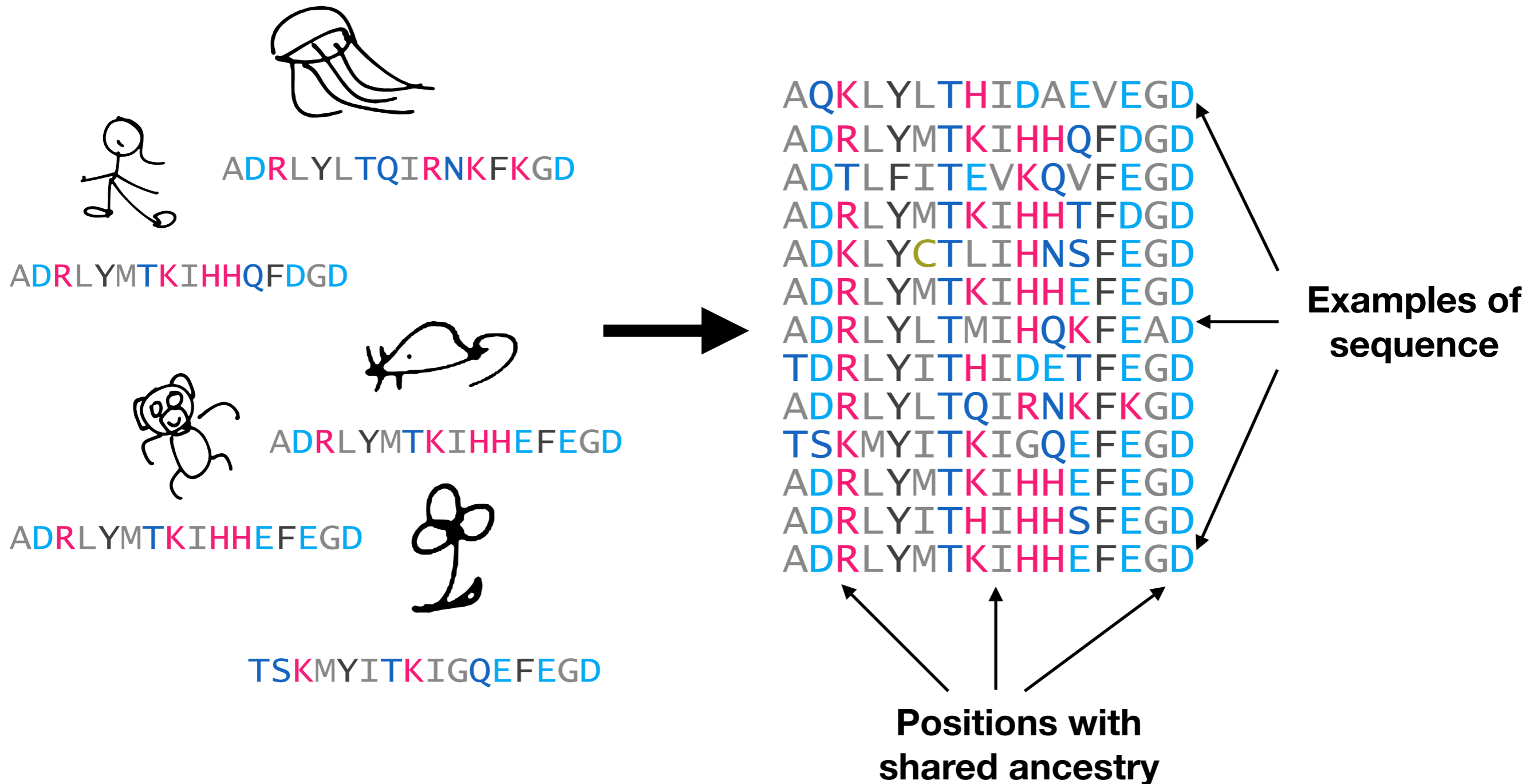
Assumption:

Present in database: Tolerated

Not in database: Deleterious

All are **functional, homologous examples** of Protein X

Natural sequences can be grouped into families via alignments





Dihydrofolate reductase
a billion years of data

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD
TSKMYITKIGQEFEGD
ADRLYMTKIHHEFEGD
ADRLYITHIHHSFEGD
ADRLYMTKIHHEFEGD

Natural sequences have been evolved to be **functional**

x

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD
TSKMYITKIGQEFEGD
ADRLYMTKIHHEFEGD
ADRLYITHIHHSFEGD
ADRLYMTKIHHEFEGD

Fit generative
model



P(x)

Generative model captures **functional constraints**

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD

↓
 $p(\mathbf{x}|\theta)$

2) Compute **Log Ratio**
for each mutant

$$\log \frac{p(\mathbf{x}_{\text{mut}}|\theta)}{p(\mathbf{x}_{\text{wild}}|\theta)}$$

“How much does this mutation look like what we’ve seen in nature?”

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio** for each mutant

Uses public data (effectively free)

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFKGD

Fast

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \theta)}{p(\mathbf{x}_{\text{wild}} | \theta)}$$

Works on almost any protein

↓
 $p(\mathbf{x} | \theta)$

Accurate

“How much does this mutation look like what we’ve seen in nature?”

First-pass mutation predictors model **evolutionary conservation**

Amino acid i



■ Positive ■ Polar ■ Hydrophobic
■ Negative ■ Cysteine ■ Aromatic

Product of **site** factors

$$P(\mathbf{x}) = p_1(x_1)p_2(x_2) \cdots p_L(x_L)$$

R → K Neutral

R → I Deleterious

How to capture **interactions**?



Pairwise undirected model

$$P(\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_{i < j} J_{ij}(x_i, x_j) + \sum_i h_i(x_i) \right)$$

a.k.a.

Markov Random Field

a.k.a.

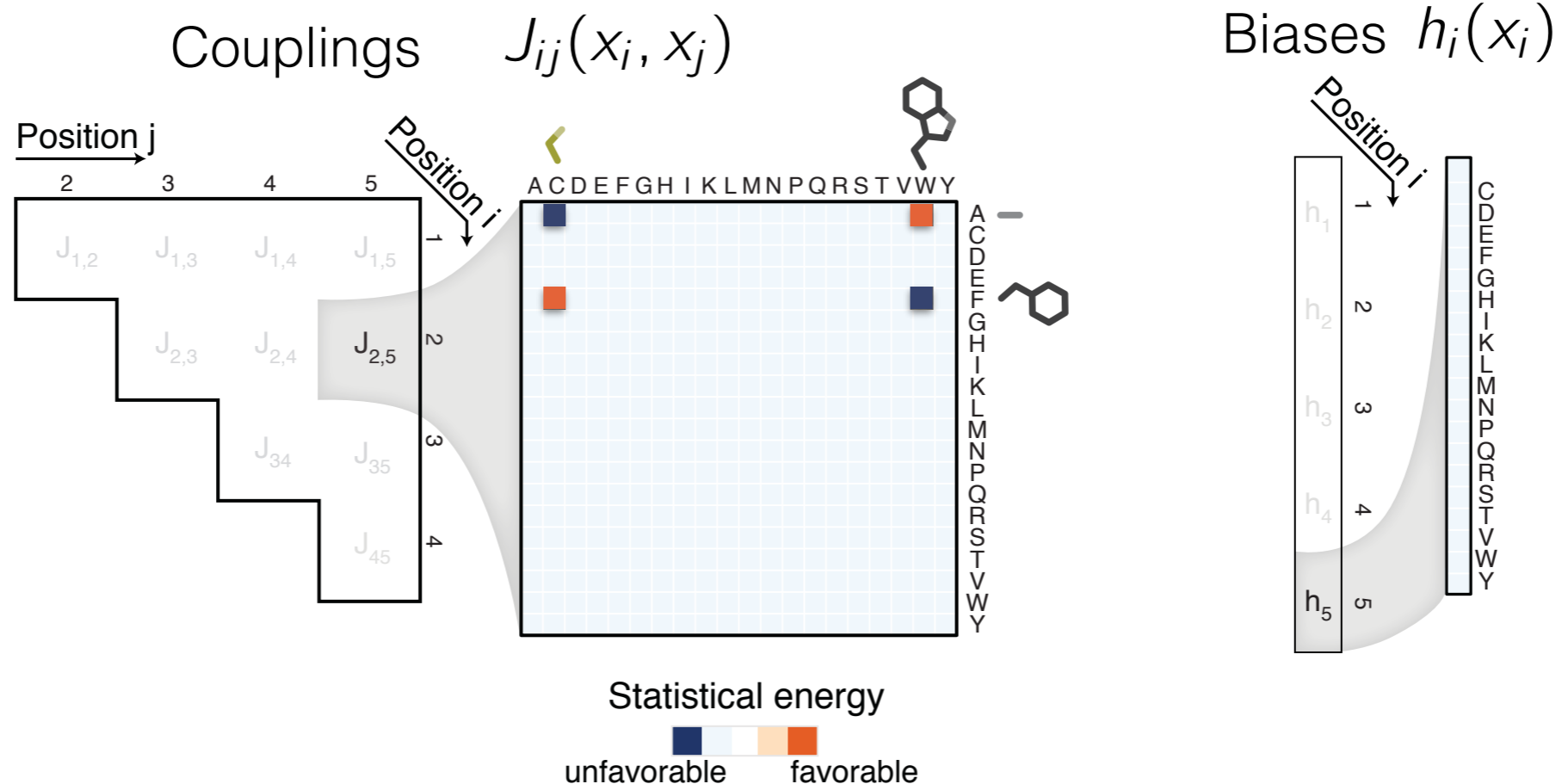
Ising (Potts) model

a.k.a.

Multinomial logistic regression

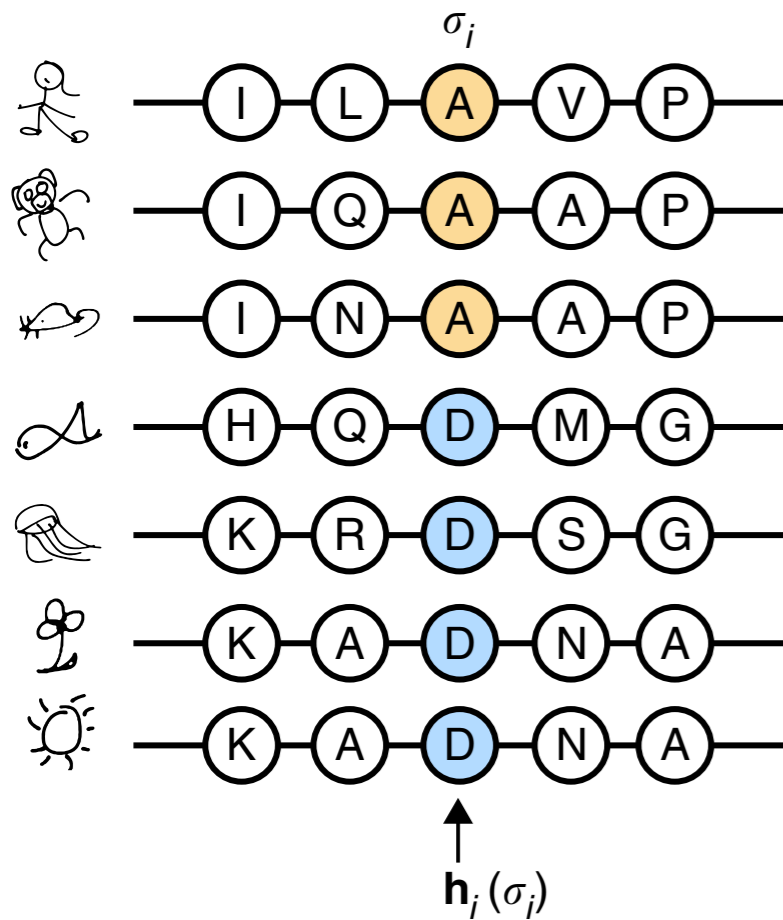
Undirected graphical model for sequences parameterizes **pairs of letters** at **pairs of positions**

$$P(\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_{i < j} J_{ij}(x_i, x_j) + \sum_i h_i(x_i) \right)$$

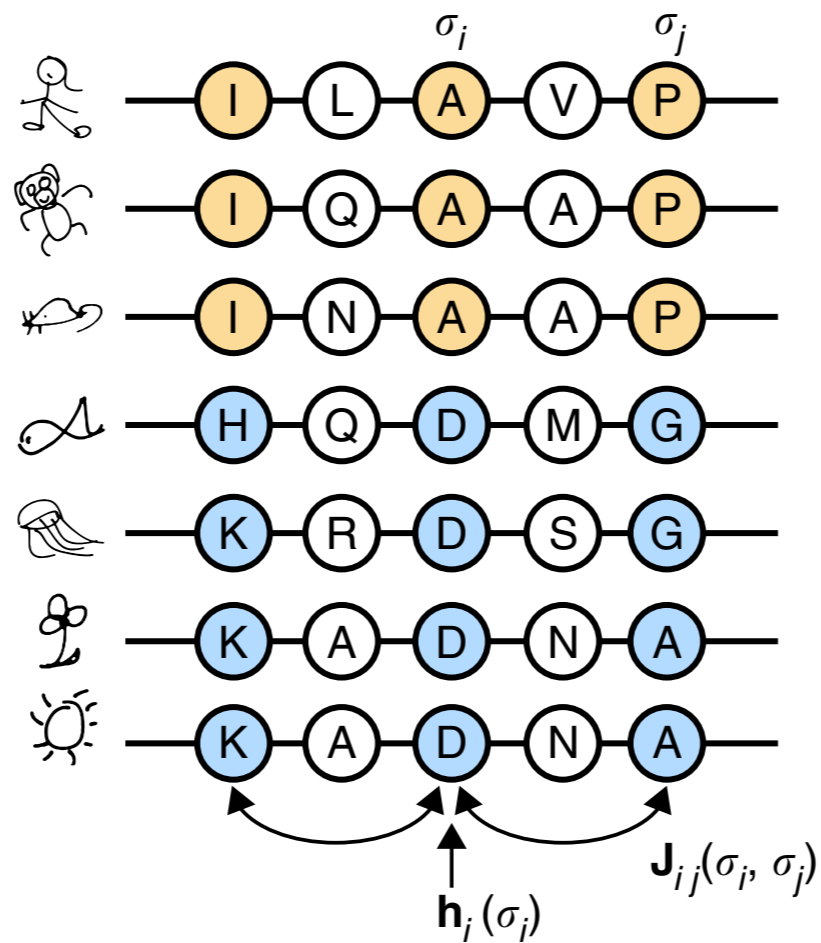


Pairwise \gg sitewise. Should we stop there?

Independence



Pairwise model



Higher order?

...

Pairwise interactions insufficient for mutation effects in **proteins**

Dihydrofolate
reductase

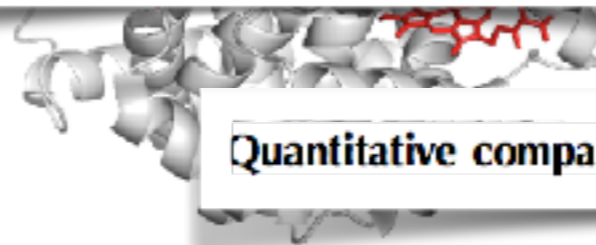
The landscape is shaped by high order genetic interactions. The fitness landscape has extensive high-order genetic interactions. A series of models of increasing complexity were constructed that described the $\log(\text{IC}_{75})$ of each genotype as a sum of parameters (equivalent to multiplying fold-changes in IC_{75}) that



Hemoglobin

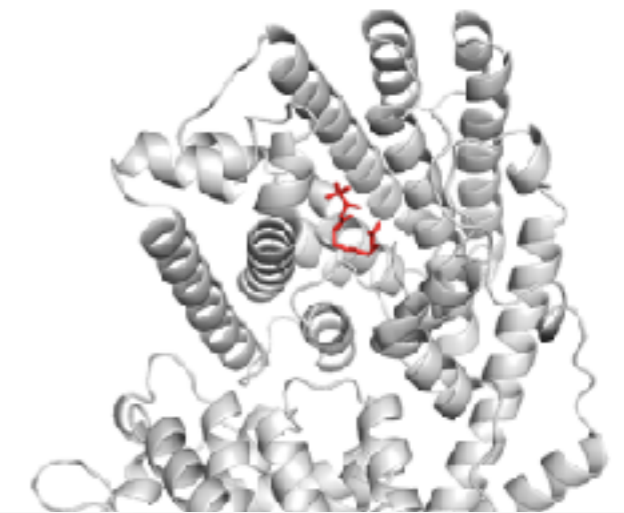


In summary, results of our mutagenesis experiments revealed pervasive epistasis among segregating amino acid variants in deer mouse Hb (Table 1). The individual and joint effects



Quantitative comparisons indicated context dependence for mutational effects.

Terpene
Synthase

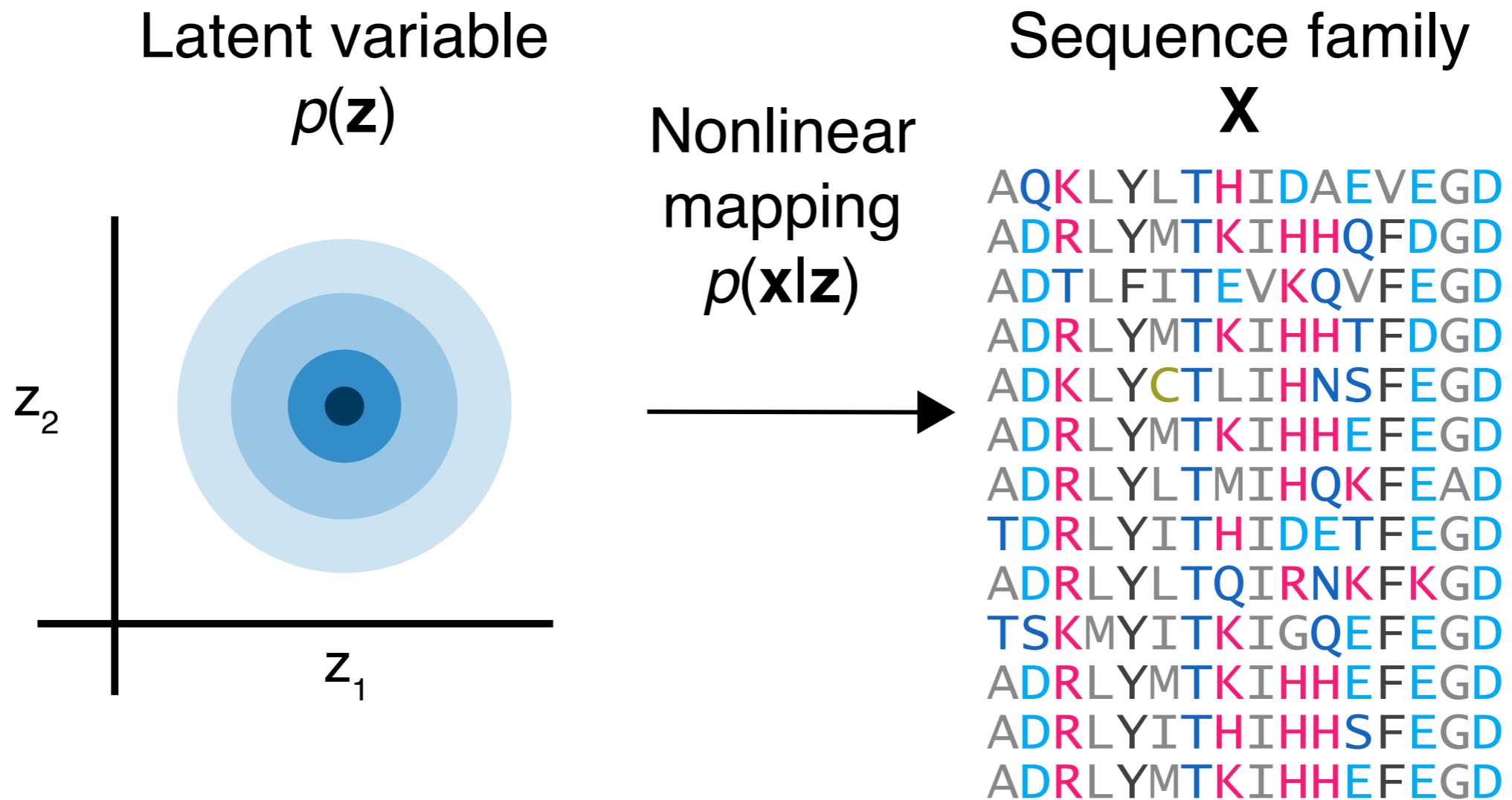


Drug resistance

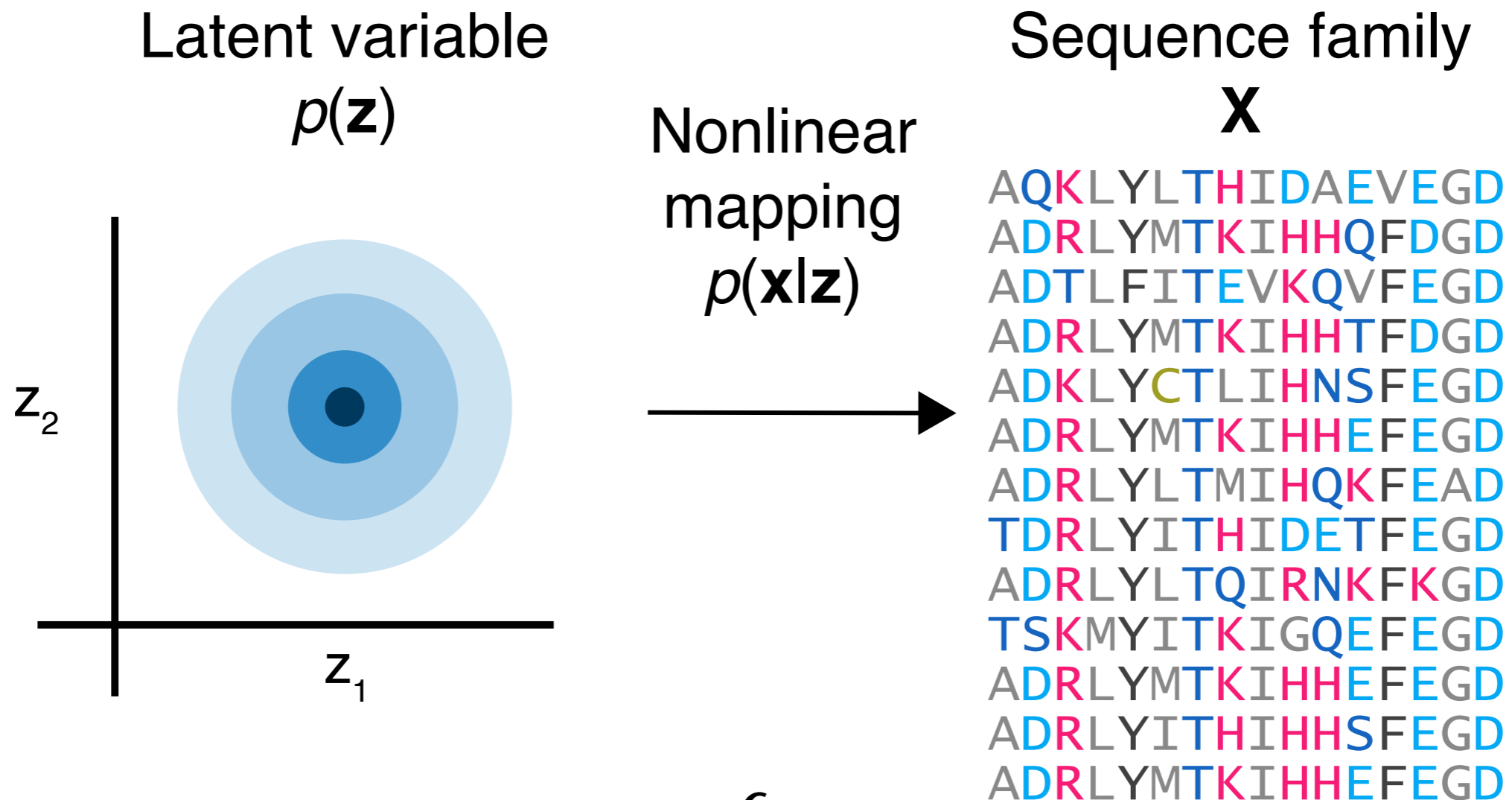
Physiology

Synthetic
biology

Neural networks make powerful latent variable models



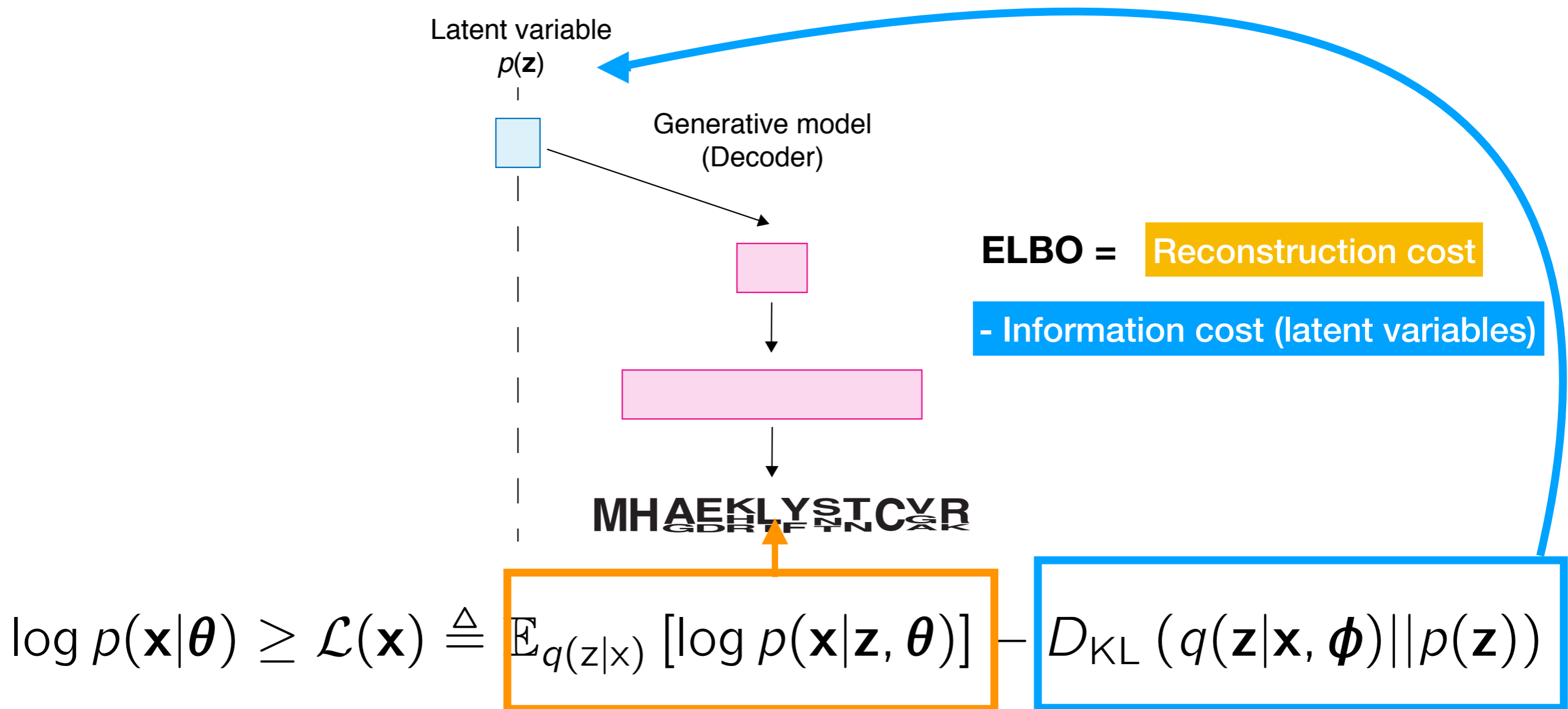
$p(\mathbf{x})$ for latent variable models is generally **intractable**



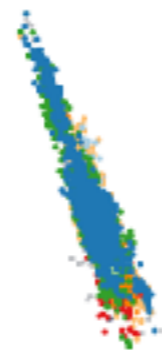
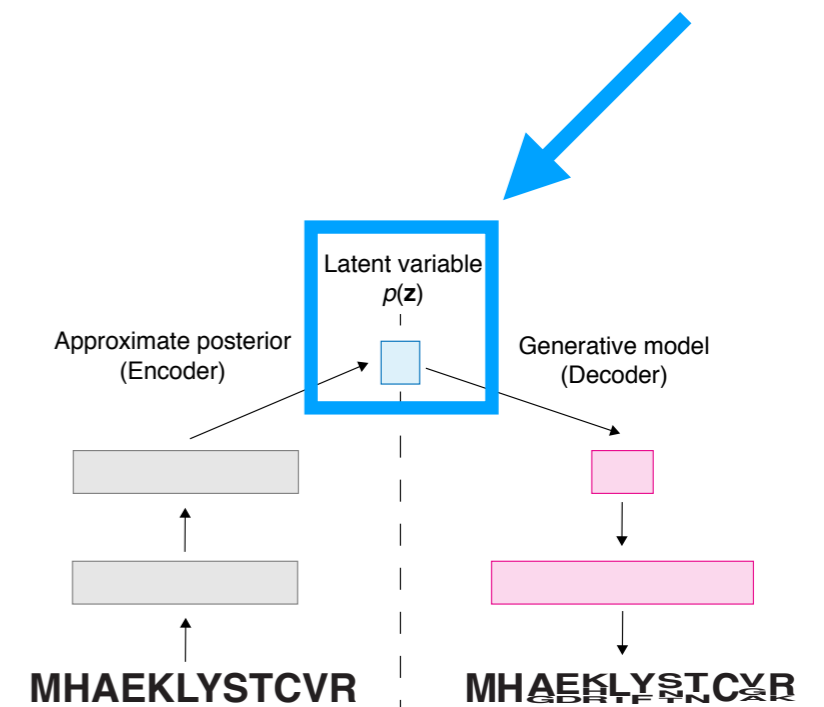
$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z}, \theta) p(\mathbf{z}) d\mathbf{z}$$

have to account for all possible \mathbf{z} for each \mathbf{x}

Variational autoencoders provide a tractable lower bound on $p(\mathbf{x})$



Latent variables are generated for each sequence in alignment

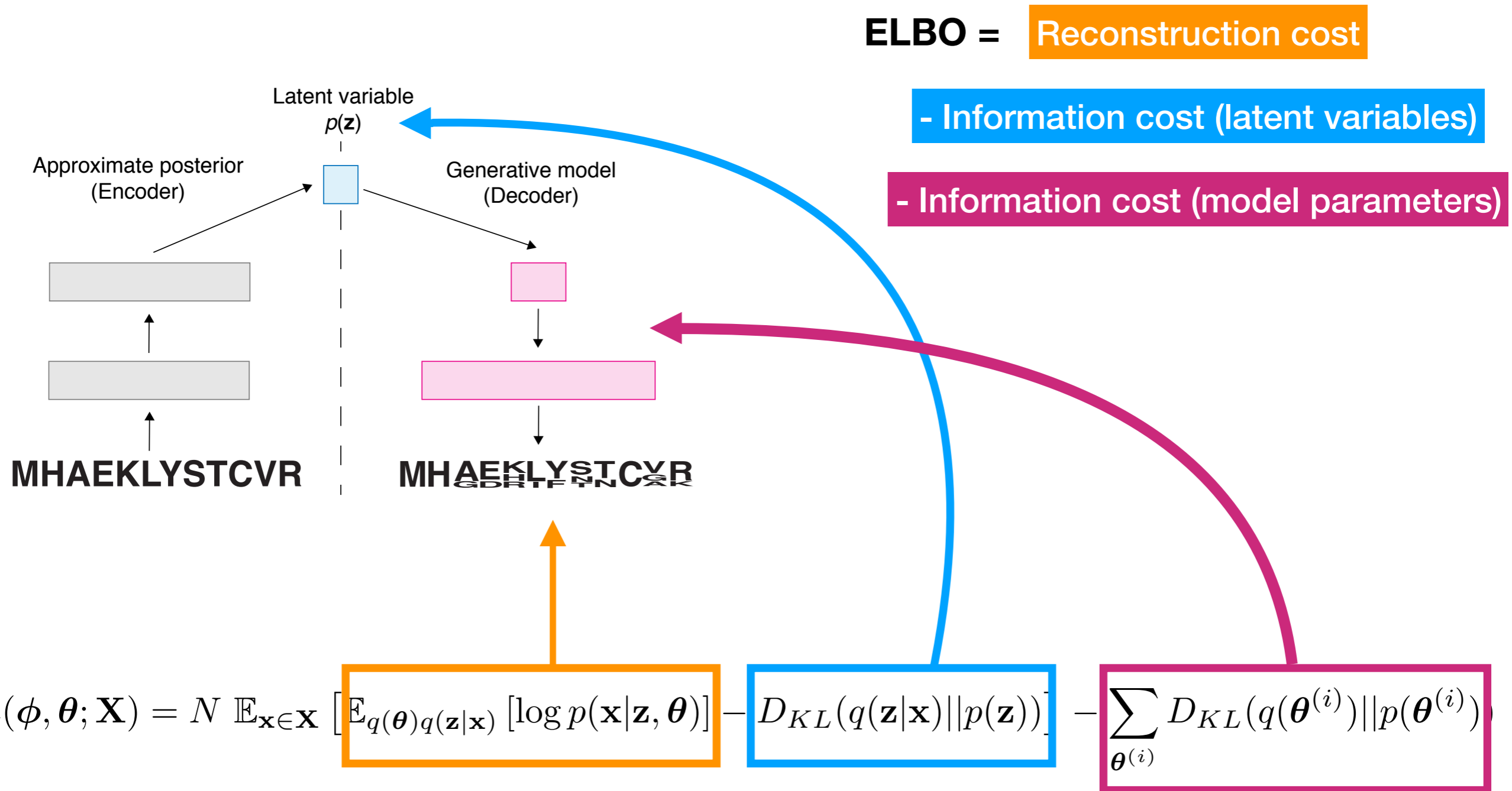


Update 10

β -lactamase sequence family

- Acidobacteria
- Actinobacteria
- Bacteroidetes
- Chloroflexi
- Cyanobacteria
- Deinococcus-Thermus
- Firmicutes
- Fusobacteria
- Proteobacteria
- Other

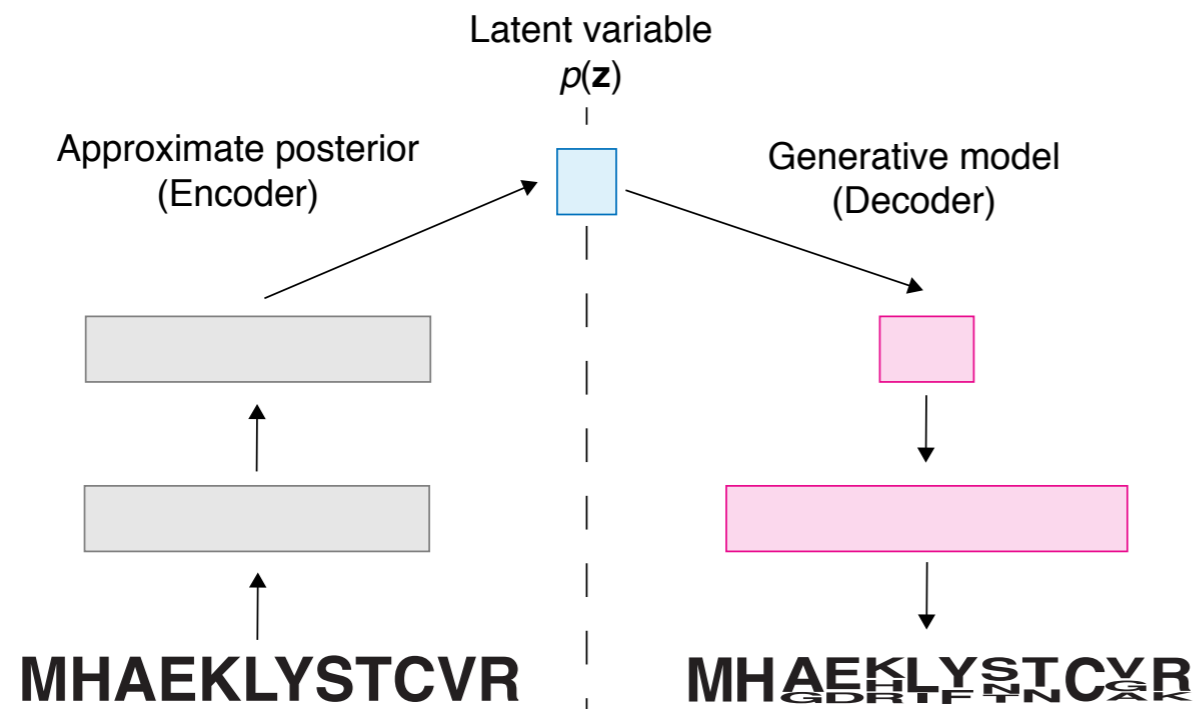
Variational inference on decoder weights prevents overfitting



Mutation prediction with a variational autoencoder

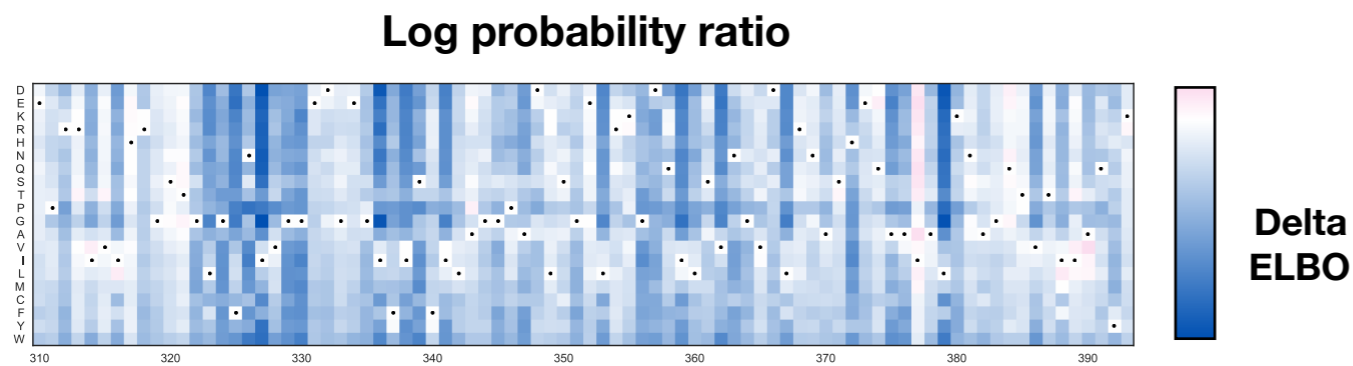
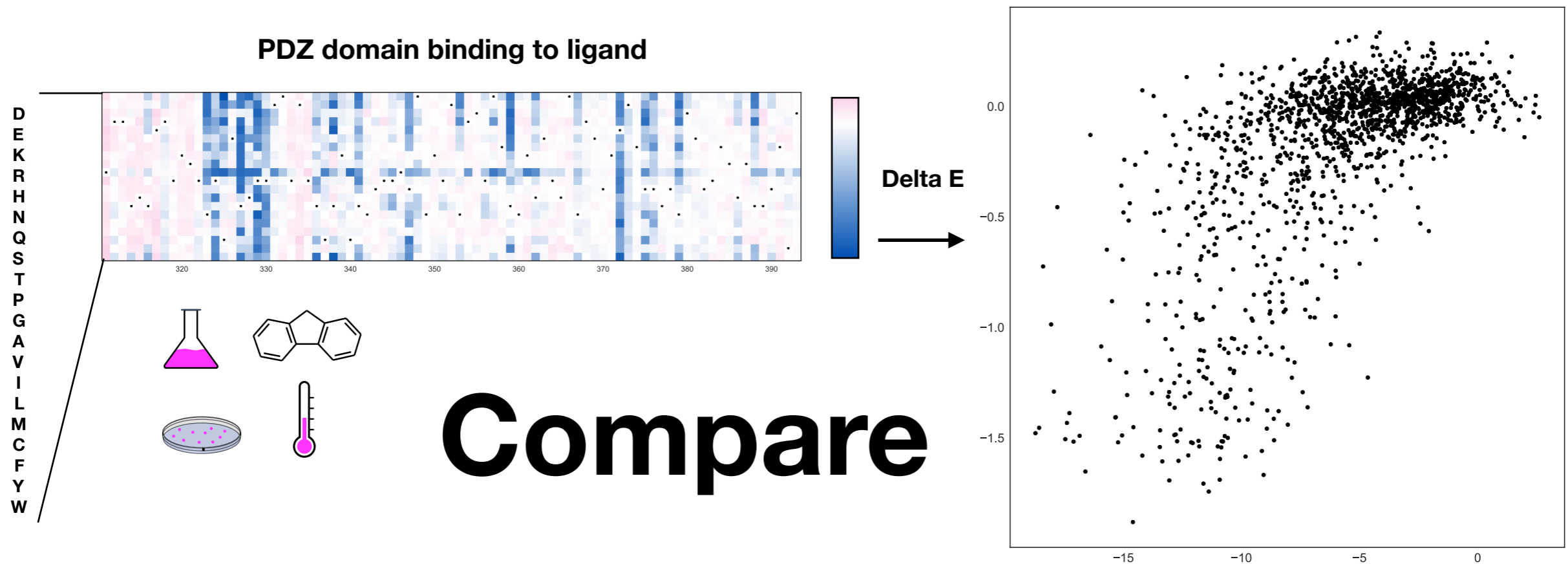
1) Infer a **generative model** of the family

2) Approximate **Log Ratio** with difference in ELBO



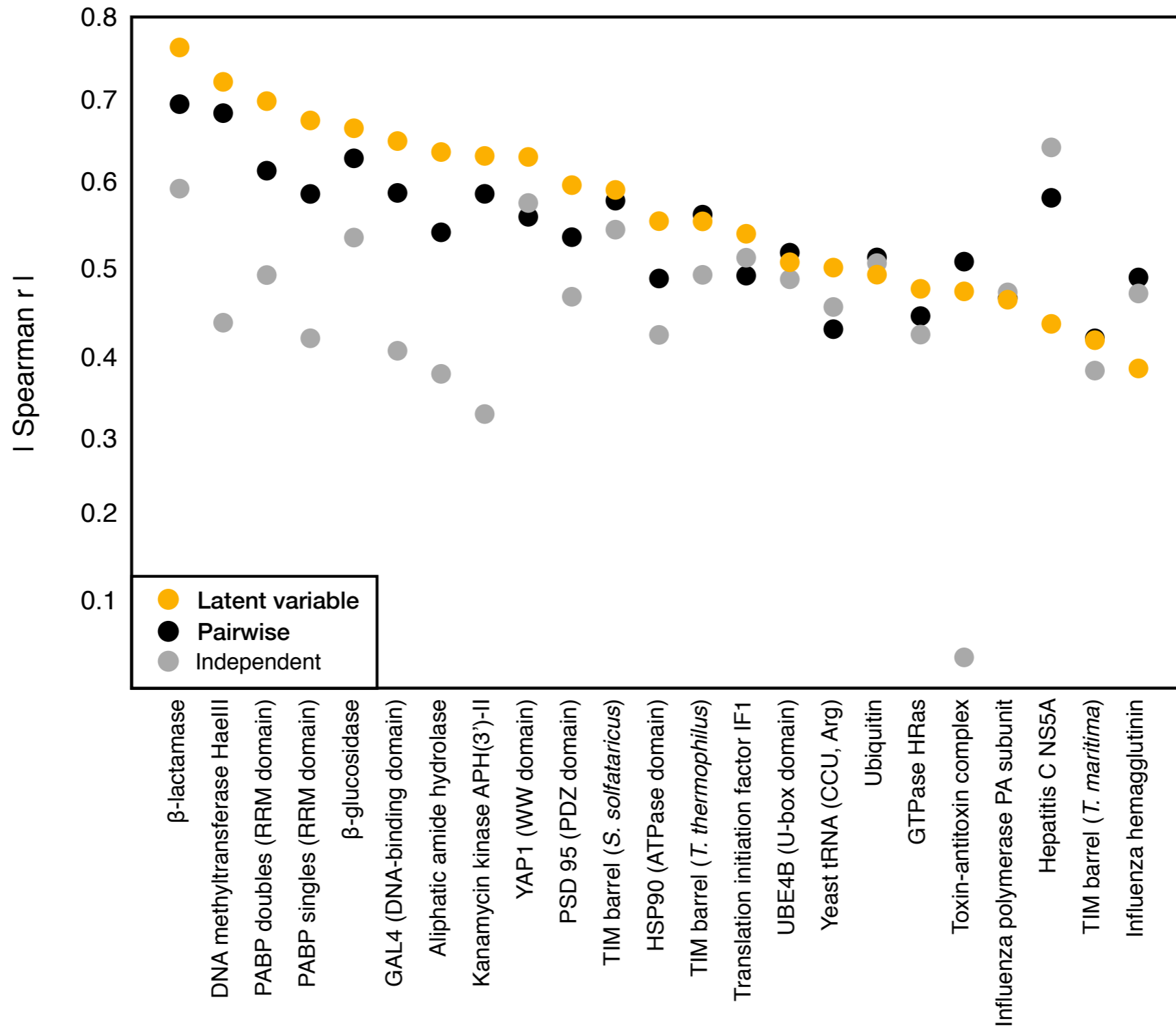
$$\log \frac{p(\mathbf{x}_{\text{mut}} | \boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}} | \boldsymbol{\theta})}$$

To evaluate performance, we collected ~30 saturation mutagenesis experiments

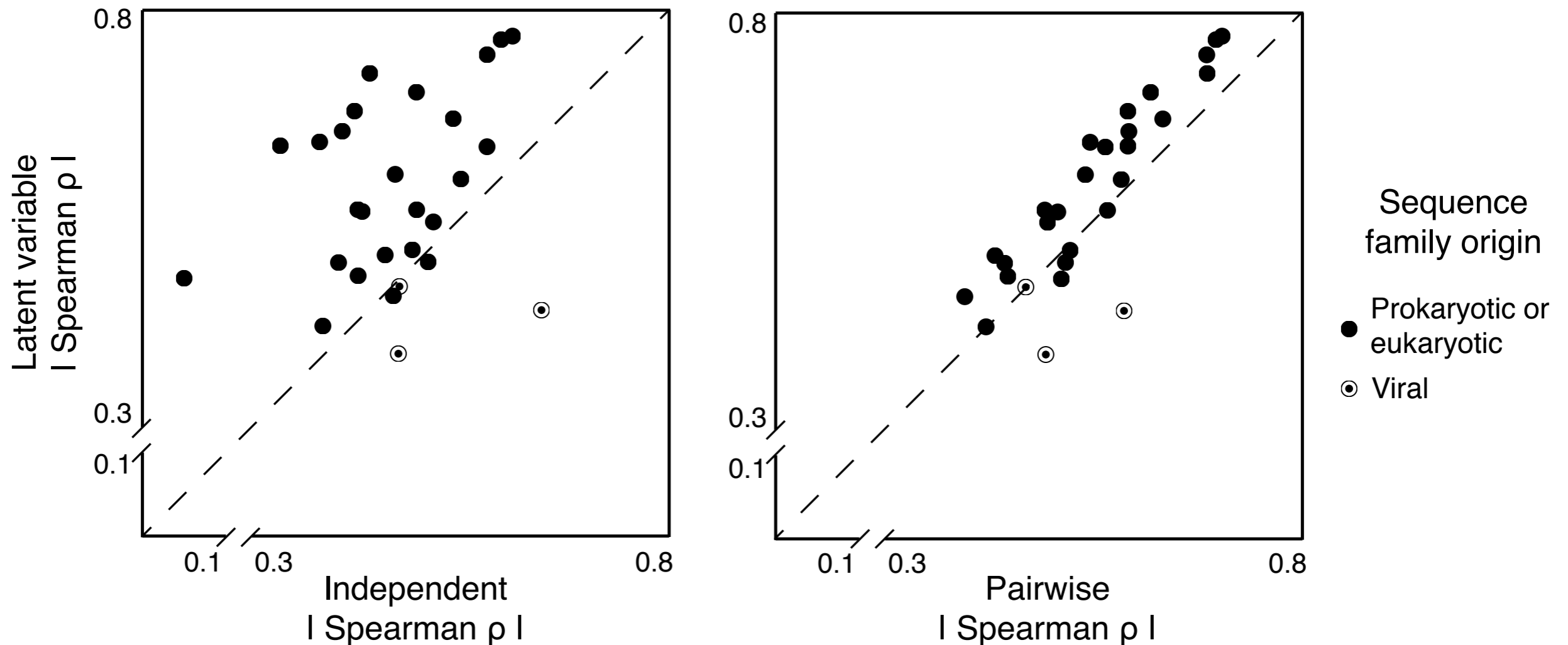


$$\log \frac{p(\mathbf{x}_{\text{mut}}|\boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}}|\boldsymbol{\theta})}$$

Latent variable model is more predictive than pairwise model



Deeper sequence alignments lead to **more predictive** models



Can we **interpret** our model by
building **biology** into the **components**

Encoding **biological knowledge** in a structured matrix with **parameter sharing**

z = latent variable Categorical VAE decoder

h = hidden vector

x = sequence

W, b = weights

$$\mathbf{h} = \text{MLP}(\mathbf{z})$$

$$p(x_i|\mathbf{z}) = \text{Softmax}(\mathbf{W}^{(i)}\mathbf{h} + \mathbf{b}^{(i)})$$

Parameterized by:

$$p(x_i|\mathbf{z}) = \text{Softmax}\left(\mathbf{D}\tilde{\mathbf{W}}^{(i)}\left(\text{Sigmoid}\left(\mathbf{s}^{(i)}\right)\odot\mathbf{h}\right) + \mathbf{b}^{(i)}\right)$$

Dictionary shared across
all positions

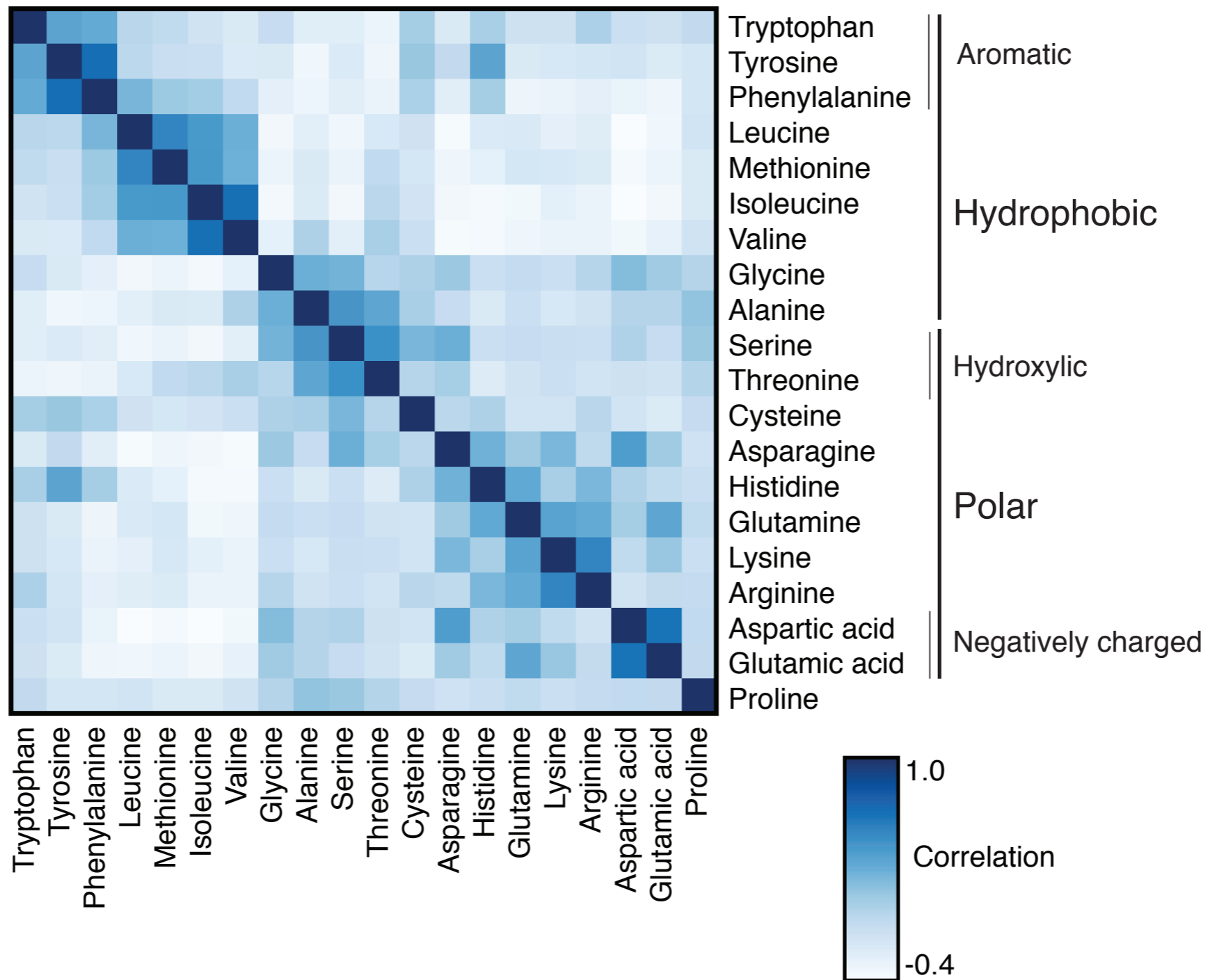
Scale shared across
a position

Biological constraints were included in model parameterization

$$p(x_i|\mathbf{z}) = \text{Softmax} \left(\mathbf{D}\tilde{\mathbf{W}}^{(i)} \left(\text{Sigmoid} \left(\mathbf{s}^{(i)} \right) \odot \mathbf{h} \right) + \mathbf{b}^{(i)} \right)$$

Dictionary shared across all positions

The dictionary encodes **amino acid preferences**



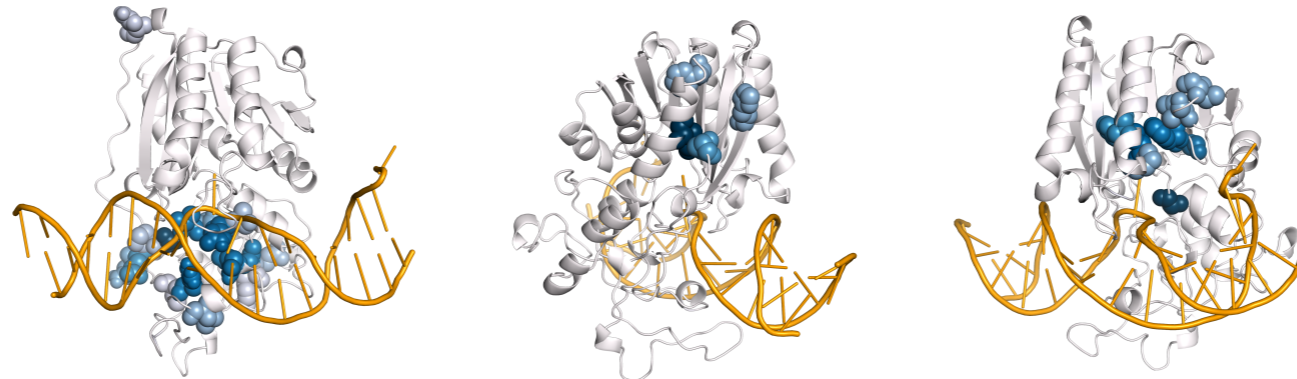
Biological constraints were included in model parameterization

$$p(x_i|\mathbf{z}) = \text{Softmax} \left(\mathbf{D}\tilde{\mathbf{W}}^{(i)} \left(\text{Sigmoid} \left(\mathbf{s}^{(i)} \right) \odot \mathbf{h} \right) + \mathbf{b}^{(i)} \right)$$

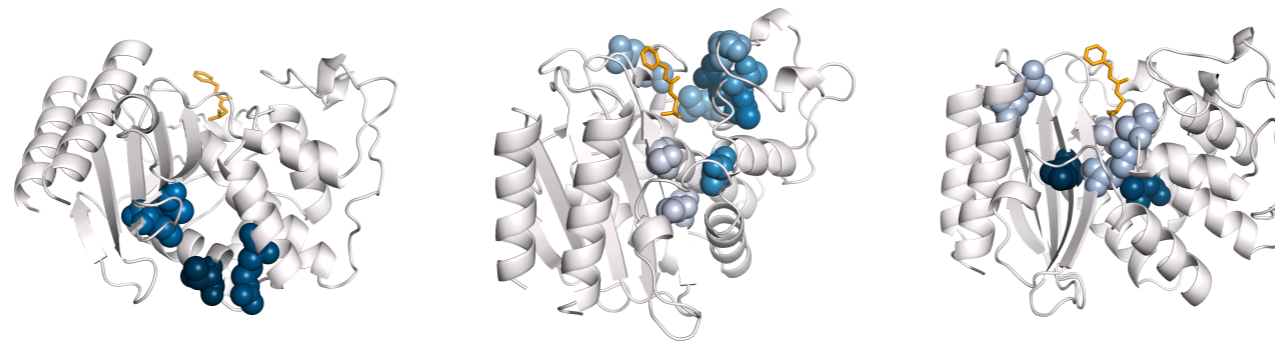
Scale shared across a position

Sparse scale factors are **localized in 3D**

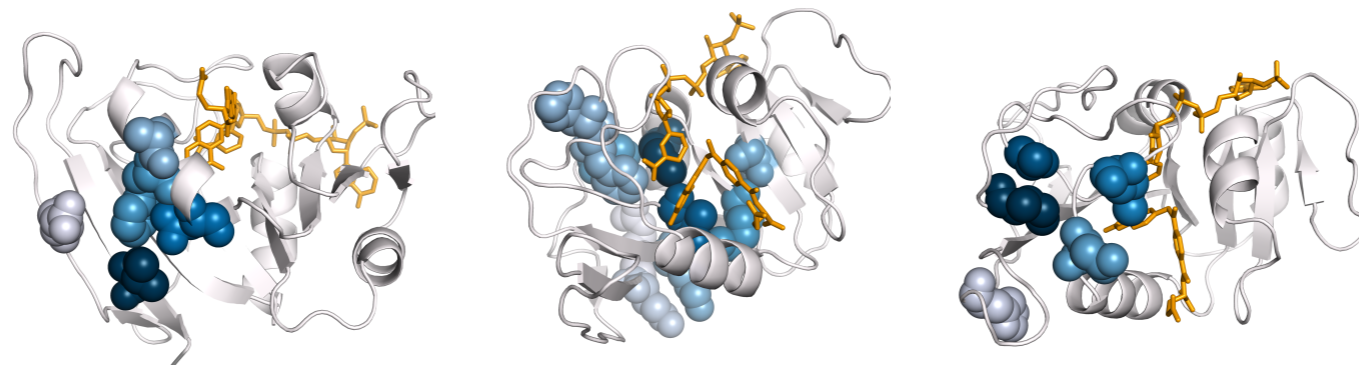
DNA methyltransferase HaeIII



β -lactamase



Dihydrofolate reductase

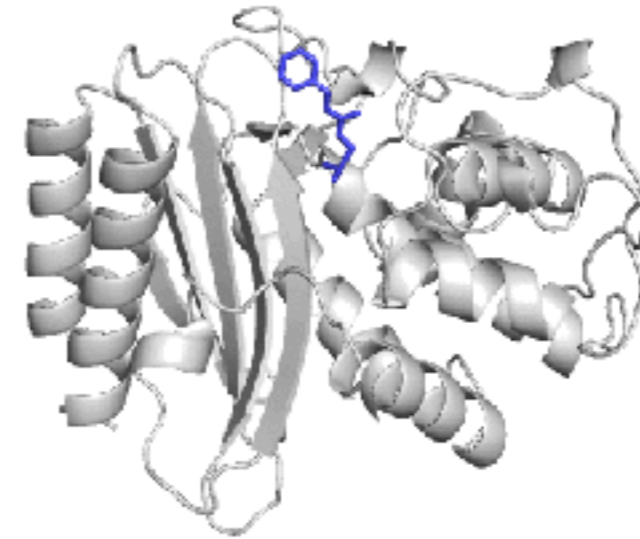


Recap

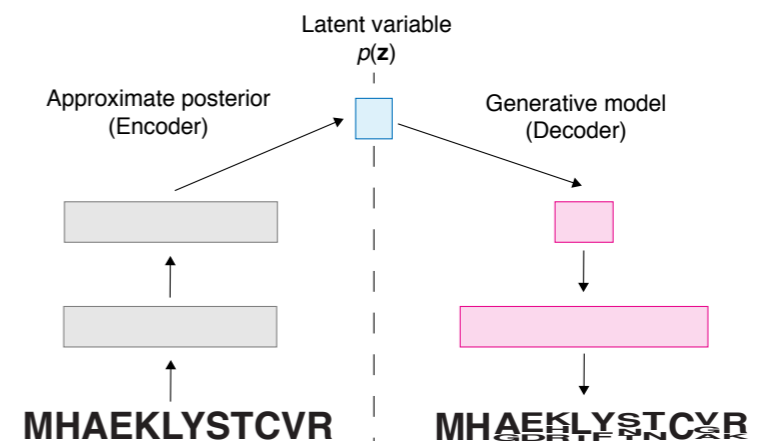
Predicting the effects of mutations is important

Building good generative models of sequence families is useful

Latent variable models predict the effect of mutations better than state-of-art



$$\log \frac{p(\mathbf{x}_{\text{mut}} | \boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}} | \boldsymbol{\theta})}$$



Thank you!



Marks Lab +

Debora Marks
John Ingraham

Chris Sander

DeepSequence github:

<https://github.com/debbiemarkslab/DeepSequence>

EVcouplings python package:

<https://github.com/debbiemarkslab/EVcouplings>

Thomas Hopf
Anna Green
Charlotta Scharfe
Benni Schubert
Eli Weinstein
Kelly Brock
Rohan Maddamsetti
David Ding
June Shin
Hailey Cambra
Agnes Toth-Petroczy
Perry Palmedo
Frank Poelwijk
Nick Gauthier

Jennie Epp

Thank you!