

Computational Systems Biology Deep Learning in the Life Sciences

6.802 6.874 20.390 20.490 HST.506

David Gifford
Lecture 14
April 4, 2019

Machine Learning Designed Therapeutics



<http://mit6874.github.io>

What's on tap today!

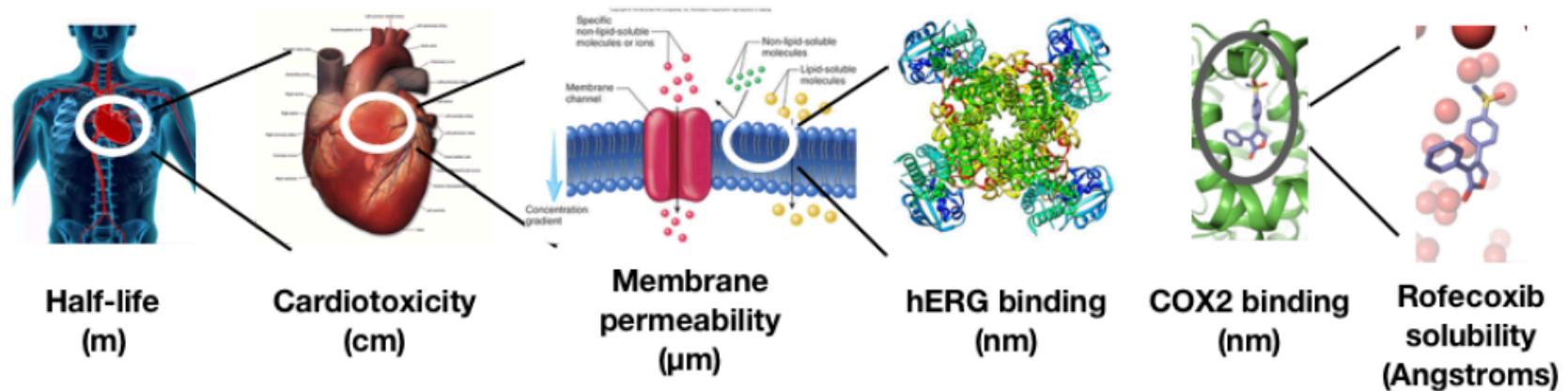
- Characterizing small molecule therapeutics
- Formulating peptide vaccines
- Designing antibodies
- Gain-of-function repairs with a DNA cut

What you should know

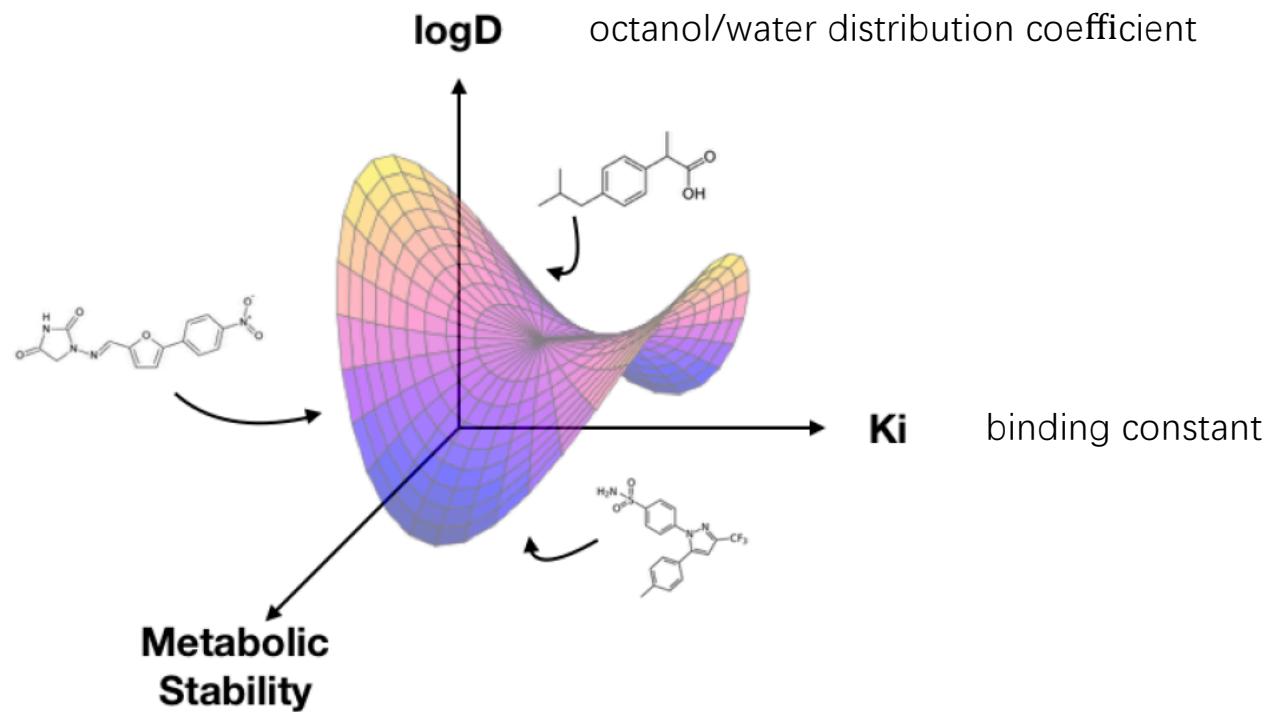
- Graph Convolutional Networks
- Language modes (ElMo)
- How to design novel antibody CDR sequences
- How to predict CRISPR outcomes
- Issues in the representation of small molecules

Graph Convolutional Neural Networks (GCNNs)

There are diverse prediction tasks for small molecules



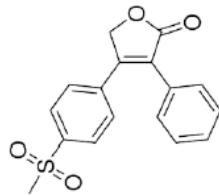
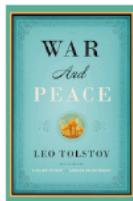
How can we find molecules on the thin manifold of acceptability?



<https://medium.com/@pandelab/step-change-improvement-in-molecular-property-prediction-with-potentialnet-f431ffa32a2c>

How can we create features for small molecules to predict key properties?

Year: 2012



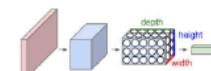
Flat Vector Featurizer

Grid of Pixels →

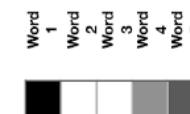
Representation



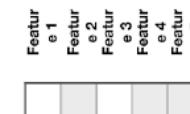
Learning Algorithm



Bag of words →

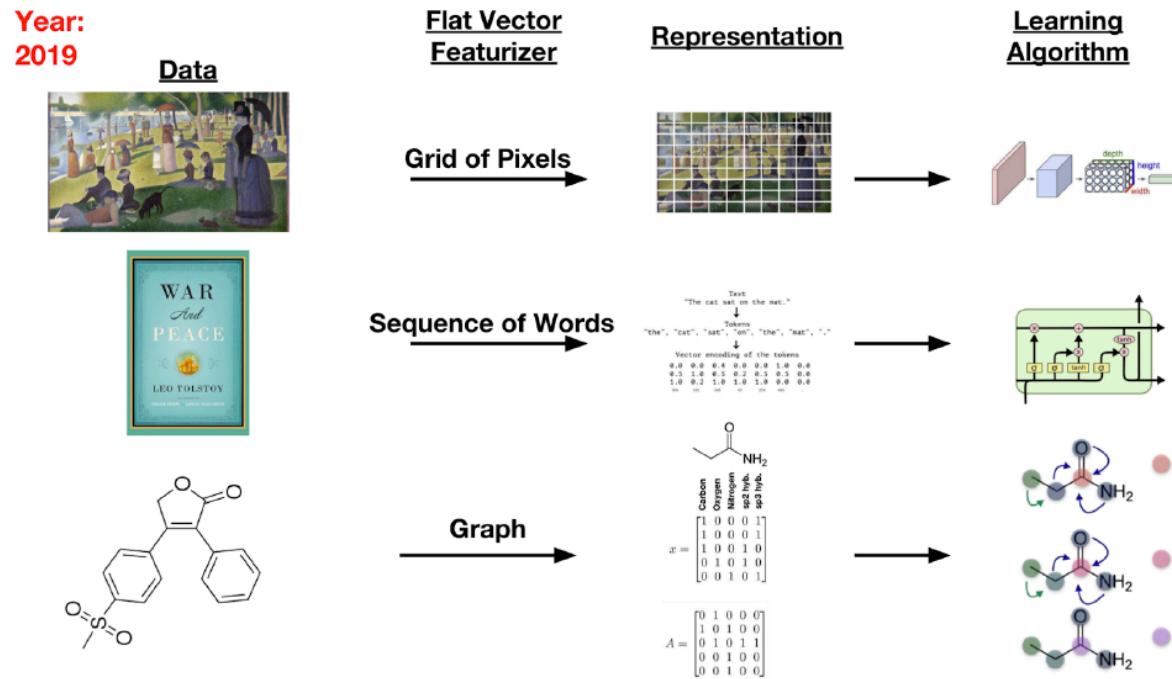


ECFP →



<https://medium.com/@pandelab/step-change-improvement-in-molecular-property-prediction-with-potentialnet-f431ffa32a2c>

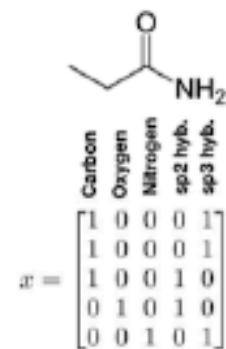
We can use graph convolutions on small molecules by representing their node features and adjacencies



<https://medium.com/@pandelab/step-change-improvement-in-molecular-property-prediction-with-potentialnet-f431ffa32a2c>

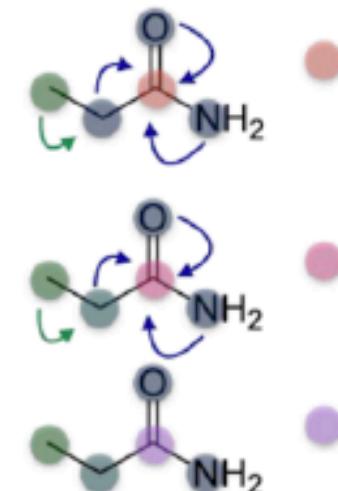
We can use graph convolutions on small molecules by representing their node features and adjacencies

X is per-atom features of each graph node as rows

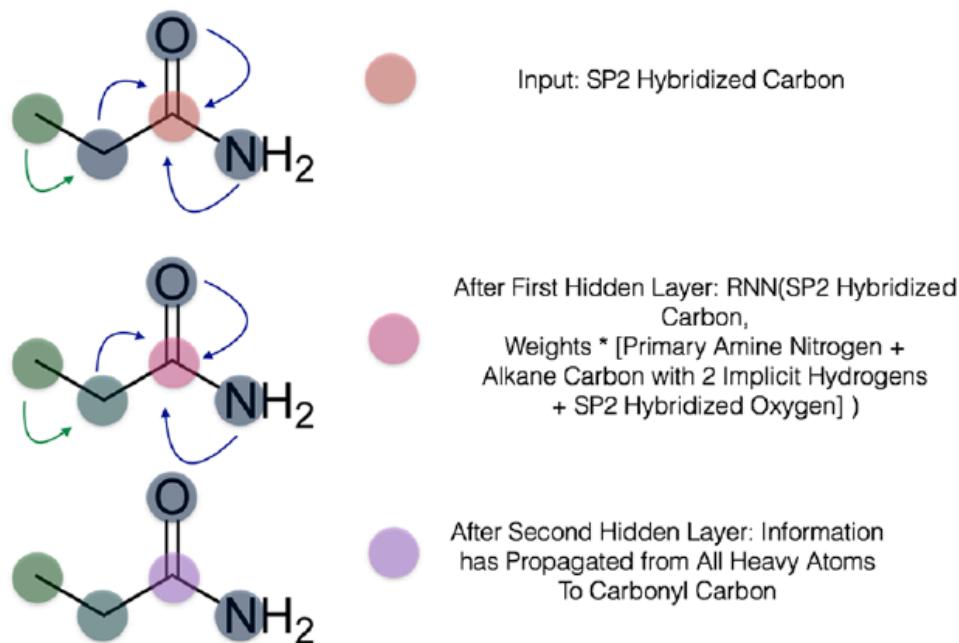


A is the per-atom adjacency matrix as rows

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$



Information propagates along adjacencies one step at a time



**Multi-Layer Perceptron /
Fully Connected Neural Network**

X is per-atom features of each graph node as rows

A is the per-atom adjacency matrix as rows

K is the number of graph convolutional layers

x is flat features of molecule

$$h^{(1)} = \text{ReLU} \left(W^{(1)} \cdot x \right)$$

$$h^{(2)} = \text{ReLU} \left(W^{(2)} \cdot h^{(1)} \right)$$

⋮

$$h^{(K)} = \text{ReLU} \left(W^{(K)} \cdot h^{(K-1)} \right)$$

**Graph Convolutional
Neural Network (GCNN)**

$$H^{(1)} = \text{ReLU} \left(W^{(1)} \cdot A \cdot X \right)$$

$$H^{(2)} = \text{ReLU} \left(W^{(2)} \cdot A \cdot H^{(1)} \right)$$

⋮

$$H^{(K)} = \text{ReLU} \left(W^{(K)} \cdot A \cdot H^{(K-1)} \right)$$

$$x^{(NN)} = \sum_{atoms} H^{(K)}$$

$$h^{(1)} = \text{ReLU} \left(W^{(1)} \cdot x^{(NN)} \right)$$

$$h^{(2)} = \text{ReLU} \left(W^{(2)} \cdot h^{(1)} \right)$$

⋮

$$h^{(K)} = \text{ReLU} \left(W^{(K)} \cdot h^{(K-1)} \right)$$

X is per-atom features of each graph node

A is the atom adjacency matrix

H is the feature map for all atoms

1 bond length away

2 bond lengths away

K bond lengths away

Sum over per atom features

$x^{(NN)}$ is convolutional “fingerprint” of entire molecule

K fully connected layers

Graph Convolutional Neural Network (GCNN)

$$H^{(1)} = \text{ReLU} \left(W^{(1)} \cdot A \cdot X \right)$$

$$H^{(2)} = \text{ReLU} \left(W^{(2)} \cdot A \cdot H^{(1)} \right)$$

⋮

$$H^{(K)} = \text{ReLU} \left(W^{(K)} \cdot A \cdot H^{(K-1)} \right)$$

$$x^{(NN)} = \sum_{atoms} H^{(K)}$$

$$h^{(1)} = \text{ReLU} \left(W^{(1)} \cdot x^{(NN)} \right)$$

$$h^{(2)} = \text{ReLU} \left(W^{(2)} \cdot h^{(1)} \right)$$

⋮

$$h^{(K)} = \text{ReLU} \left(W^{(K)} \cdot h^{(K-1)} \right)$$

PotentialNet

h_i is the feature map for atom i

1 bond length away; edge type specific processing;
Gated Recurrent Unit enables selecting adding

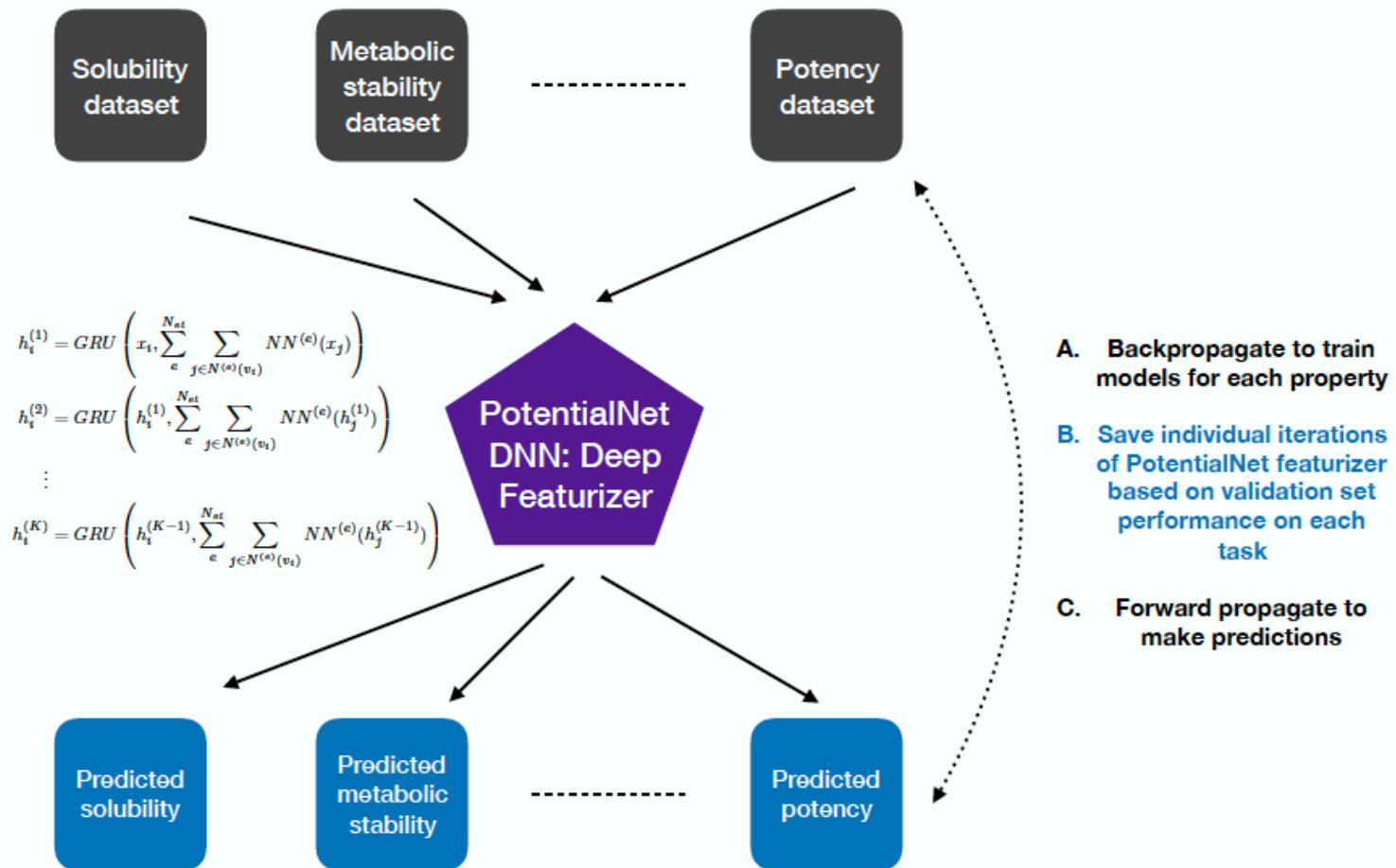
K bond length away; edge type specific processing

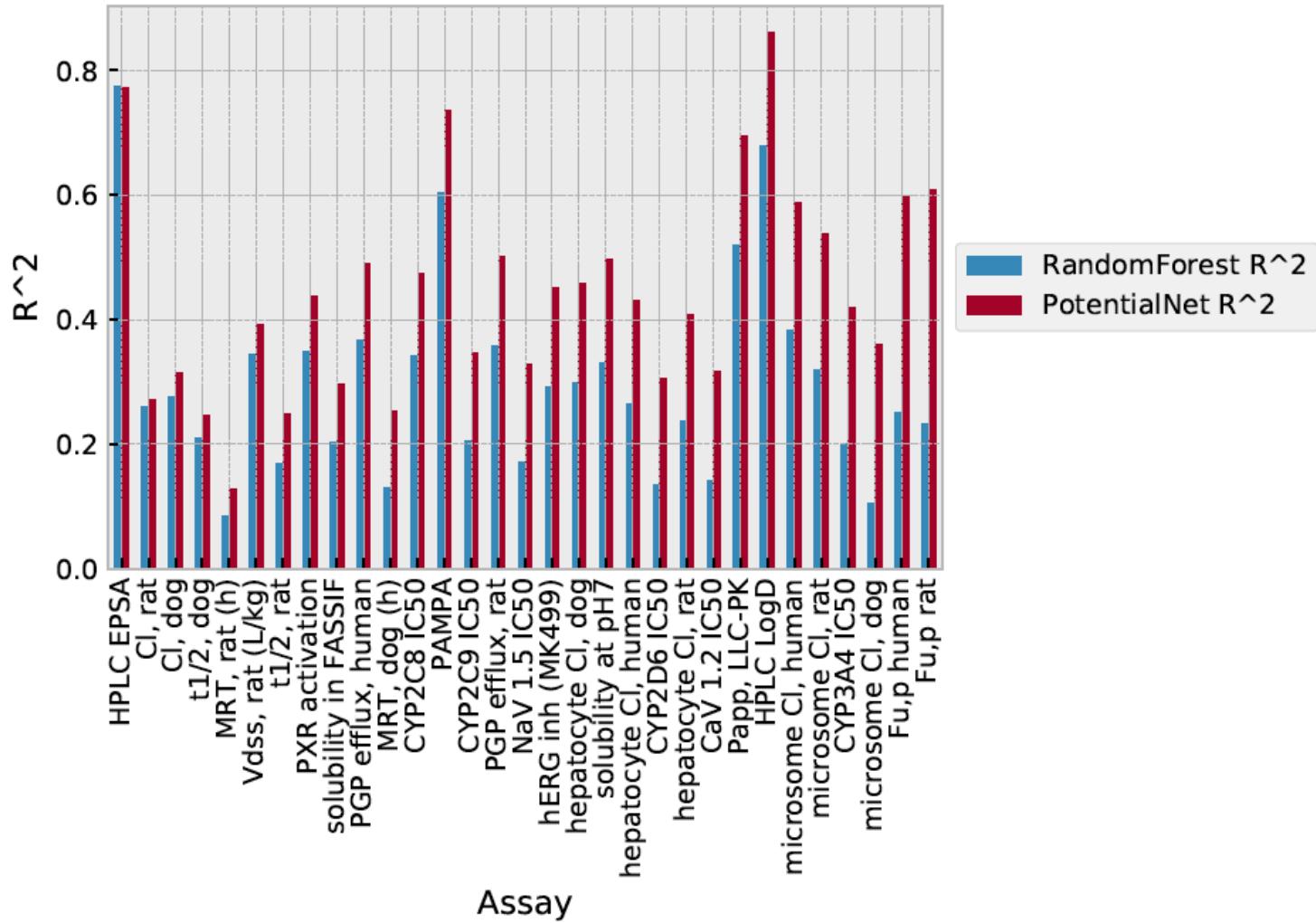
$h^{(NN)}$ is convolutional “fingerprint” of entire molecule

Fully connected layers

$$\begin{aligned}
 h_i^{(1)} &= GRU \left(x_i, \sum_e^{N_{\text{et}}} \sum_{j \in N^{(e)}(v_i)} NN^{(e)}(x_j) \right) \\
 &\vdots \\
 h_i^{(K)} &= GRU \left(h_i^{(b_{K-1})}, \sum_e^{N_{\text{et}}} \sum_{j \in N^{(e)}(v_i)} NN^{(e)}(h_j^{(b_{K-1})}) \right) \\
 h^{(NN)} &= \sigma \left(i(h^{(K)}, x) \right) \odot \left(j(h^{(K)}) \right) \\
 h^{(FC_0)} &= \sum_{j=1}^{N_{\text{Lig}}} h_j^{(NN)} \\
 h^{(FC_1)} &= \text{ReLU} \left(W^{(FC_1)} h^{(FC_0)} \right) \\
 &\vdots \\
 h^{(FC_K)} &= W^{(FC_K)} h^{(FC_{K-1})}
 \end{aligned}$$

Training and Prediction

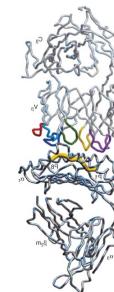




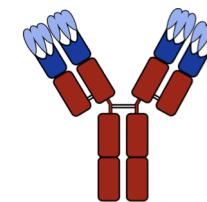
Computational immunotherapeutics

We wish to both identify therapeutic targets and therapeutic molecules with the help of machine learning

- Identifying target peptide-MHC molecules

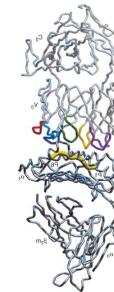


- Designing antibody Complementarity Determining Regions (CDRs)

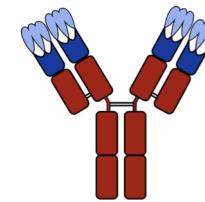


We wish to both identify therapeutic targets and therapeutic molecules with the help of machine learning

- Identifying target peptide-MHC molecules

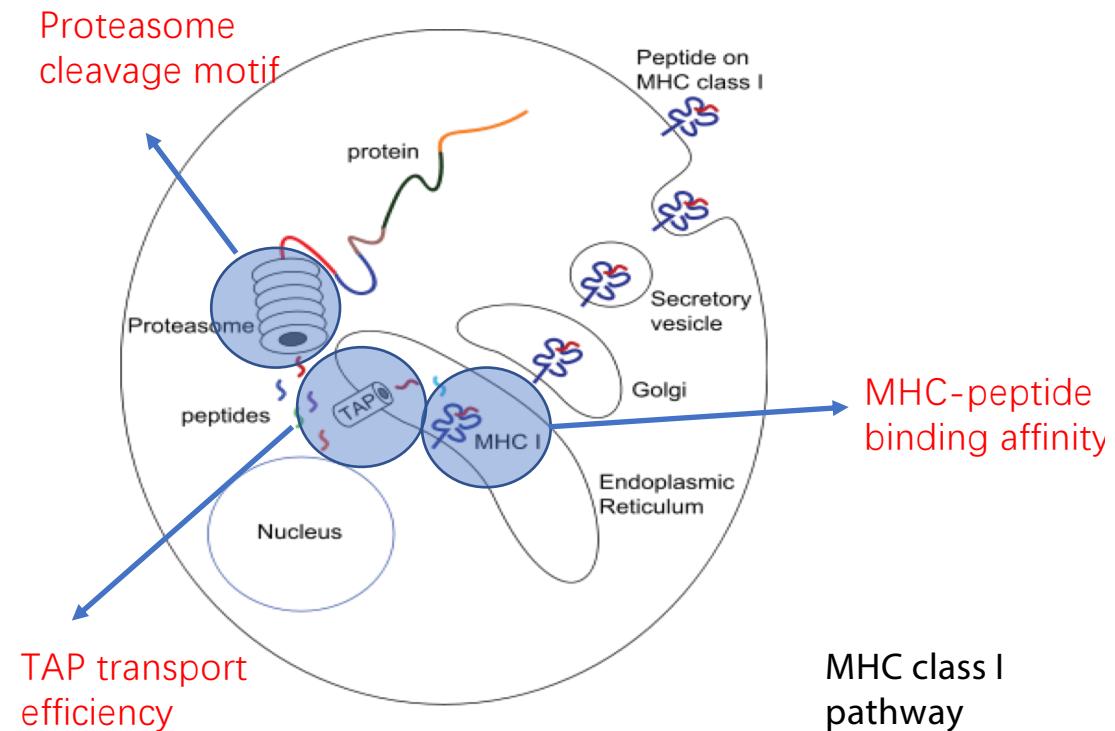


- Designing antibody Complementarity Determining Regions (CDRs)

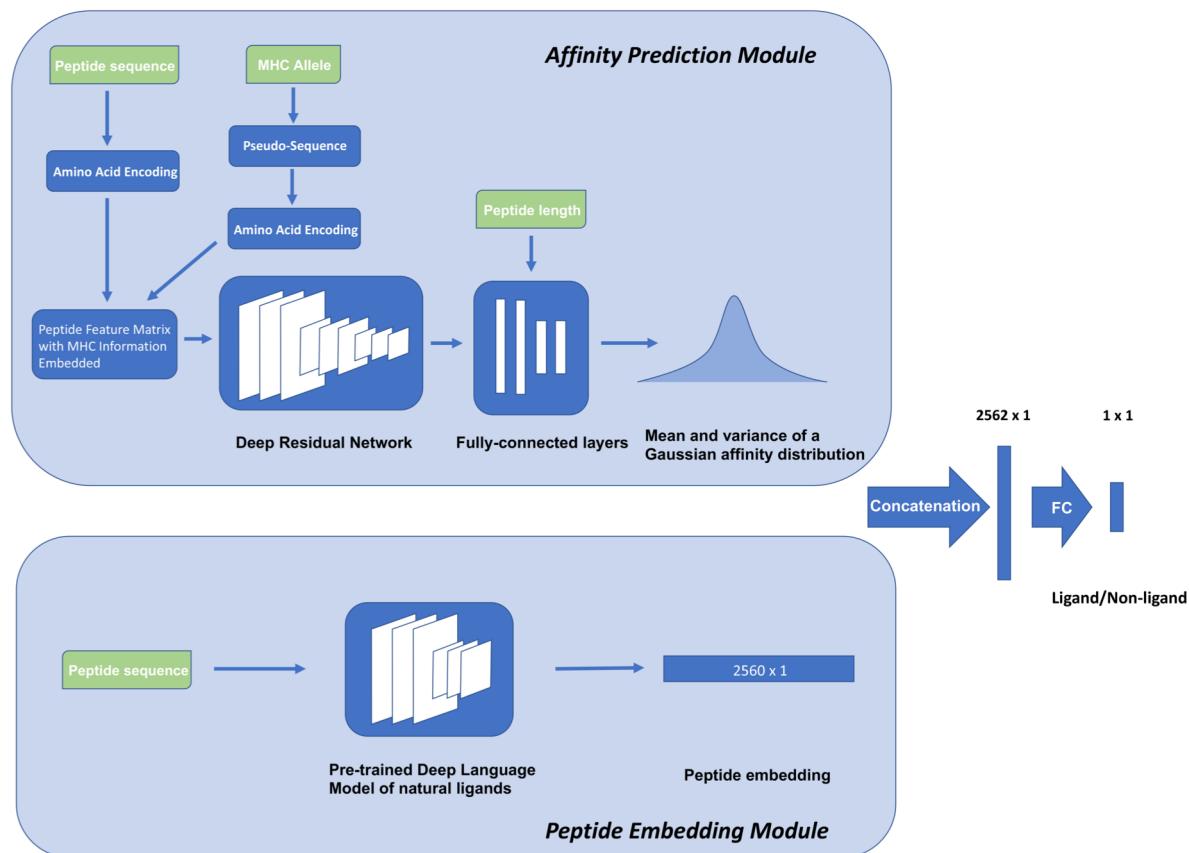


Most existing methods for peptide display focus on modeling MHC binding affinity

- Immune Epitope Database (IEDB) contains a large collection of binding affinity datasets curated from literature
- However, models trained on affinity data are not able to consider other factors in MHC ligand selection



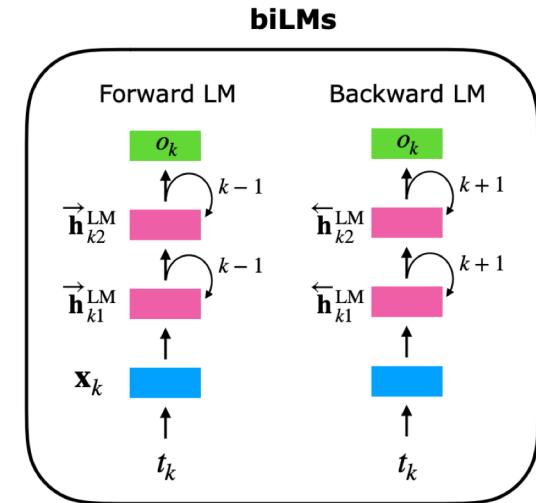
DeepLigand predicts MHC class I peptide presentation (552,252 positive examples, 2.5M negative, 192 MHC alleles)



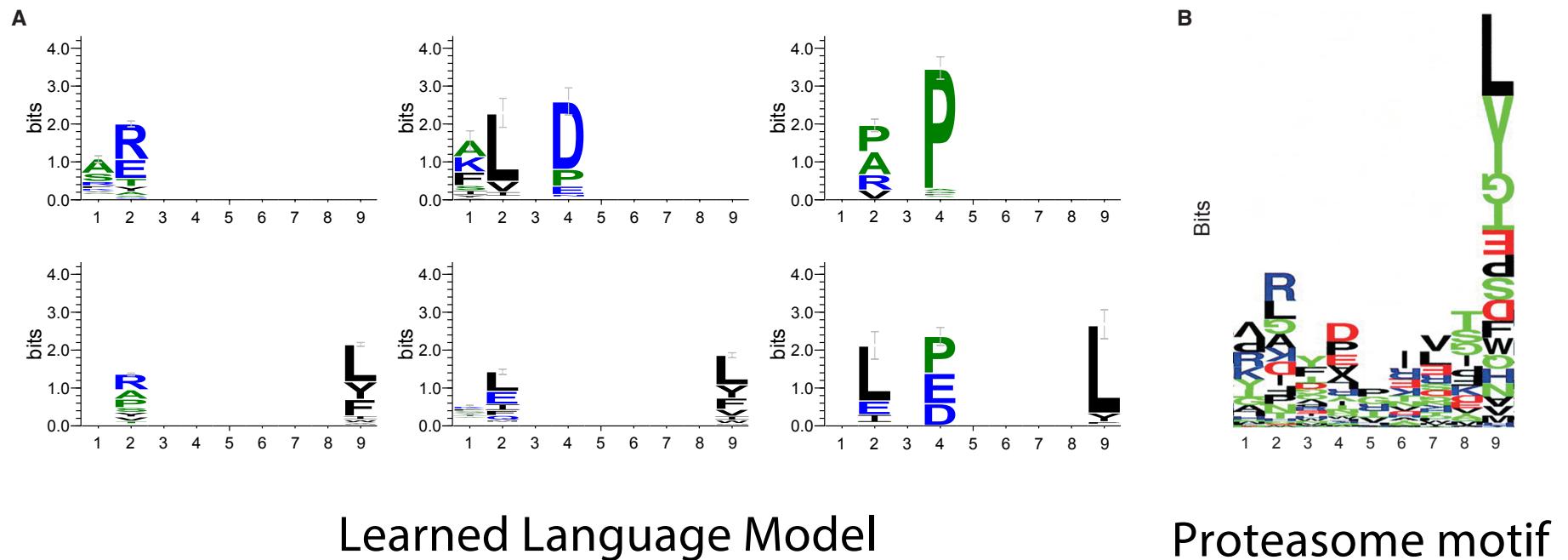
ELMo for learning contextualized word embedding

ELMo represents a word t_k as a linear combination of corresponding hidden layers (inc. its embedding)

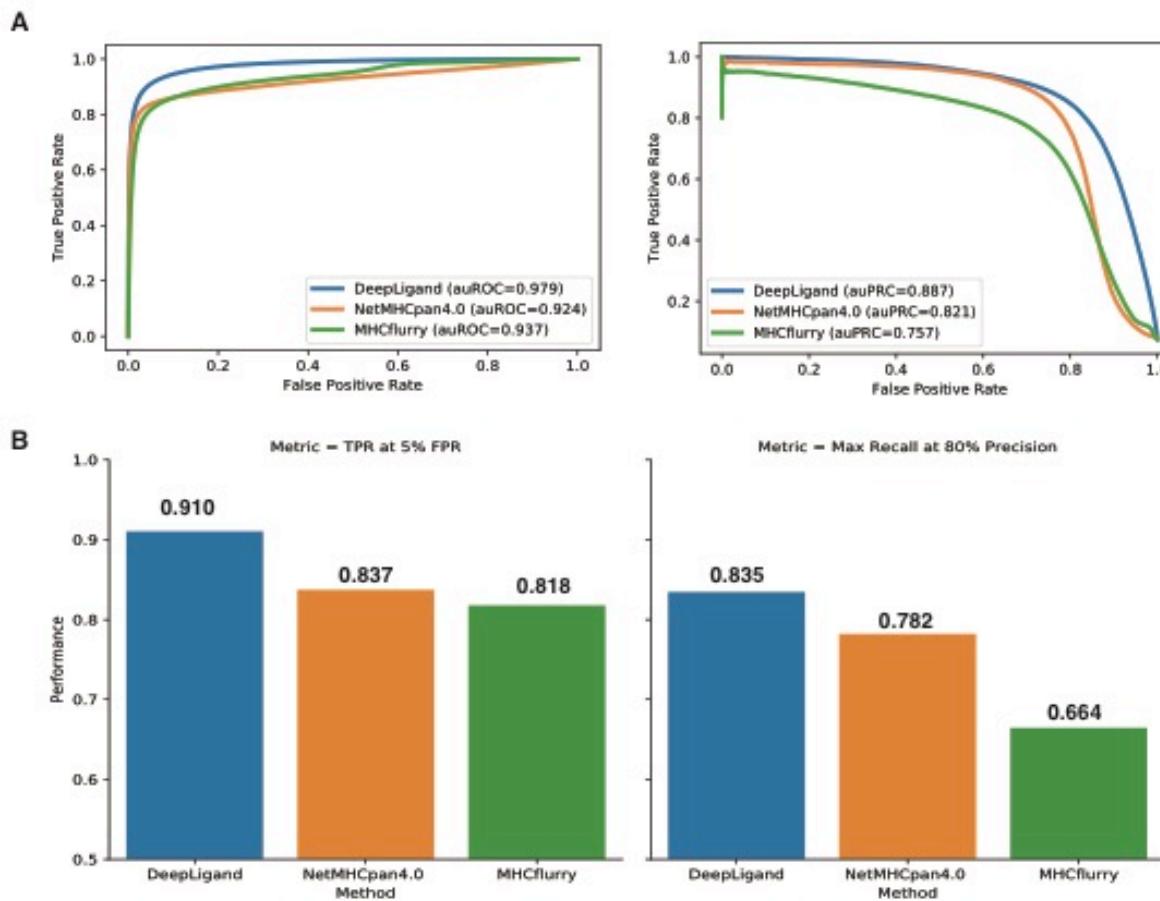
$$\text{ELMo}_k^{\text{task}} = \gamma^{\text{task}} \times \sum \left\{ \begin{array}{l} s_2^{\text{task}} \times \mathbf{h}_{k2}^{\text{LM}} \quad \text{pink} | \text{pink} \\ s_1^{\text{task}} \times \mathbf{h}_{k1}^{\text{LM}} \quad \text{pink} | \text{pink} \\ s_0^{\text{task}} \times \mathbf{h}_{k0}^{\text{LM}} \quad \text{blue} | \text{blue} \\ ([\mathbf{x}_k; \mathbf{x}_k]) \end{array} \right. \begin{array}{l} \text{Concatenate hidden layers} \\ \leftarrow \\ [\vec{\mathbf{h}}_{kj}^{\text{LM}}; \overleftarrow{\mathbf{h}}_{kj}^{\text{LM}}] \end{array}$$



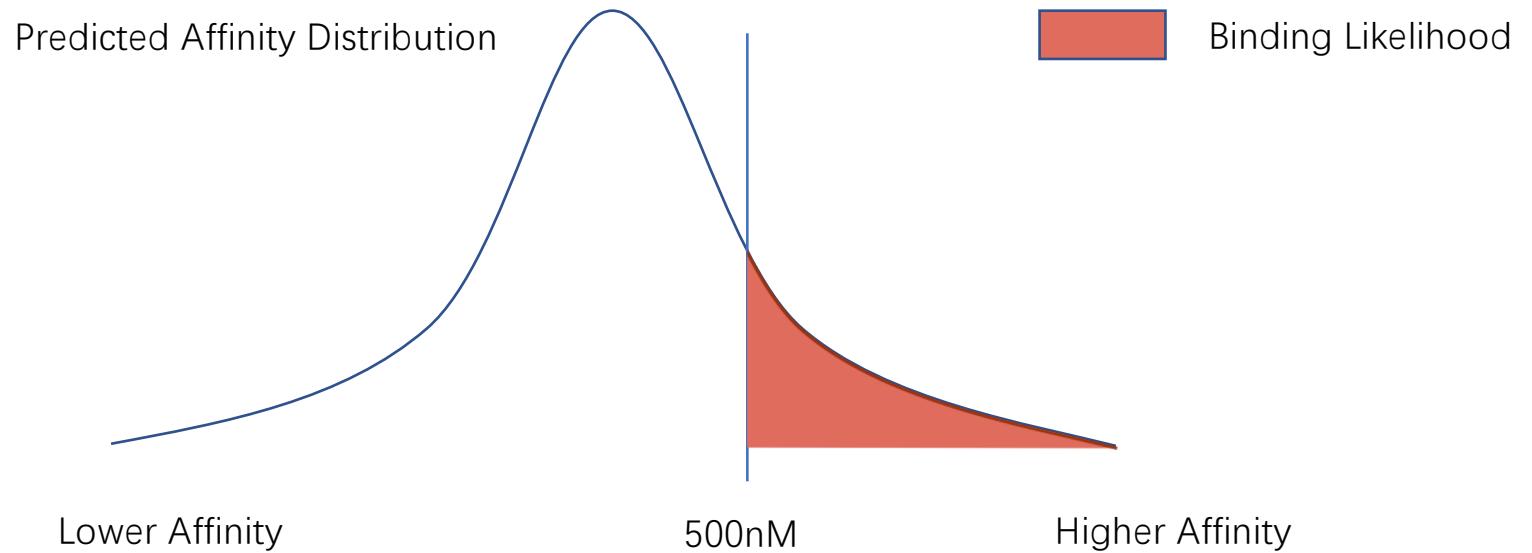
Class I learned language model is consistent with the known proteasome cleavage motif



DeepLigand outperforms existing methods (Class I)

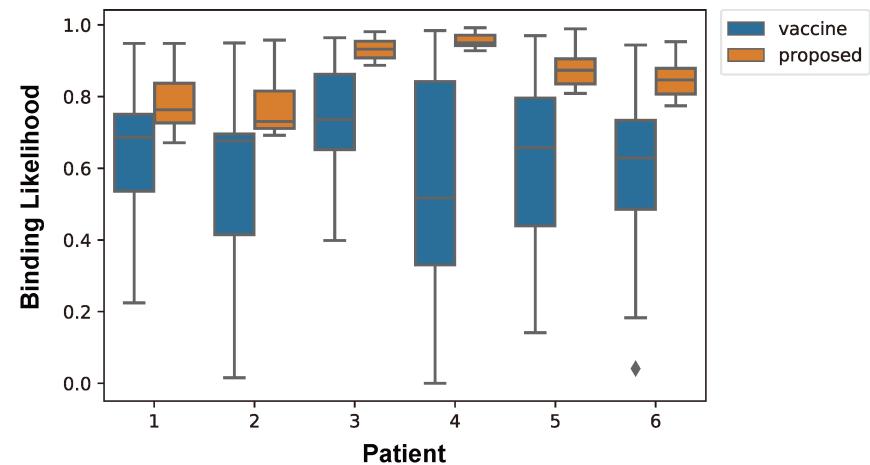
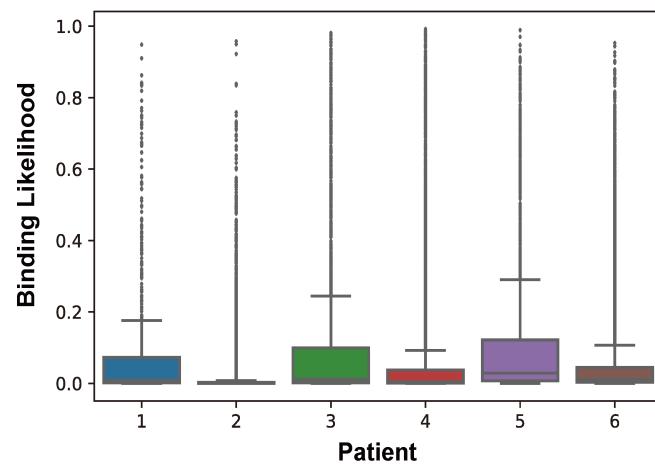


MHC class I uncertainty metrics enable binding likelihood for peptide-MHC molecules



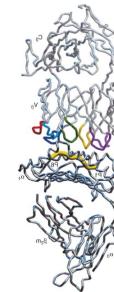
Published MHC class I peptide vaccine formulations can be improved by uncertainty metrics

- We examined the binding likelihood of neo-antigen peptides designed by Otta et al (2017) as well as all the mutation-spanning peptides in each patient

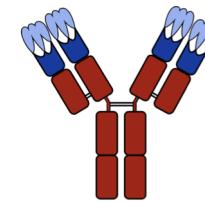


We wish to both identify therapeutic targets and therapeutic molecules with the help of machine learning

- Identifying target peptide-MHC molecules

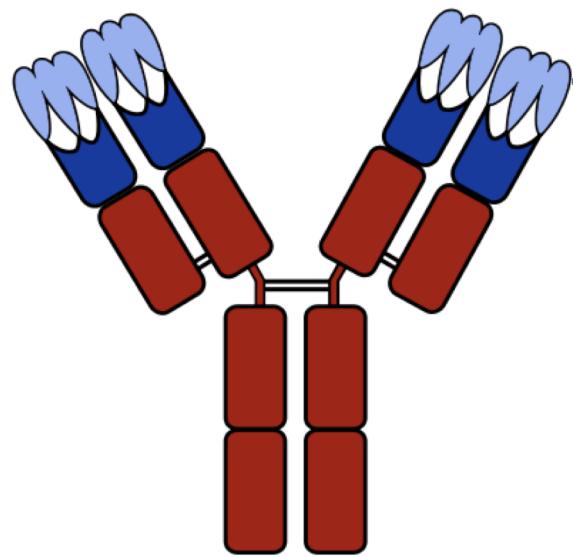


- Designing antibody Complementarity Determining Regions (CDRs)

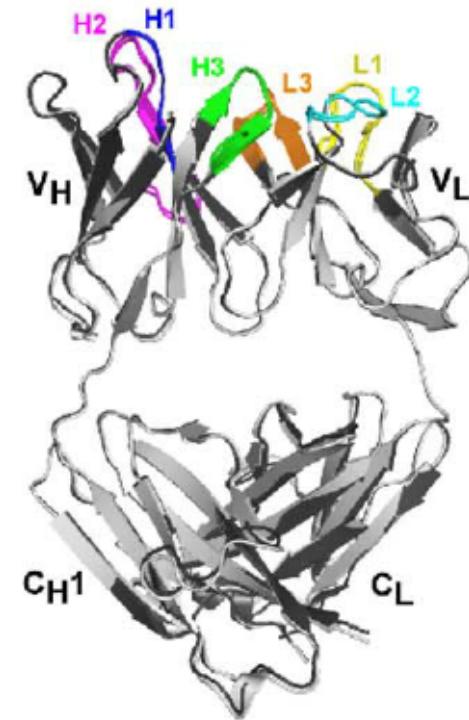


How can we design CDRs that meet multiple objectives?

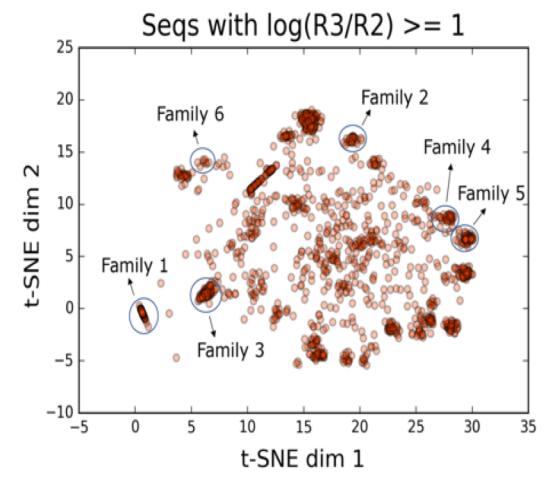
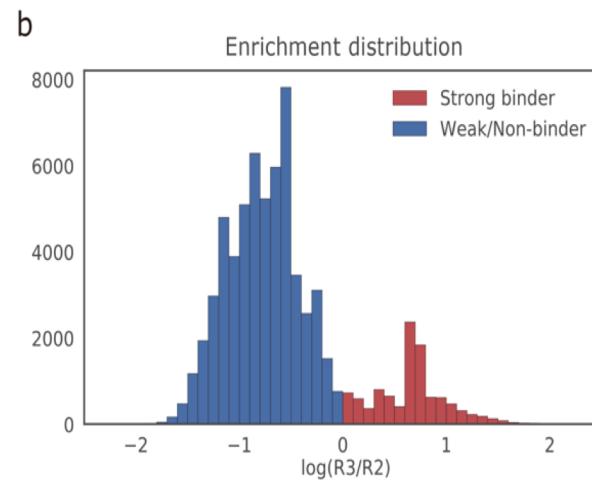
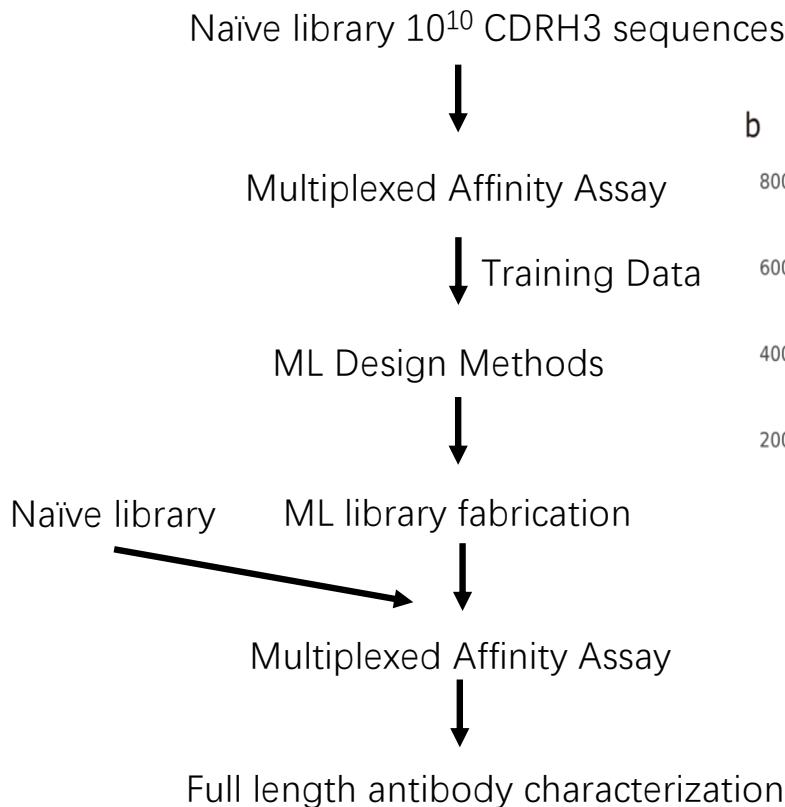
Complementarity-determining regions (CDRs) largely determine target affinity



Six CDRs in total
for each Fab

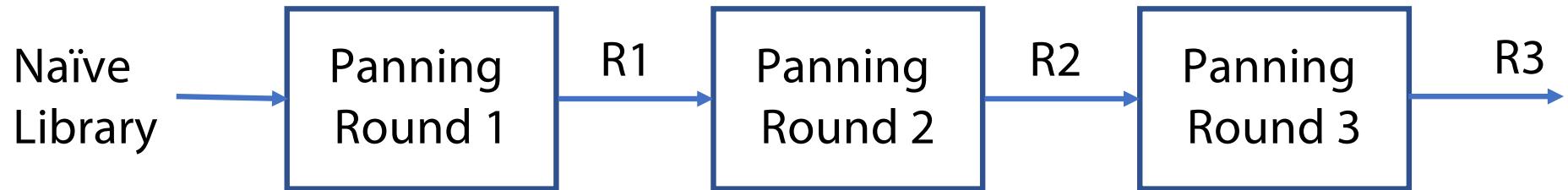


Design flow



Training Data

Enrichment is defined by the output of three panning rounds



$$\text{Enrichment}_{R3/R1} = \log_{10} (\text{Frequency } R3 / \text{Frequency } R1)$$

$$\text{Enrichment}_{R3/R2} = \log_{10} (\text{Frequency } R3 / \text{Frequency } R2)$$

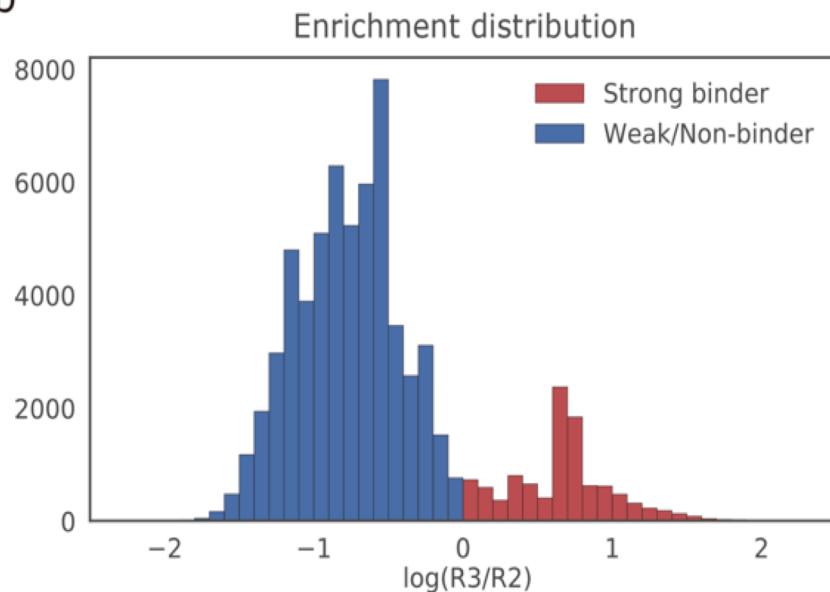
Train on CDR-H3 sequence and enrichment

Sequence	Log(R3/R2)	Enbrel_a	Avastin_a	Herceptin_a
ADGAFDAYMDY	-0.9561	-0.5989	-0.9730	-1.2414
ADGYRVYYYAMDY	1.2253	1.4830	0.9872	1.1175
ADRRPPLIFFDY	0.8519	0.8458	1.9072	1.9057
ADWLSLLYRFDY	-0.4779	-0.9202	-0.7767	-0.8339
AEHVAYHPRYSFDY	-0.9474	-0.7291	-0.9730	-0.8649
AGRYWWLLDY	0.3242	0.3843	1.7872	0.6588
AGYHQTWPYGLDY	1.0482	0.8792	0.9135	-0.2221
AKRRRQYVYHPIYFDY	1.6727	1.4852	1.9769	2.0698
AKYADTYGLDY	0.4839	0.2024	-0.2996	0.9655
AKYGSYYGFDY	0.5650	0.3526	0.3929	0.5801
DAYPGWDLWPDPYPDFY	0.2757	-0.0151	-0.0842	0.4879
DDIHHLLYYFDY	0.9610	1.1010	0.9135	1.5183
DDQYVGFYGEGGLDY	-0.2620	0.0897	0.3372	0.1532
DDVKGHSKQDLRVFDY	0.7702	-0.0341	1.7246	0.1893
DDVYWIAAFDY	-0.5247	0.8792	-0.8859	-0.4439
DDWYGGLERGLIQLFDY	0.2621	-0.0544	1.3027	0.3120

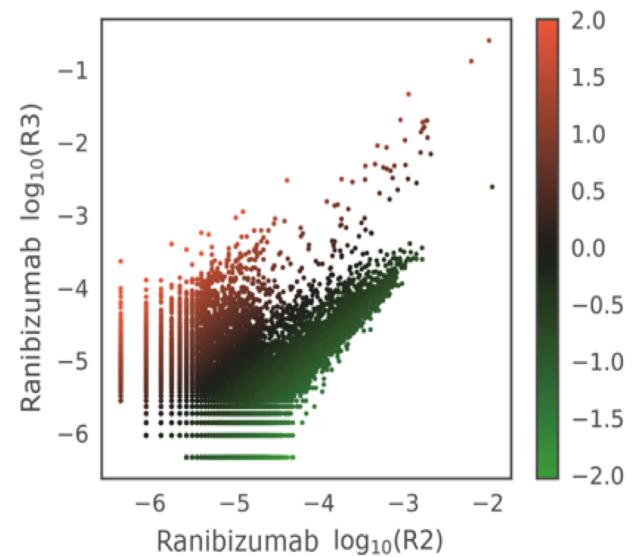


High enrichment suggests high affinity sequences (Ranibizumab, 67769 sequences)

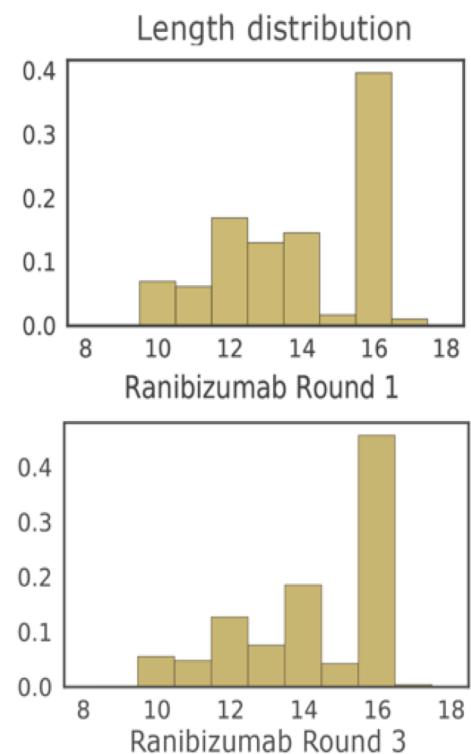
b



c



Ranibizumab binders have preferred CDR-H3 lengths



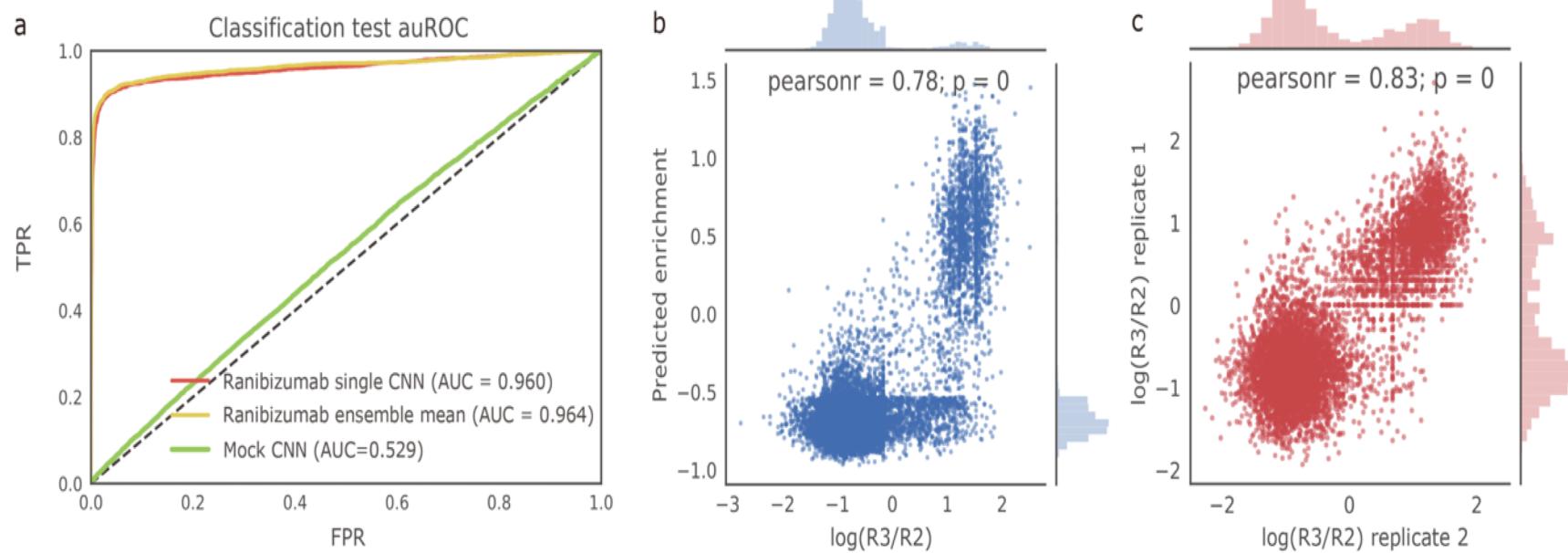
We used six different model architectures

	# of Convolutional layers	# of Convolutional filters	Convolutional filter size	# of Fully connected layer	# of Fully connected neurons	# of parameters in total
2fc	0	0	0	2	32	13954
1conv(32*5)+1fc	1	32	5	1	16	8402
2conv(32*5_64*5)+1fc	2	32,64	5	1	16	18706
1conv(64*5)+1fc	1	64	5	1	16	16754
1conv(32*3)+1fc	1	32	3	1	16	7122
2conv(8*1_64*5)+1fc	2	8, 32	1, 5	1	16	13082

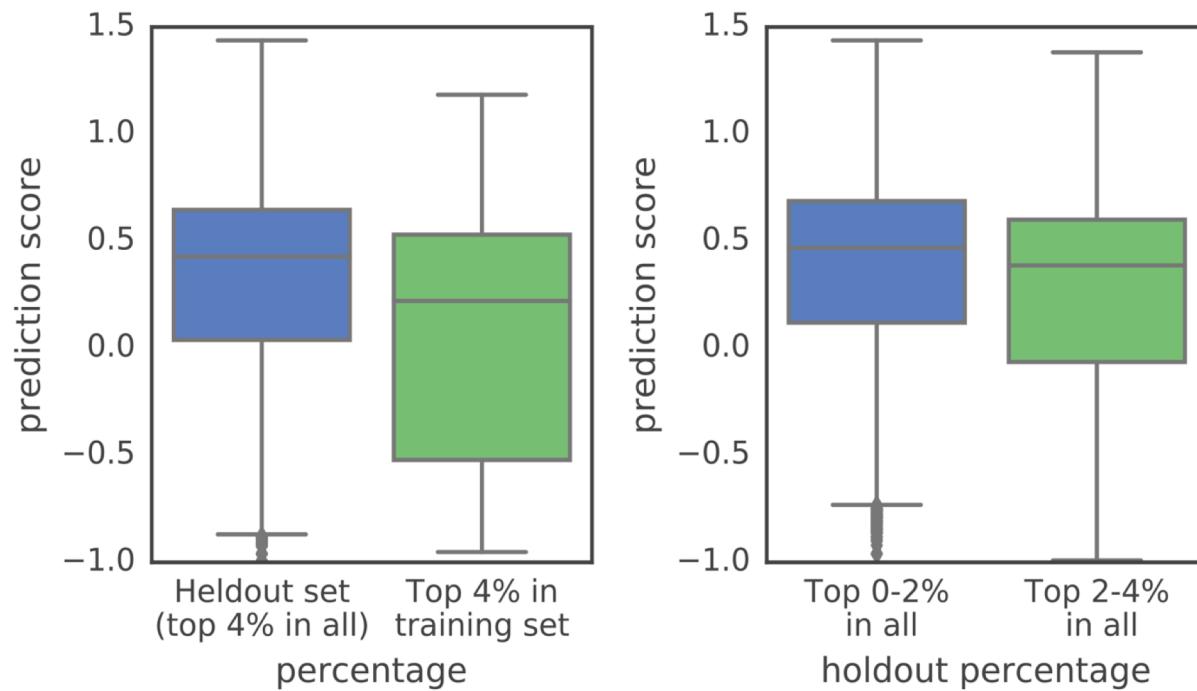
Output layer: Classification – binary cross entropy loss

Regression – mean squared error

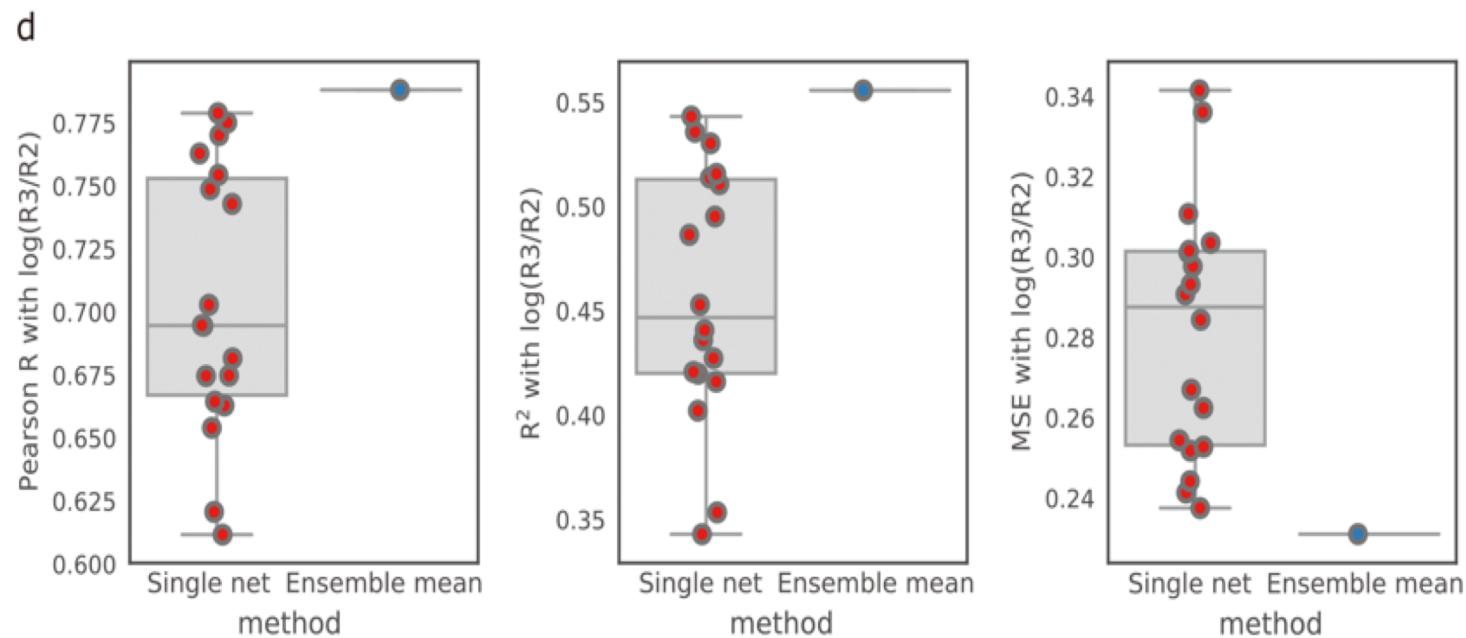
Regression performance is comparable to replicate experiment performance



CNNs produce better scores than they have seen in training for top sequences

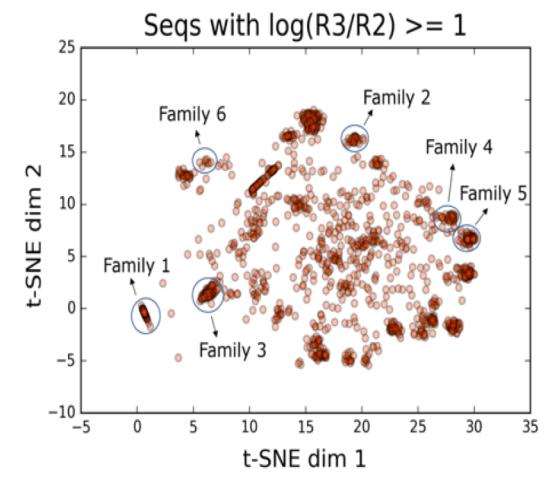
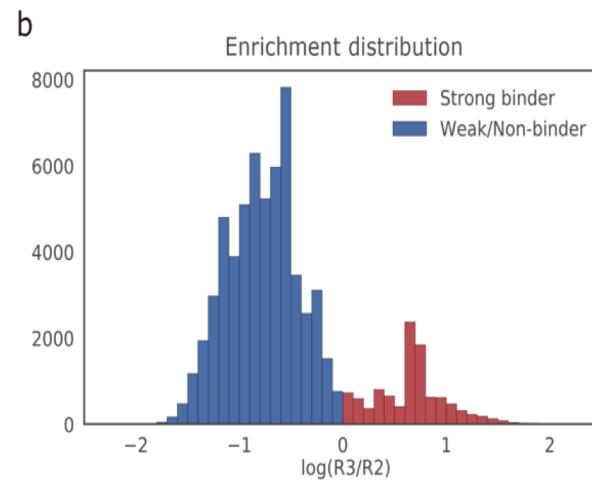
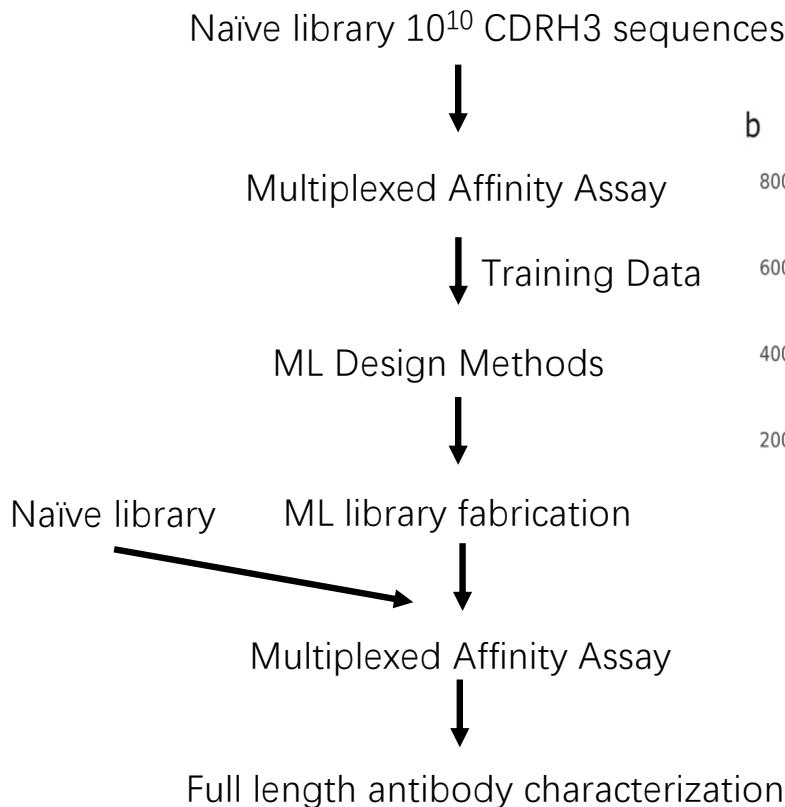


An ensemble of 24 networks is more robust than the individual networks



How can we optimize CDRs?

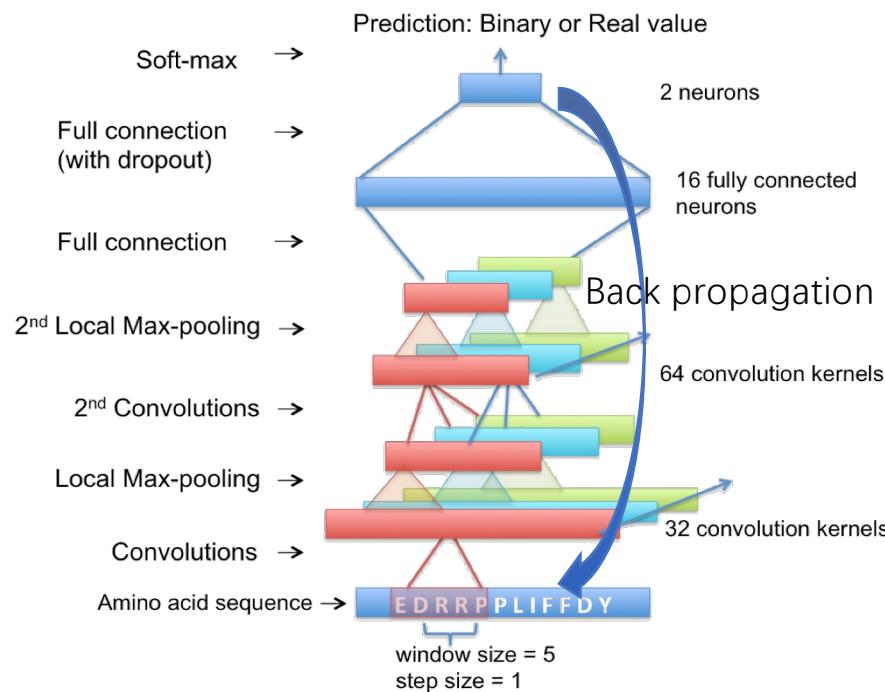
Design flow



Training Data

Our model from sequence to enrichment is differentiable

Method 1 - Optimization with gradients



Projecting continuous representation into one-hot representation

I	0	1	0	0
L	0	0	1	0
V	0	0	0	0
:	:	:	:	:
D	1	0	0	0
K	0	0	0	0
R	0	0	0	1

D I L R

Seed Sequence

Optimization
Gradient ascent

Optimization in continuous space
Gradient ascent

I	0.6	-2	0.2	-5
L	1.2	-1	4.6	0.3
V	-2	0.1	-1	0.7
:	:	:	:	:
D	0.2	3.4	1.1	2.2
K	-4	0.2	-3	-1
R	-1	1.2	-2	6.7

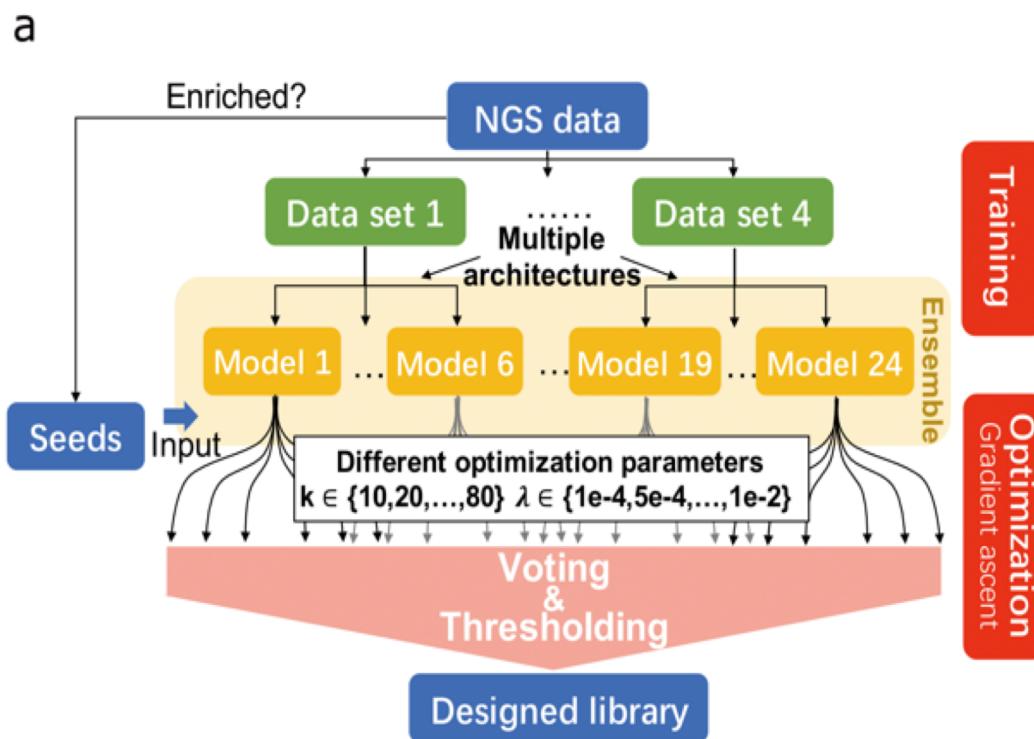
Projection
Every k iteration

I	0	0	0	0
L	1	0	1	0
V	0	0	0	0
:	:	:	:	:
D	0	1	0	0
K	0	0	0	0
R	0	0	0	1

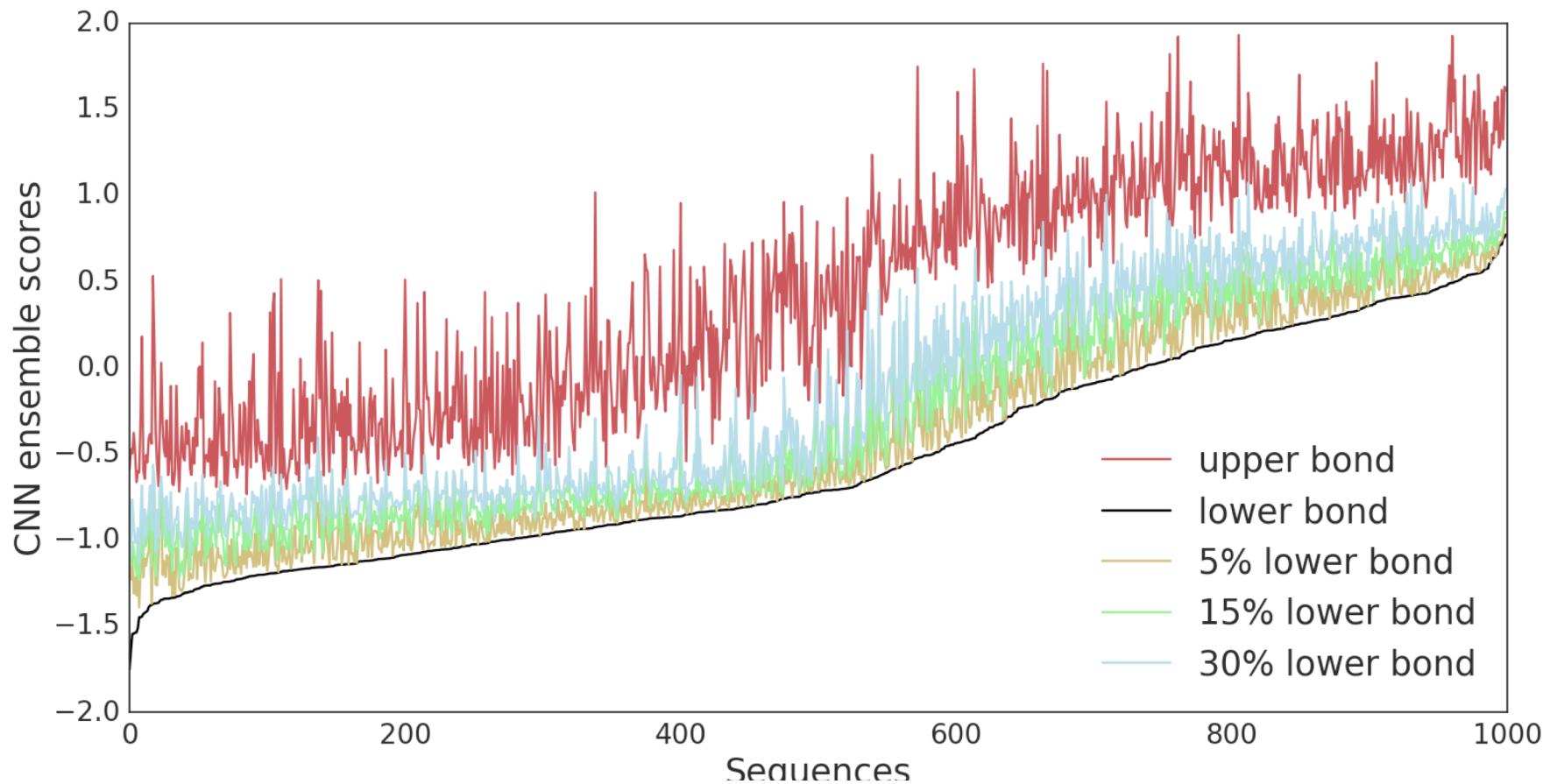
L D L R

New Sequence

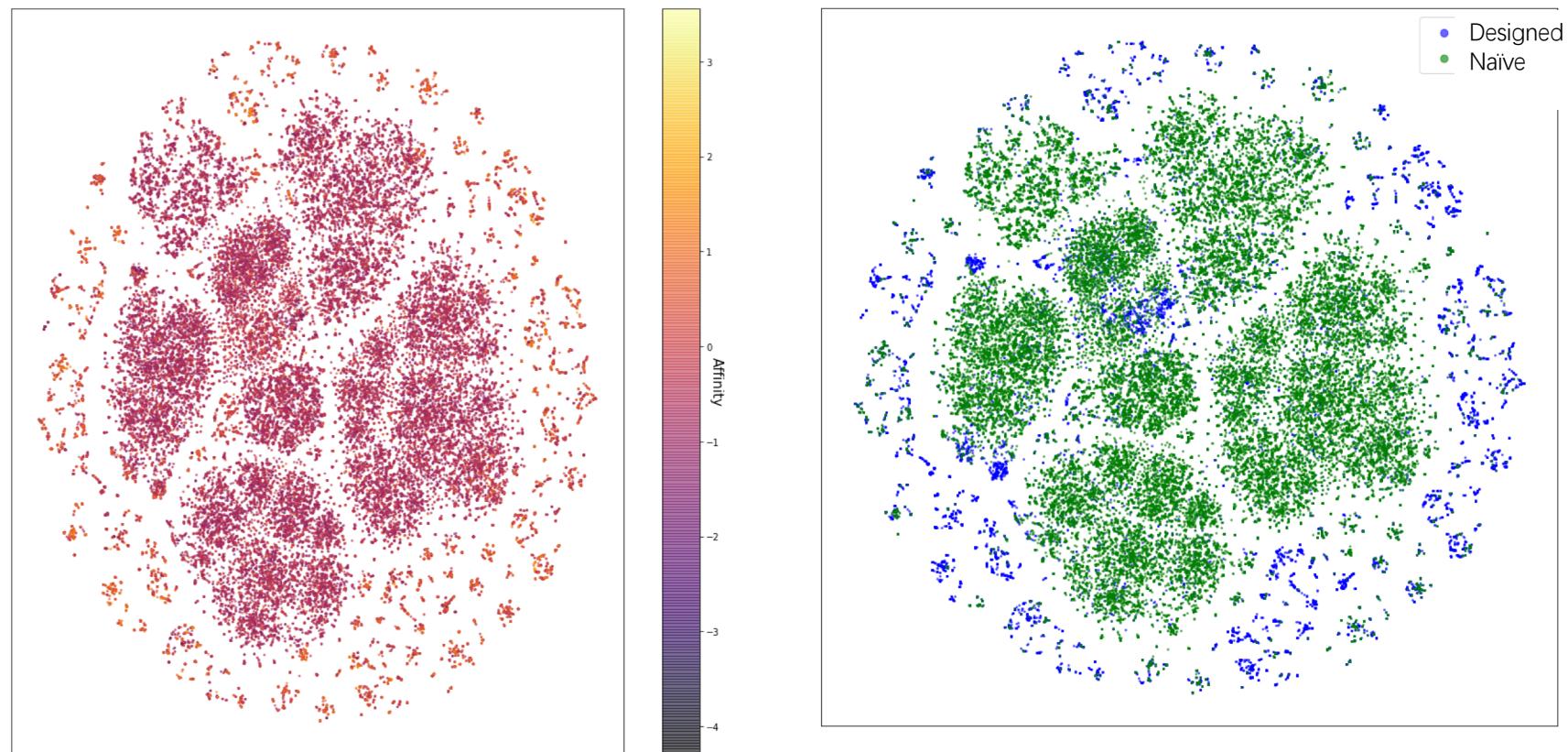
Ens-Grad uses voting across ensembles and hyper-parameters to choose sequences



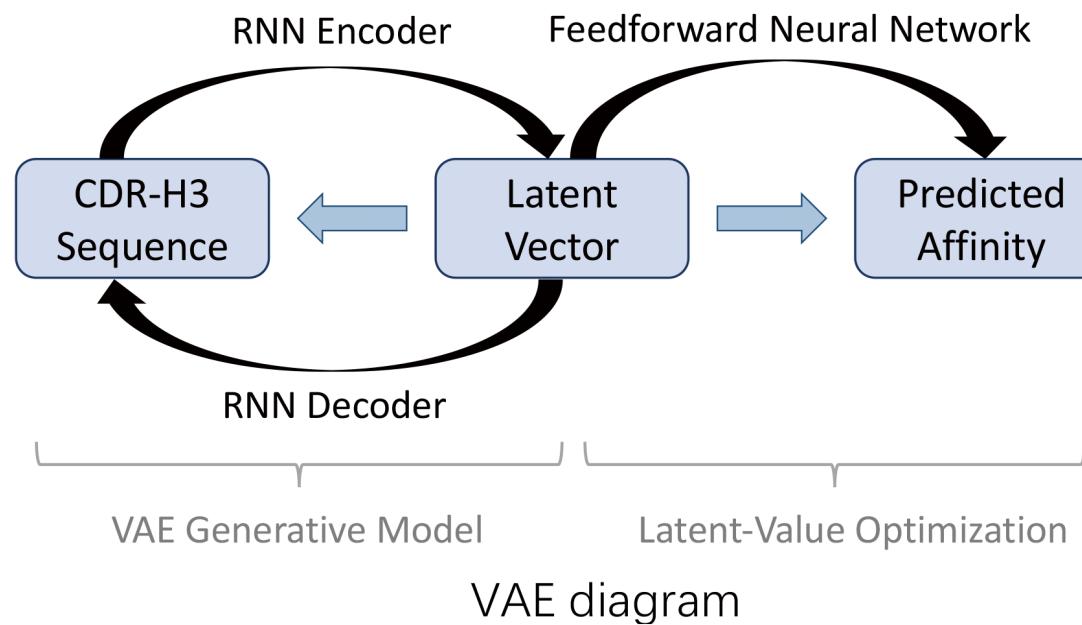
Ranked sequences are filtered by the lower bound
from all 18 networks



Designed sequences appear in islands of enrichment

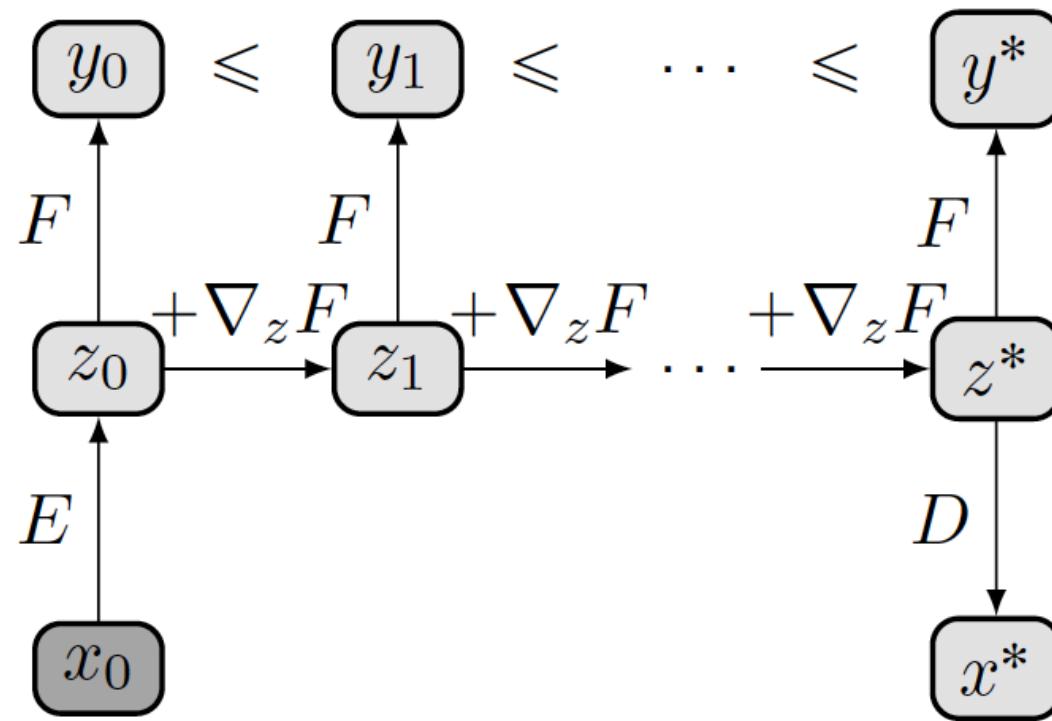


Method 2- Optimize in latent space with a variational autoencoder

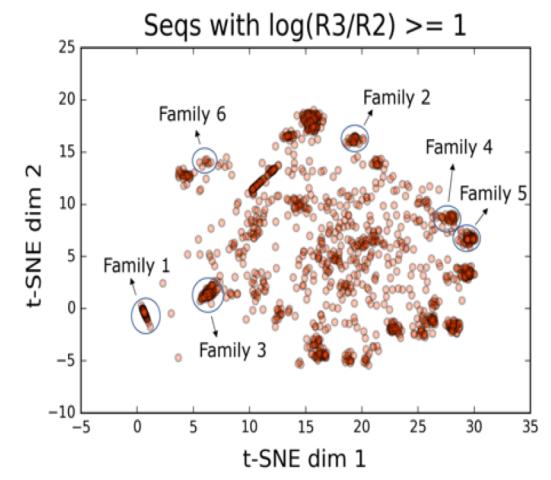
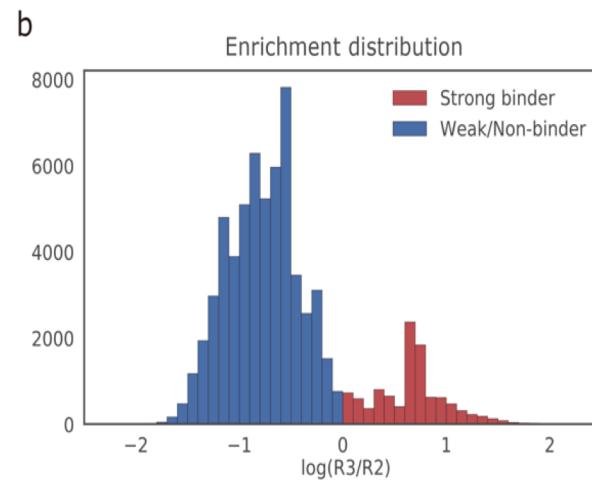
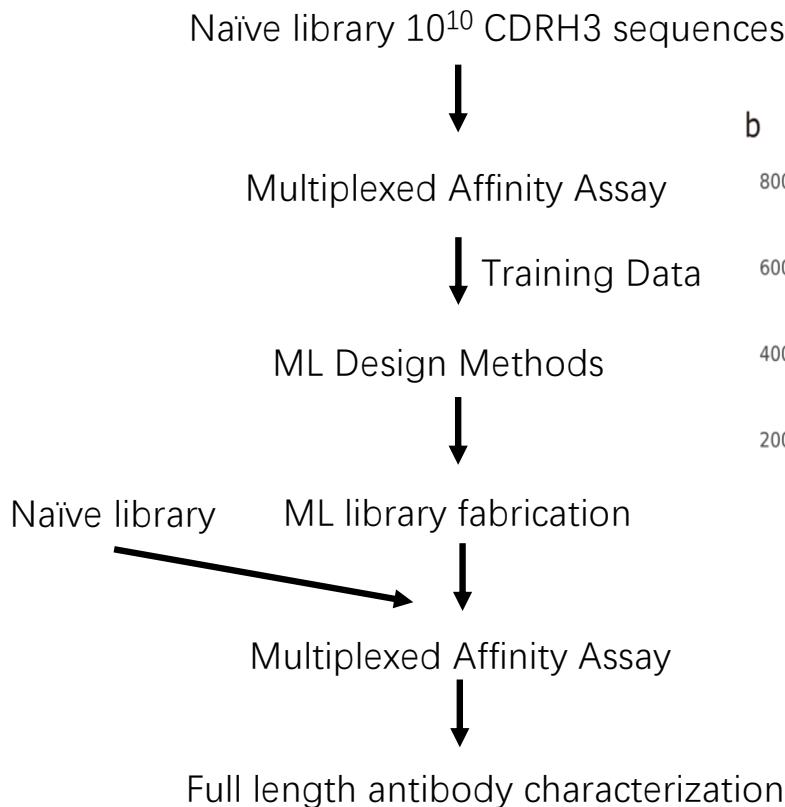


Jonas Mueller, David Gifford, Tommi Jaakkola ;
Proceedings of the 34th International Conference on Machine Learning, PMLR 70:2536-2544, 2017.

Method 2- Optimize in latent space with a variational autoencoder



Design flow

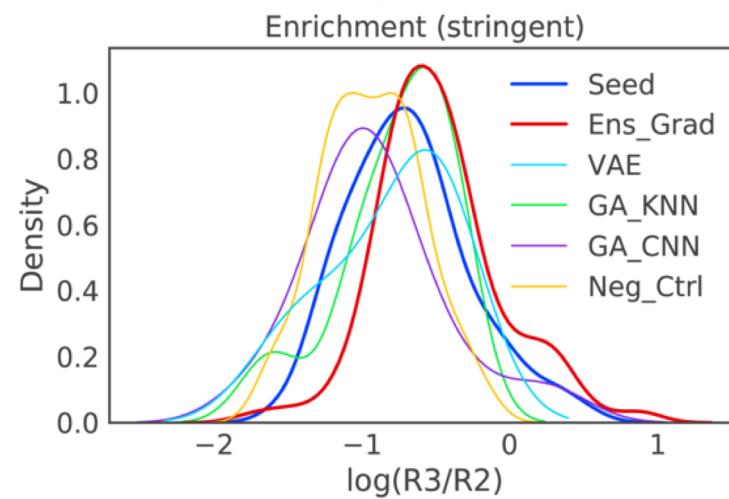
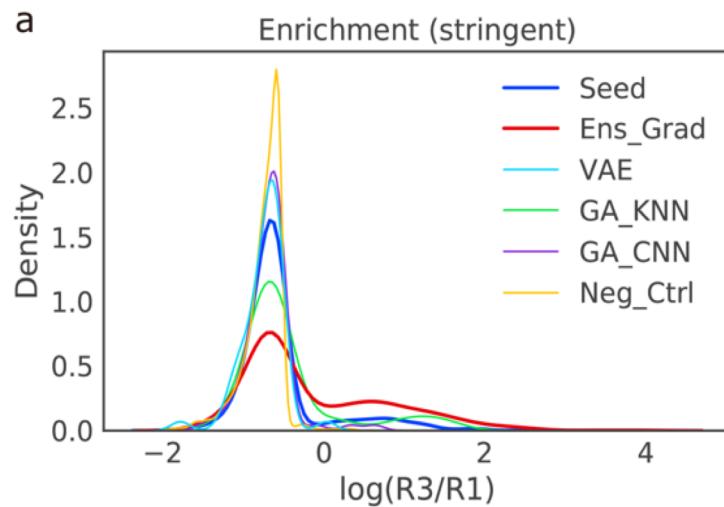


Training Data

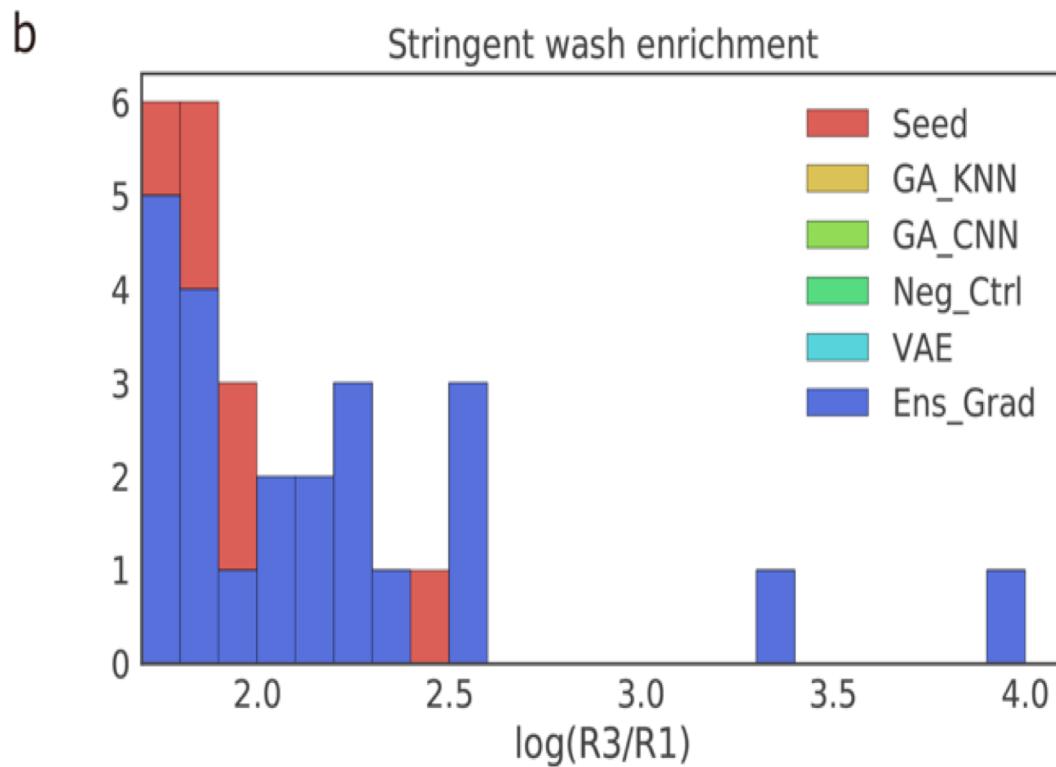
Testing of Fab sequences by direct synthesis

- We computed 77,596 novel machine learning proposed CDR-H3 sequences (Ens-Grad 5,467 sequences)
- We added 26,939 controls and synthesized a total of 104,525 oligonucleotides encoding CDR-H3 sequences
- The oligonucleotides were cloned into a Fab framework and expressed on phage
- Our library with complexity 10^5 was mixed 1:100 into a native library of complexity $\sim 10^{10}$
- The combined library was subject to rounds of panning

Ens-Grad sequences are on average more enriched than seeds and the synthetic results of other ML methods

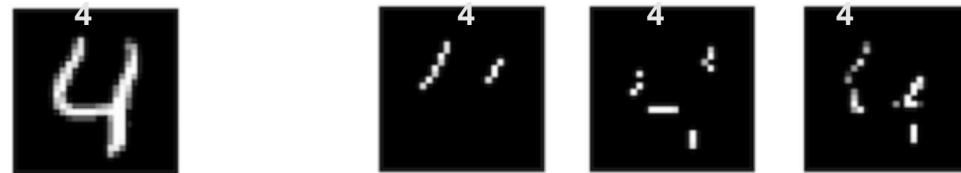


Top scoring synthetic CDR_H3 sequences were designed
by Ens Grad (Stringent Wash)



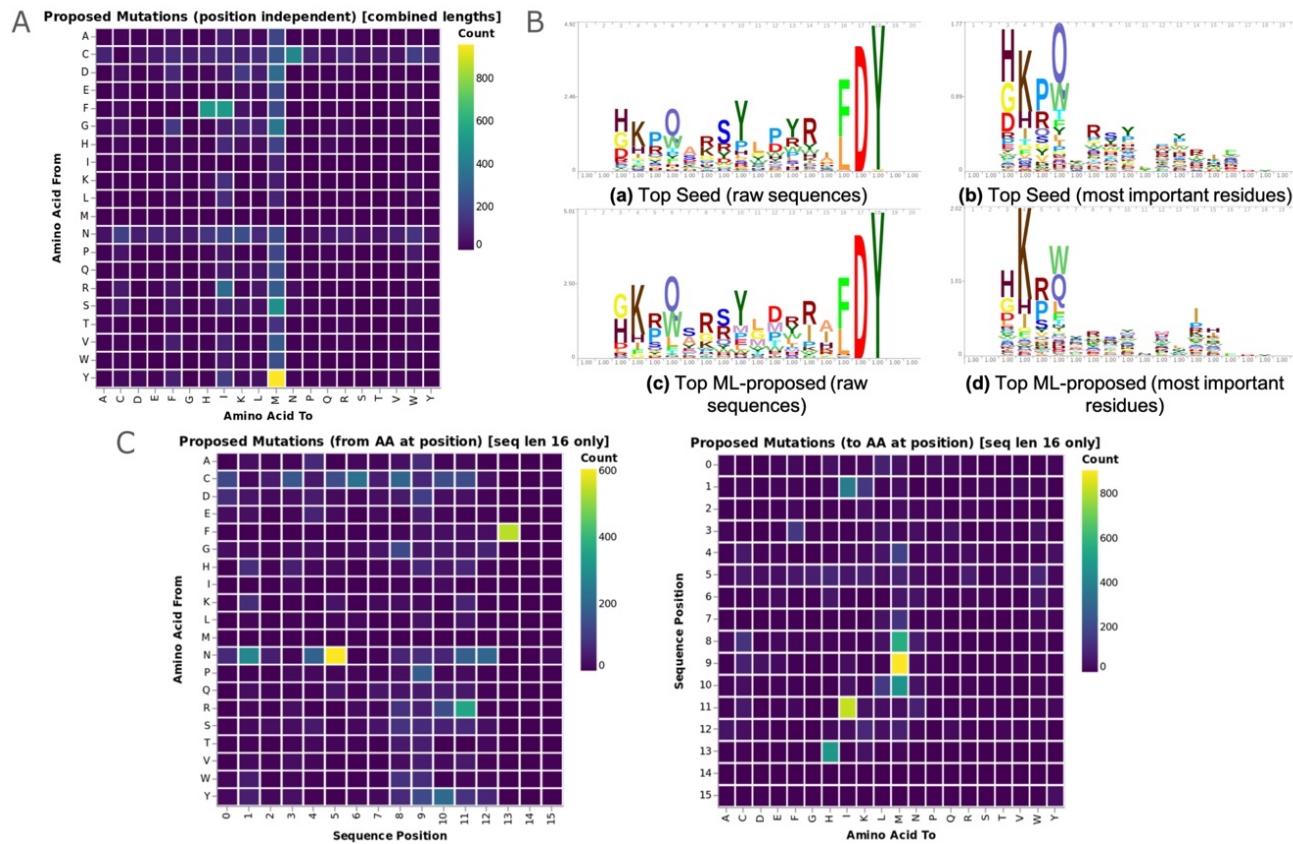
Sufficient Input Subsets provide model interpretation

- One simple rationale for ***why*** a black-box decision is reached is a sparse subset of the input features whose values form the basis for the decision
- A ***sufficient input subset*** (SIS) is a minimal feature subset whose values alone suffice for the model to reach the same decision (even without information about the rest of the features' values)



	CDR-H3 Sequence	Group	R ²	EC50(nM)	Standard log(R3/R1)	Stringent log(R3/R1)
Family 1	HKPQAKSYLPLRLLDY	Ens_Grad	0.99	0.47	3.369	2.399
	HKPQAISYL PYRLLDY	Ens_Grad	0.998	0.5	2.61	2.577
	HKPQAISYLPYRILDY	Seed	0.993	0.62	2.418	2.467
	HKPQAKSYLPMRLLDY	Ens_Grad	0.98	0.93	2.409	0.836
	HKPQAVSYLPYRILDY	Ens_Grad	0.994	0.98	2.915	2.561
	HKPQAKSYL PYRLLDY	Seed	0.996	1.48	2.693	1.128
	HKPQAKSYL PYRTL LDY	Seed	0.993	2.49	2.371	1.986
	HKPQSKSYL PYRLLDY	Seed	0.995	4.78	2.634	0.445
Family 2	HKPQAKSYL PYRILDY	Seed	0.992	6.55	1.41	1.112
	YRSPHHRGGATWQFDY	Seed	0.992	5.79	-0.037	0.036
Family 3	DLFRYYYFMWPLDY	Ens_Grad	0.986	34.05	2.638	0.523
	DLFRYYYFFWPLDY	Seed	0.99	109.5	2.988	1.283
Family 4	MHYYDIGVFPWDTFDY	Ens-Grad	0.971	0.29	2.089	3.381
	GHYYDIGVFPWDTFDY	Seed	0.99	0.49	0.703	1.593
Family 5	WQQWAGYPRQKYSF DY	Seed	0.986	3.31	2.657	1.888
	WQQWSGYPRQKYSF DY	Seed	0.975	66.81	0.264	-0.219
Family 6	GKSLYGOETTWP HF DY	Seed	0.99	0.67	2.002	0.946

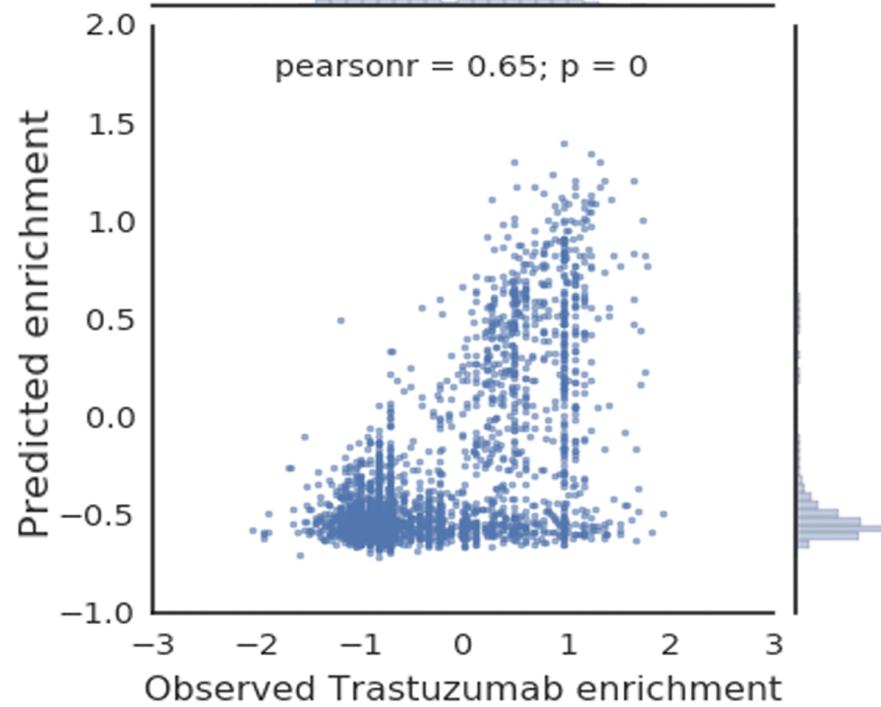
What ML changes facilitate enrichment improvement?



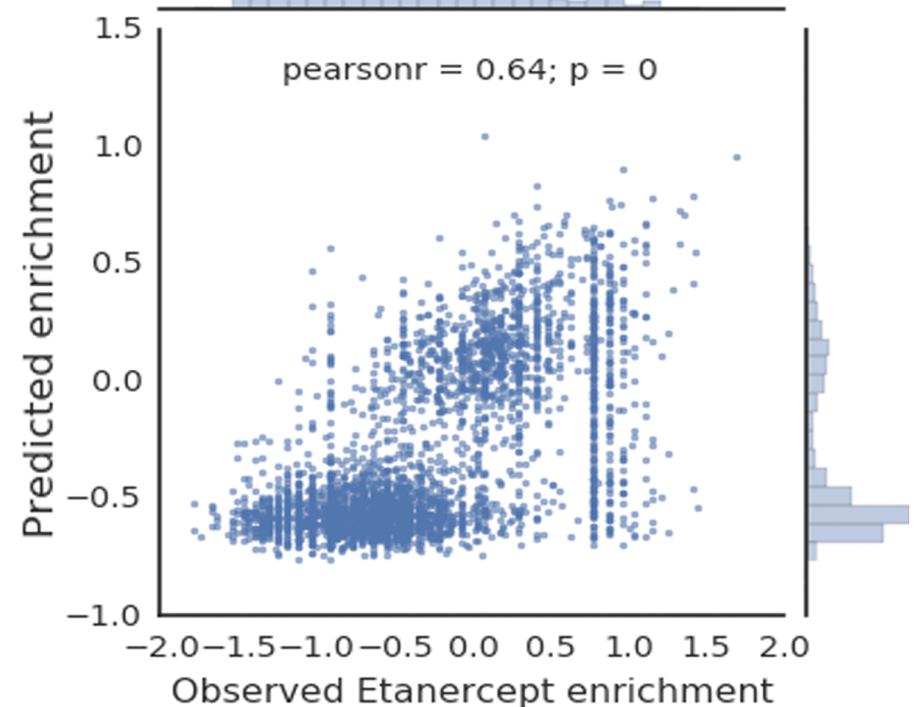
Multi-objective optimization for specificity

We can build a joint model of Trastuzumab and Entanercept binding

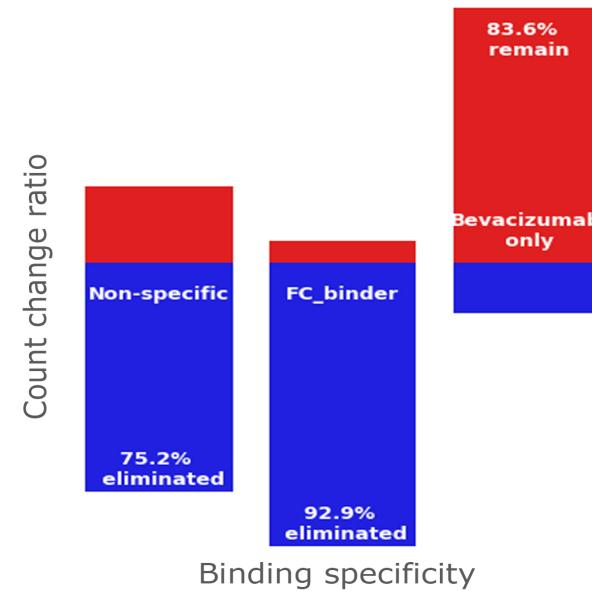
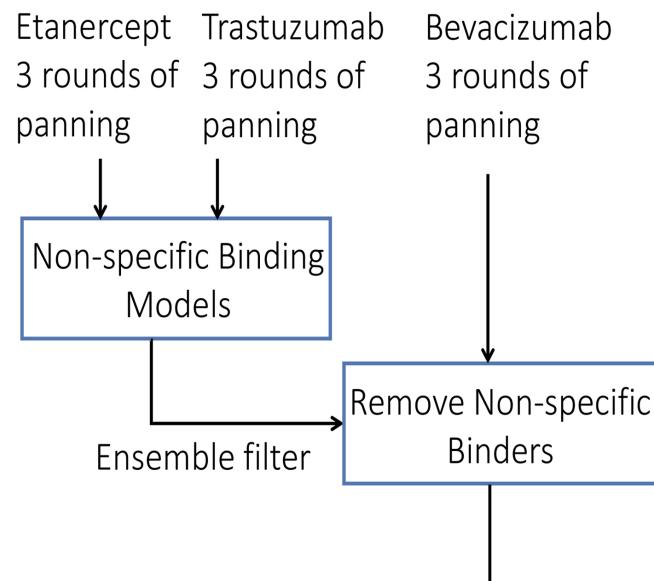
C



D

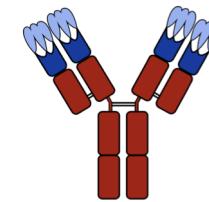
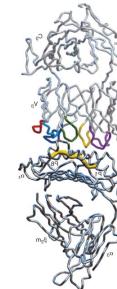


We can remove non-specific binders



We wish to both identify therapeutic targets and therapeutic molecules with the help of machine learning

- Identifying target peptide-MHC molecules
 - Best prediction model
 - Likelihood metrics can optimize vaccine formulation
- Designing antibody Complementarity Determining Regions (CDRs)
 - Differentiable models from sequence to objective
 - Enables multi-objective optimization



Acknowledgements

MIT

Saber (Ge Liu)
Haoyang Zeng
Jonas Mueller
Brandon Carter
Ziheng Wang (Tony)
Nathan Hunt
Michael Birnbaum

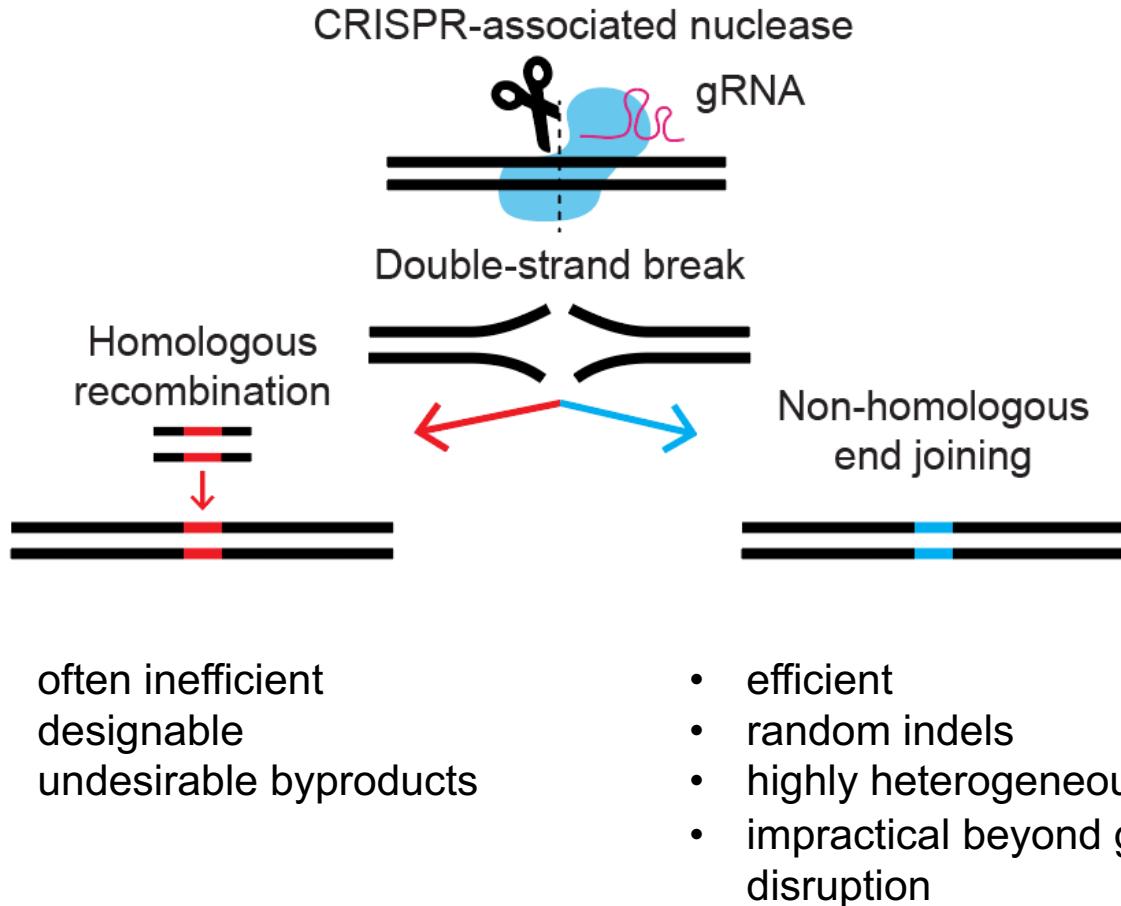
Novartis

Stefan Ewert
Jonas Schilz
Geraldine Horny



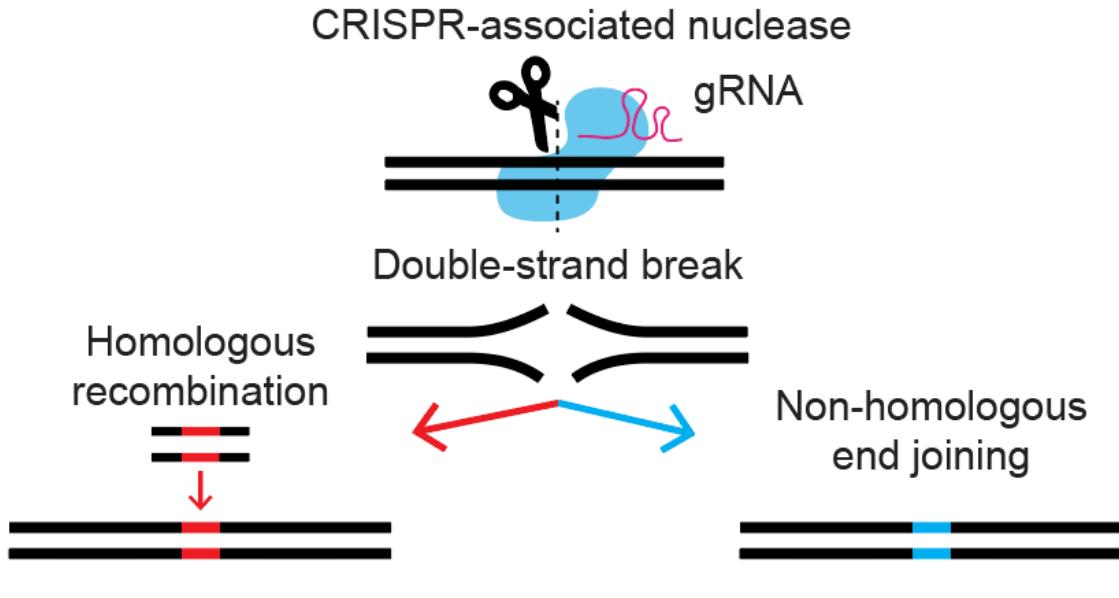
Predictable and precise template-free editing of pathogenic mutations by CRISPR-Cas9 nuclease

The state of CRISPR genome editing



reference
CATGGGGGTATAGGGCTAAATGT 2986
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 2388
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 11960
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1961
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1956
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 465
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 281
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 3426
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 7266
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 26791
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 17111
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1918
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 3673
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 972
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 646
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 4192
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1278
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 435
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 7388
CGTATAGATTTGGATATGGCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 2658
CATGGGGGTATAGGGCTAAATGT 234
CATGGGGGTATAGGGCTAAATGT 4588
CATGGGGGTATAGGGCTAAATGT 876
CATGGGGGTATAGGGCTAAATGT 745
CATGGGGGTATAGGGCTAAATGT 933
CATGGGGGTATAGGGCTAAATGT 5290
CATGGGGGTATAGGGCTAAATGT 754
CATGGGGGTATAGGGCTAAATGT 1028
CATGGGGGTATAGGGCTAAATGT 1682
CATGGGGGTATAGGGCTAAATGT 935
CATGGGGGTATAGGGCTAAATGT 5347
CATGGGGGTATAGGGCTAAATGT 3129
CATGGGGGTATAGGGCTAAATGT 1852
CATGGGGGTATAGGGCTAAATGT 1254
CATGGGGGTATAGGGCTAAATGT 492
CATGGGGGTATAGGGCTAAATGT 28310
CATGGGGGTATAGGGCTAAATGT 28216
CATGGGGGTATAGGGCTAAATGT 116
CATGGGGGTATAGGGCTAAATGT 4485
CATGGGGGTATAGGGCTAAATGT 2385
CATGGGGGTATAGGGCTAAATGT 1096
CATGGGGGTATAGGGCTAAATGT 1533
CATGGGGGTATAGGGCTAAATGT 377
CATGGGGGTATAGGGCTAAATGT 194
CATGGGGGTATAGGGCTAAATGT 175
CATGGGGGTATAGGGCTAAATGT 188
CATGGGGGTATAGGGCTAAATGT 181
CATGGGGGTATAGGGCTAAATGT 426
CATGGGGGTATAGGGCTAAATGT 357
CATGGGGGTATAGGGCTAAATGT 4572
CATGGGGGTATAGGGCTAAATGT 1136
CATGGGGGTATAGGGCTAAATGT 484
CATGGGGGTATAGGGCTAAATGT 252
CATGGGGGTATAGGGCTAAATGT 1929
CATGGGGGTATAGGGCTAAATGT 423
CATGGGGGTATAGGGCTAAATGT 3700
CATGGGGGTATAGGGCTAAATGT 610
CATGGGGGTATAGGGCTAAATGT 254
CATGGGGGTATAGGGCTAAATGT 218
CATGGGGGTATAGGGCTAAATGT 320
CATGGGGGTATAGGGCTAAATGT 458
CATGGGGGTATAGGGCTAAATGT 104
CATGGGGGTATAGGGCTAAATGT 1354
CATGGGGGTATAGGGCTAAATGT 308
CATGGGGGTATAGGGCTAAATGT 720
CATGGGGGTATAGGGCTAAATGT 351
CATGGGGGTATAGGGCTAAATGT 53856
CATGGGGGTATAGGGCTAAATGT 178
CATGGGGGTATAGGGCTAAATGT 652
CATGGGGGTATAGGGCTAAATGT 3194
CATGGGGGTATAGGGCTAAATGT 328
CATGGGGGTATAGGGCTAAATGT 149

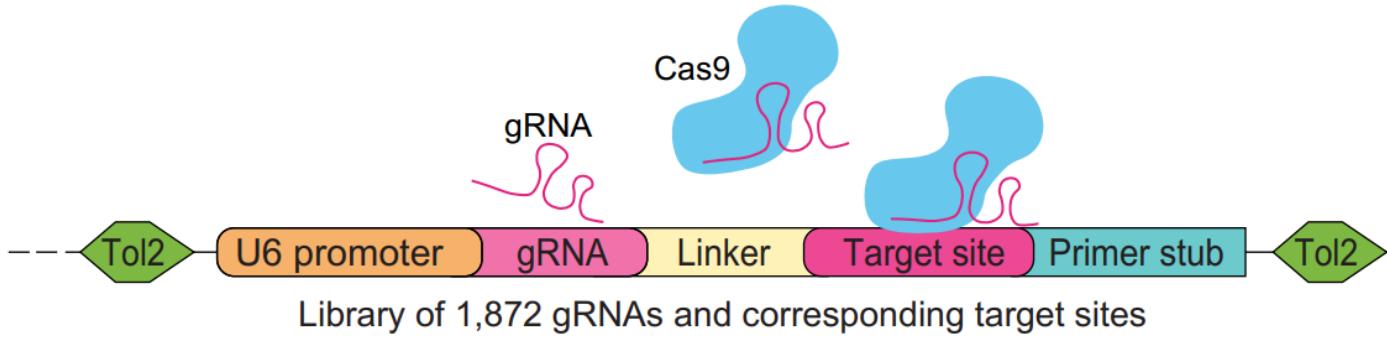
The state of CRISPR genome editing



- often inefficient
- designable
- predictable byproducts
- efficient
- predictable indels
- can be homogeneous
- practical: repair of pathogenic alleles to wild-type

reference
CATGGGGGTATAGGGCTAAATGT 2986
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 2388
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 11960
CGTATAGATTTGGATATGGCCCATG-AAGT
CATGGGGGTATAGGGCTAAATGT 1961
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1956
CGTATAGATTTGGATATGGCCCATGTA
CATGGGGGTATAGGGCTAAATGT 465
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 281
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 3426
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 7266
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1711
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 1918
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 3673
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 972
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 646
CGTATAGATTTGGATATGGCCCATGTA
CATGGGGGTATAGGGCTAAATGT 4192
CGTATAGATTTGGATATGGCCCATGTA
CATGGGGGTATAGGGCTAAATGT 1278
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 435
CGTATAGATTTGGATATGGCCCATGTAGTA
CATGGGGGTATAGGGCTAAATGT 7388
CGTATAGATTTGGATATGGCCCATGT
CATGGGGGTATAGGGCTAAATGT 2658
CATGGGGGTATAGGGCC-
CATGGGGGTATAGGGCTAAATGT 234
CATGGGGGTATAGGGCC-
CATGGGGGTATAGGGCTAAATGT 4588
CATGGGGGTATAGGGCTAAATGT 870
CATGGGGGTATAGGGCTAAATGT 745
CATGGGGGTATAGGGCTAAATGT 933
CATGGGGGTATAGGGCTAAATGT 5290
CATGGGGGTATAGGGCTAAATGT 754
CATGGGGGTATAGGGCTAAATGT 1828
CATGGGGGTATAGGGCTAAATGT 1682
CATGGGGGTATAGGGCTAAATGT 935
TATAGGTGGCCTAAATGT 5347
CGGTATAGGTGGCCTAAATGT 3129
ATCATGGGGGTATAGGTGGCCTAAATGT 1852
TAGGTGGCCTAAATGT 1254
GGTGGCCTAAATGT 492
CGGGCGTATAGGTGGCCTAAATGT 28310
GTATAGGTGGCCTAAATGT 28216
GTATAGGTGGCCTAAATGT 116
CGGGCGTATAGGTGGCCTAAATGT 4485
GTGGCCTAAATGT 2385
ATAGGTGGCCTAAATGT 1096
GGCCTAAATGT 1533
CGTAAATGT 377
GGTGGCCTAAATGT 194
AGGTGGCCTAAATGT 175
GGCCTAAATGT 188
GGTGGCCTAAATGT 181
TATAGGT 426
GGTGGCCTAAATGT 357
AGGTGGCCTAAATGT 4572
TGT 1136
CCTAAATGT 484
GT 480
CGGTAAATGT 252
TGGCCTAAATGT 1929
AGGTGGCCTAAATGT 423
TGGCCTAAATGT 3700
GGTGGCCTAAATGT 610
CGGGCGTATAGGTGGCCTAAATGT 254
TAAATGT 218
CCTAAATGT 320
CTAAATGT 458
ATAGGTGGCCTAAATGT 104
ATAGGTGGCCTAAATGT 1354
AAATGT 720
AAATGT 351
GGTGGCCTAAATGT 53856
GTGGCCTAAATGT 178
GT 652
CCGTATAGGTGGCCTAAATGT 3194
GT 328
ATGT 149

High-throughput genome-integrated assay of Cas9-mediated DNA repair



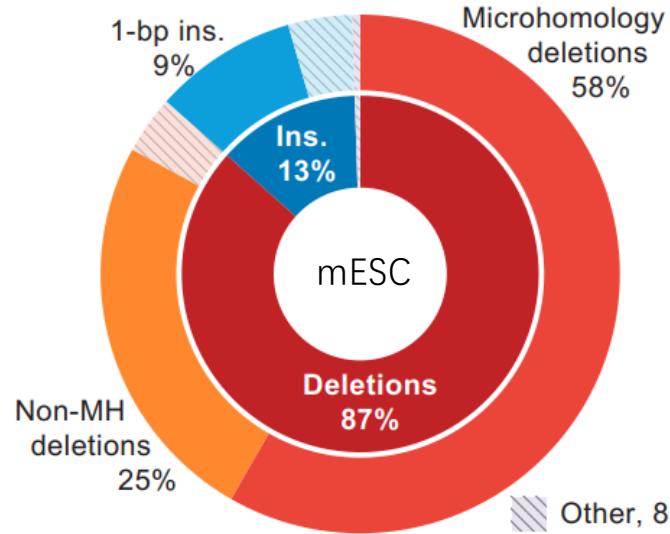
- 96 target sites in largest previous study
- Designed 1,872 target sites (55-bp) based on the human genome

Target Site Sequence	Count	Reference
CGTATAGATTTGGATATGGGCATGTAGTA	2986	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	2388	-ATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	11860	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1961	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1956	CA-CGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	465	AT-GGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	281	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	3426	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	7266	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	26791	-ATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1711	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1918	-CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	3673	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	972	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	646	CATACATAAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	4192	-TCGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1278	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	435	-CGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	7388	-GCCGTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	2658	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	234	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	4588	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	878	-ATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	745	-AGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	933	-TATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	5200	-CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	754	-AGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	1828	-CGTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCAA-	1682	CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCAA-	935	-CATGGCCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCCA-	5347	TATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCCA-	3129	-CCGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	1852	-ATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCCATGT-	1254	-TAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCCCATGT-	492	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	28310	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	28216	-GTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	116	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	4485	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	2395	-GTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	1096	-ATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGCC-	1533	-GGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	377	-GCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	194	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	175	-AGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	188	-GGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	181	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	426	-TATAGGT
CGTATAGATTTGGATATGGGC-	357	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGGGC-	4572	-AGGTGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	1136	-TGT
CGTATAGATTTGGATATGGGCATGT-	484	-CCTAAATGT
CGTATAGATTTGGATATGGGCATGT-	486	-GT
CGTATAGATTTGGATATGGGCAT-	252	-CCTAAATGT
CGTATAGATTTGGATATGGG-	1929	-TGGCCTAAATGT
CGTATAGATTTGGATATGG-	423	-AGGTGGCCTAAATGT
CGTATAGATTTGGATATGG-	3700	-TGGCCTAAATGT
CGTATAGATTTGGATATGGGCATGTAGTA	610	-GGTGGCCTAAATGT
CGTATAGATTTGGATATGG-	131	-CGGGCGTTATAGGTGGCCTAAATGT
CGTATAGATTTGGATATGG-	254	-TAAATGT
CGTATAGATTTGGATATGG-	218	-CTAAATGT
CGTATAGATTTGGATATGG-	328	-CGCTAAATGT

62

Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Cas9 primarily causes microhomology deletions in genome-integrated and endogenous settings

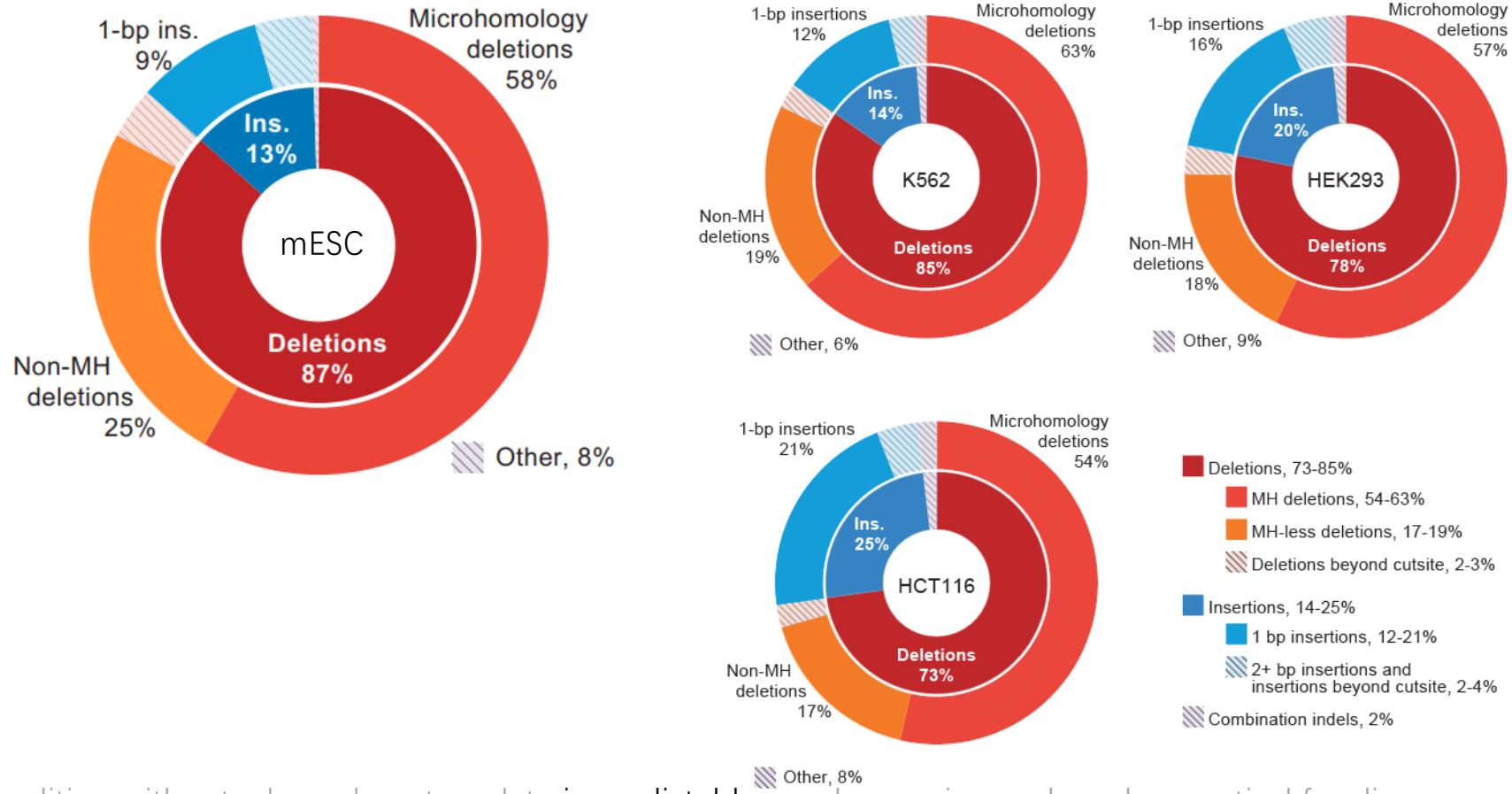


A microhomology deletion is a deletion with multiple equal-scoring alignments



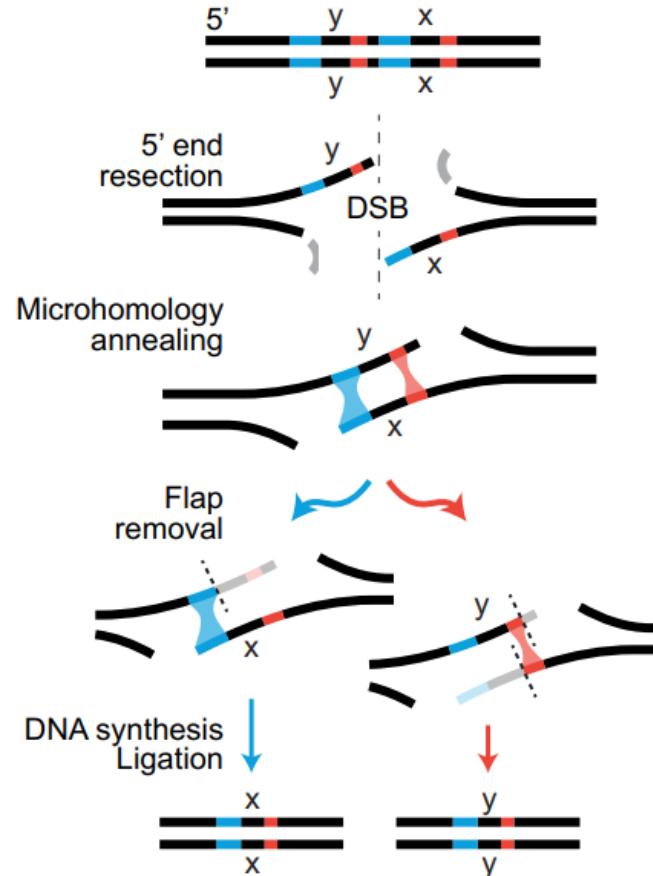
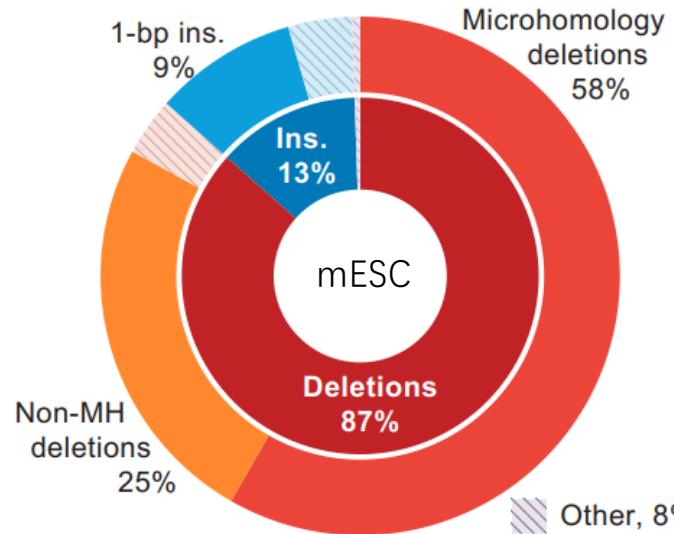
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Cas9 primarily causes microhomology deletions in genome-integrated and endogenous settings



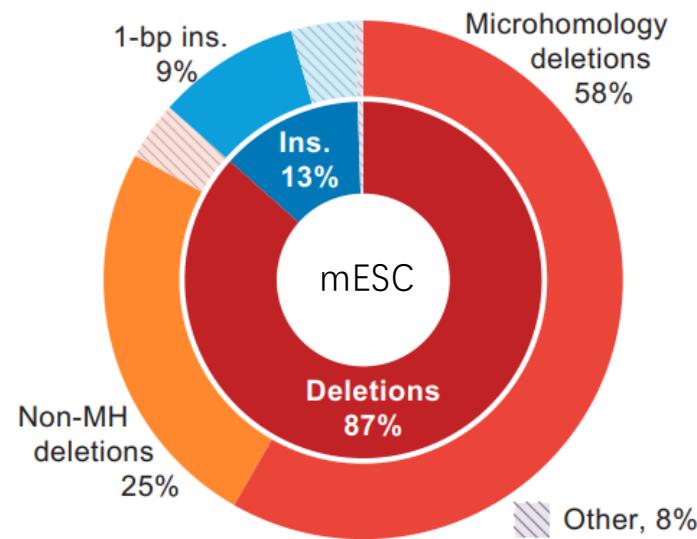
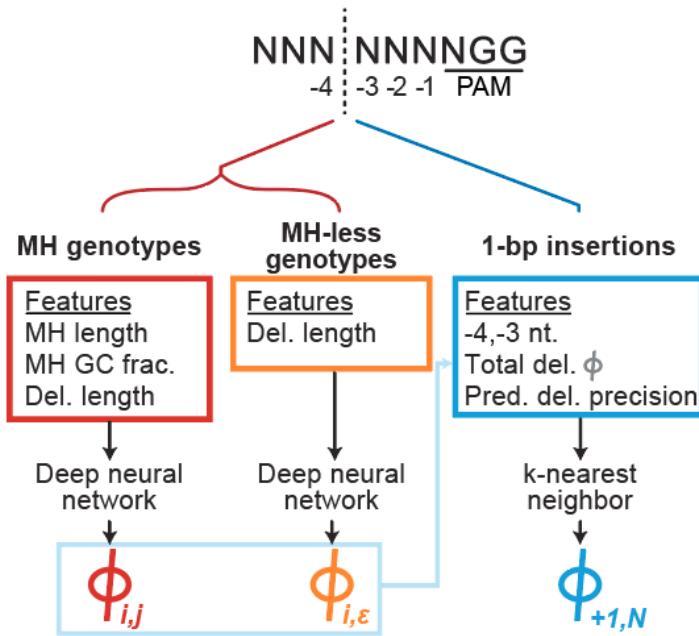
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Majority of repair products arise from microhomology-mediated end-joining (MMEJ)

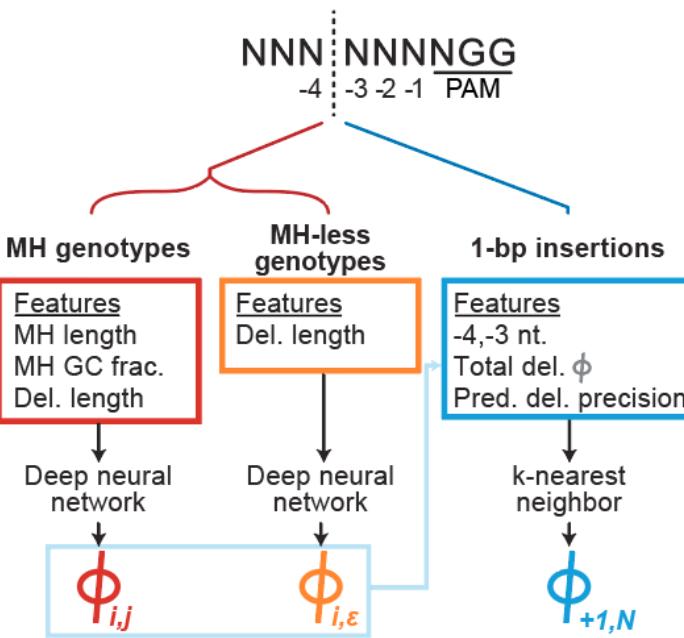


Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

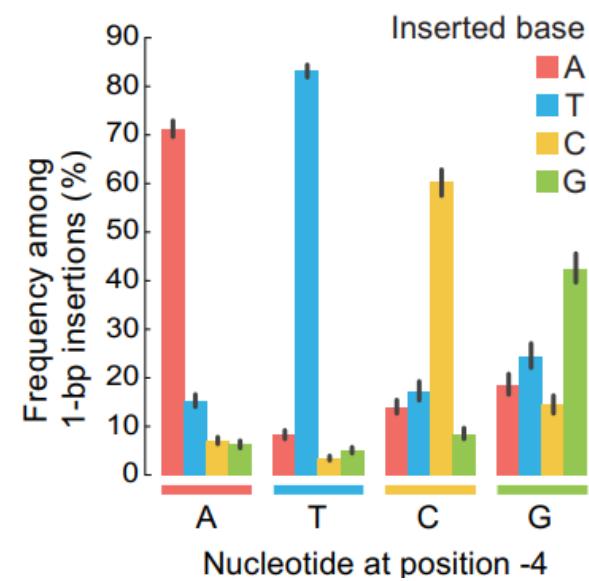
inDelphi predicts 90% of repair products from 3 major repair classes



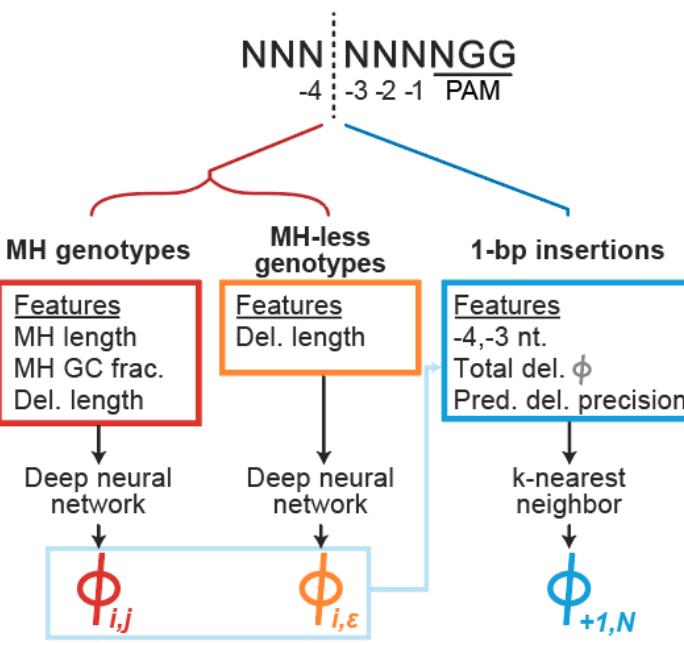
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease



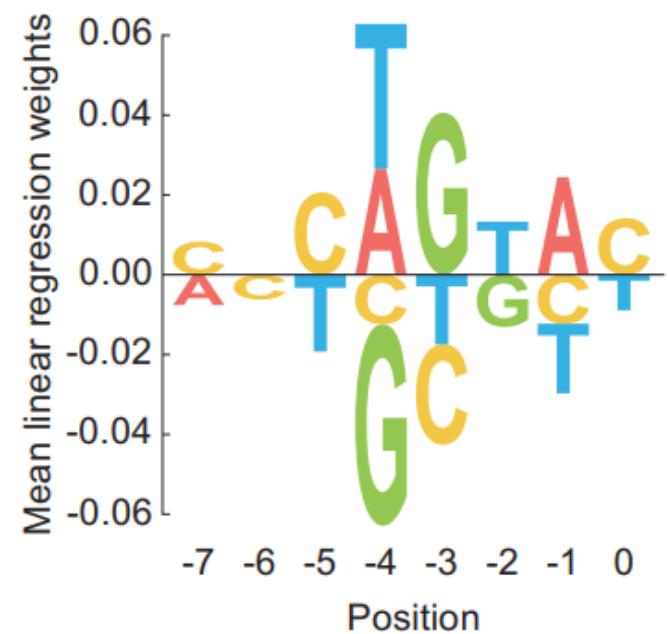
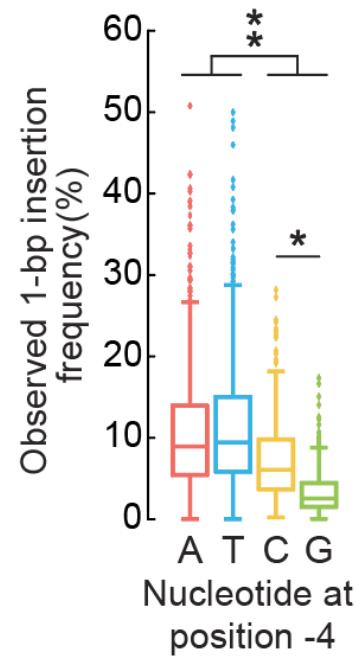
1-bp insertions copy the adjacent nucleotide



Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

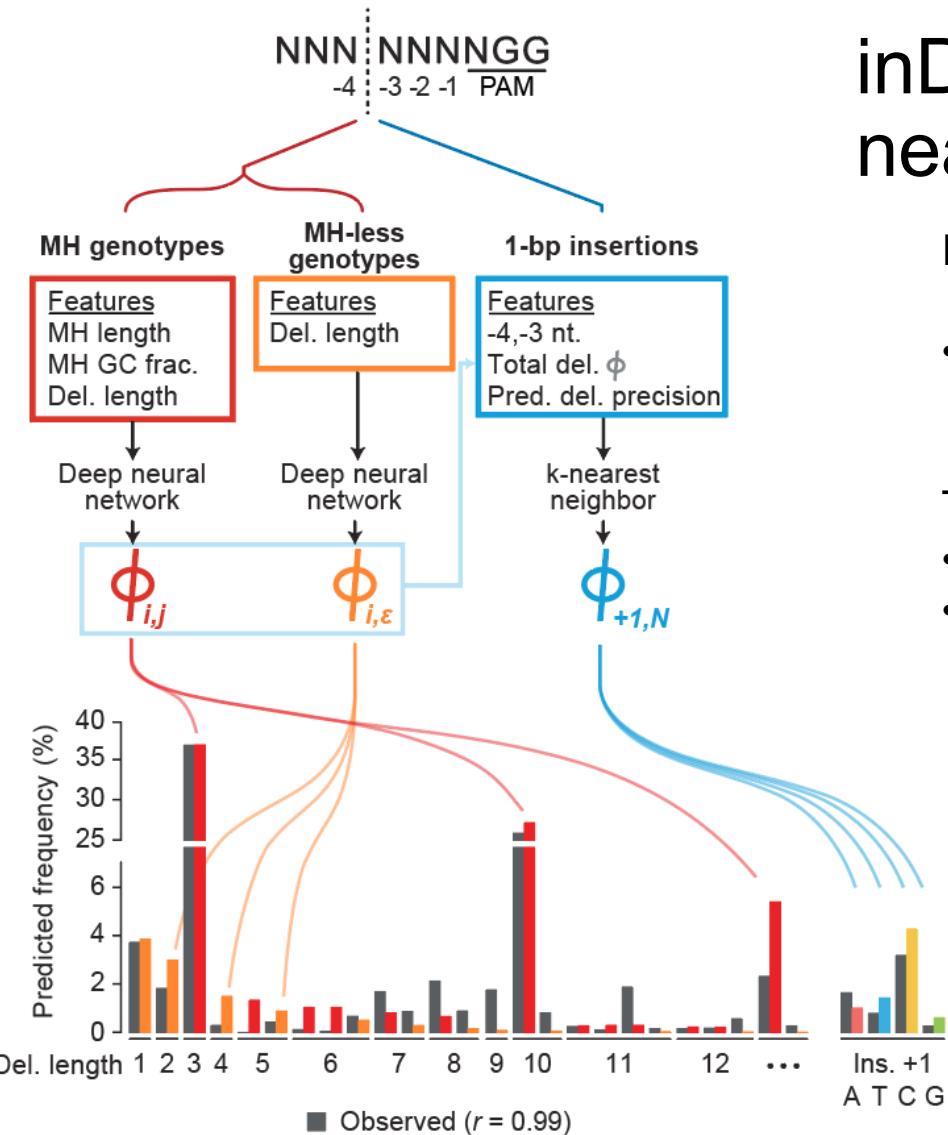


1-bp insertion frequency depends on local sequence context



Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

inDelphi accurately predicts nearly all repair outcomes

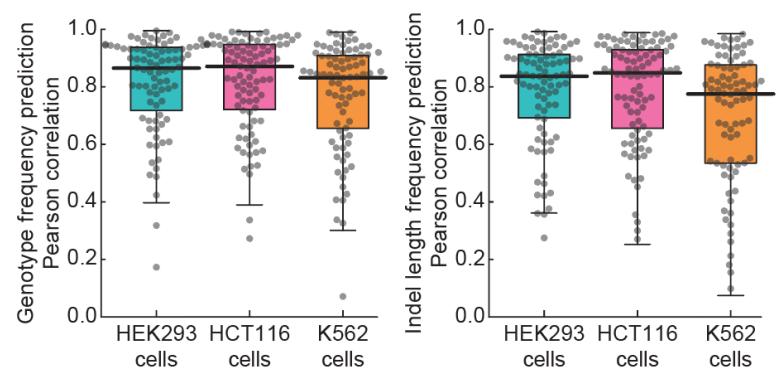


Input: Sequence, cutsite

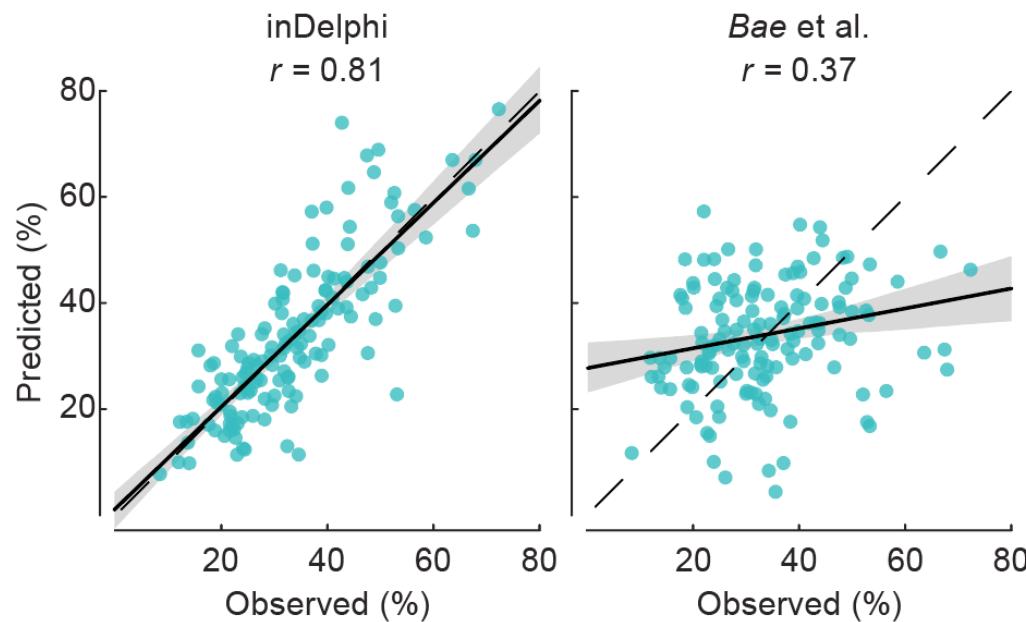
- Predicts 90% of observed repair outcomes
 - 70% at single-base resolution

Training & testing on held-out cell-types

- Median $r = 0.87$ on genotype prediction
- Median $r = 0.84$ on indel length prediction

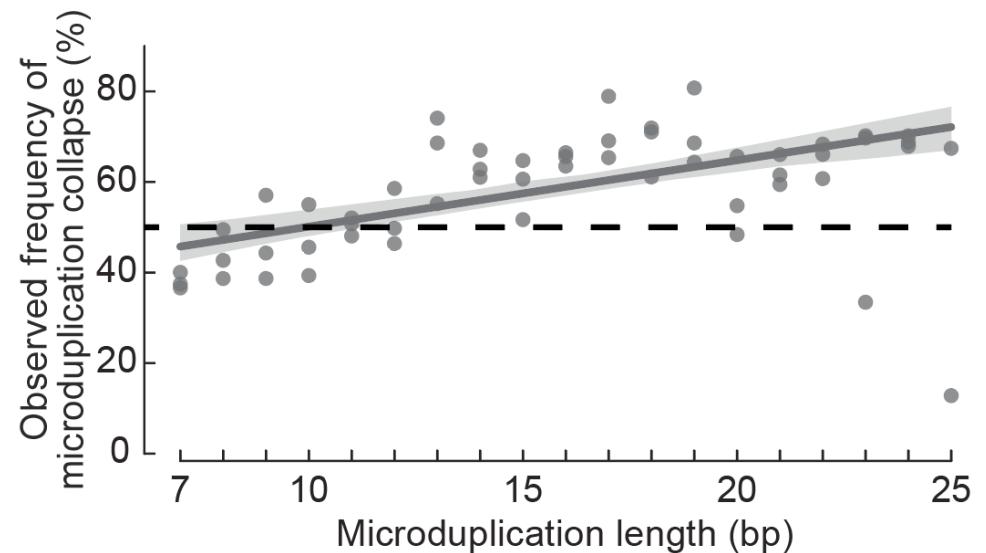
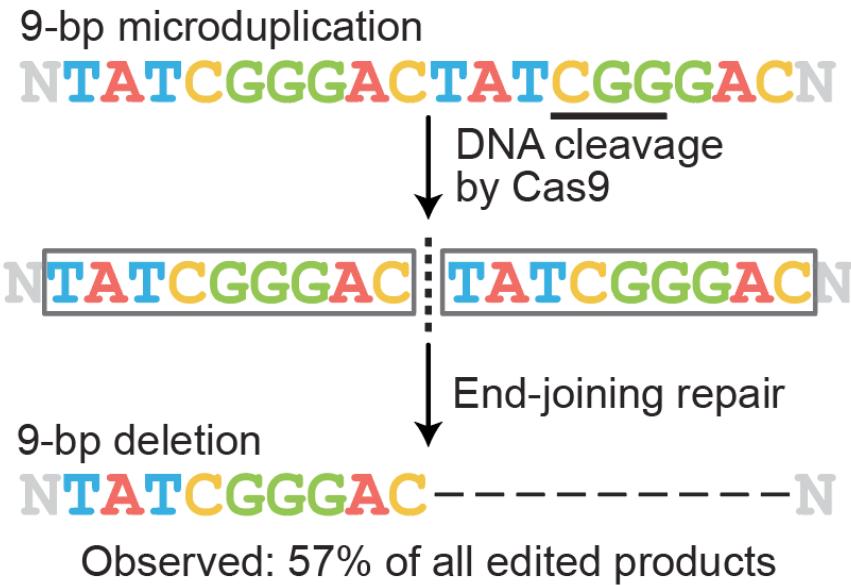


inDelphi accurately predicts frameshifts



Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Target sites yielding a single deletion repair genotype >50% of the time



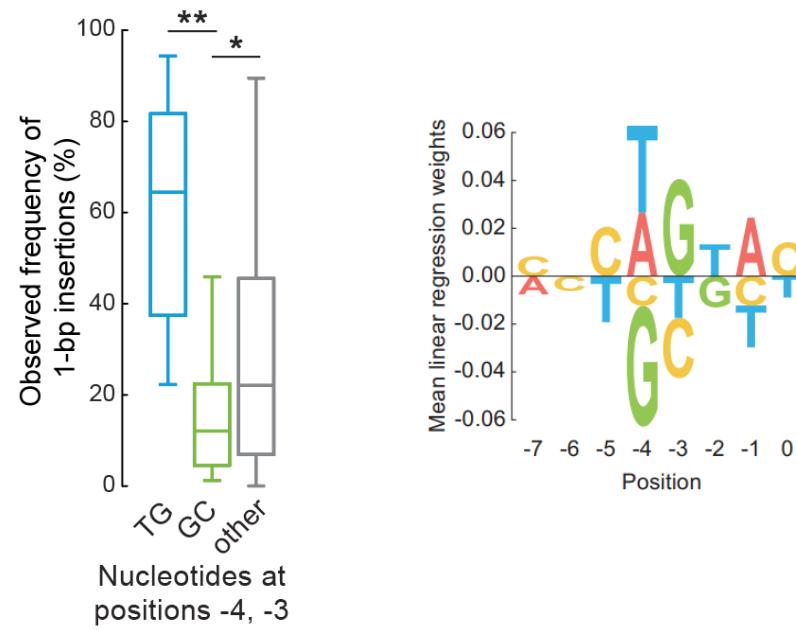
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Target sites yielding a single insertion repair genotype >50% of the time

Weak microhomology {

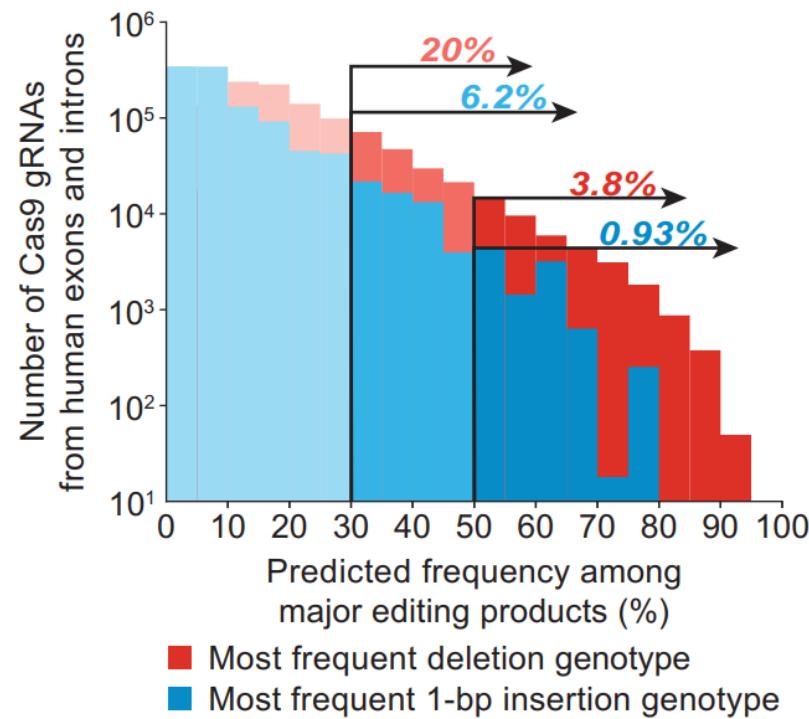
GCCCCCACAAAGCGGACCAGCCAGC NN NN TAGGGAATAAGGCCAACTTGCACCC
AGTGCACGAGTGACGGGATAATGTG NN NN CAGGTAAATTGATAGGCTTCCAACCTTA
CCGAGACCTACCGAACCAGAAATCG NN NN CCGGTGCGCAAGCAGGGCTGGCGTCC

Local sequence context {



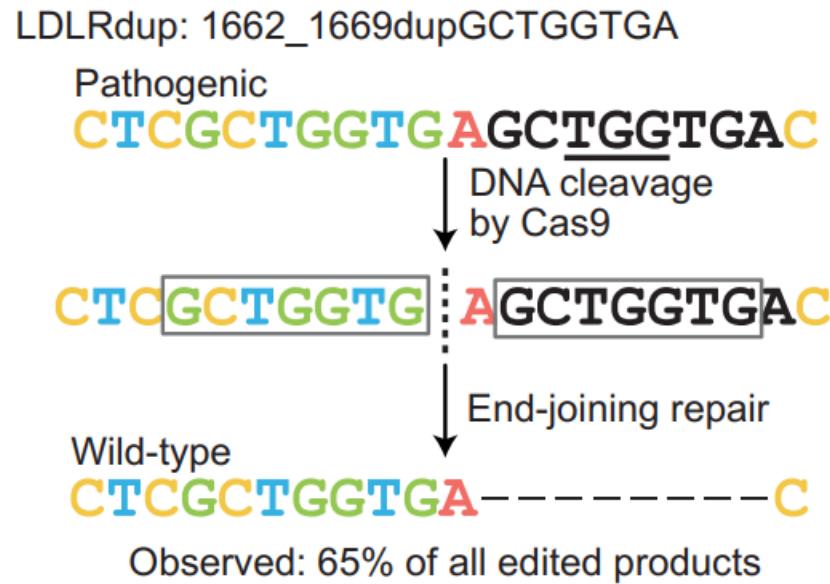
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

inDelphi predicts that 5% of gRNAs yield a single repair genotype the majority of the time



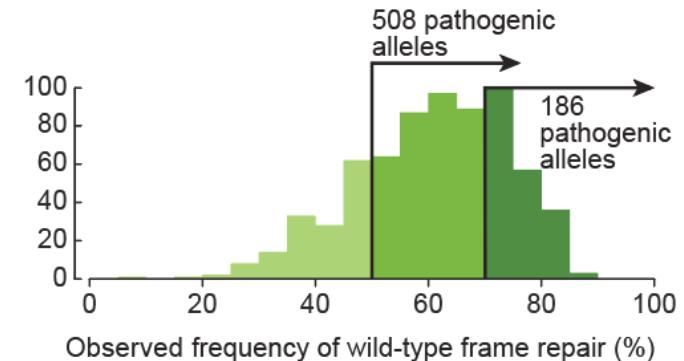
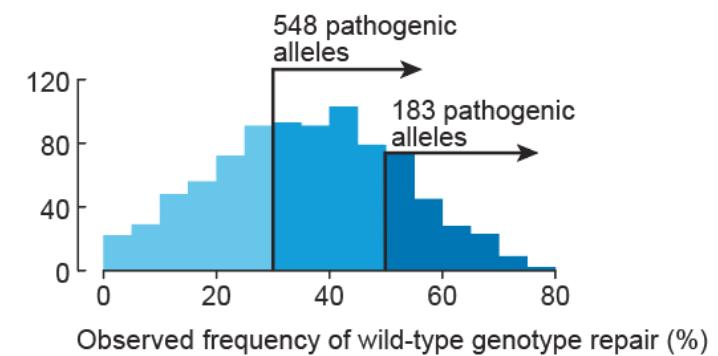
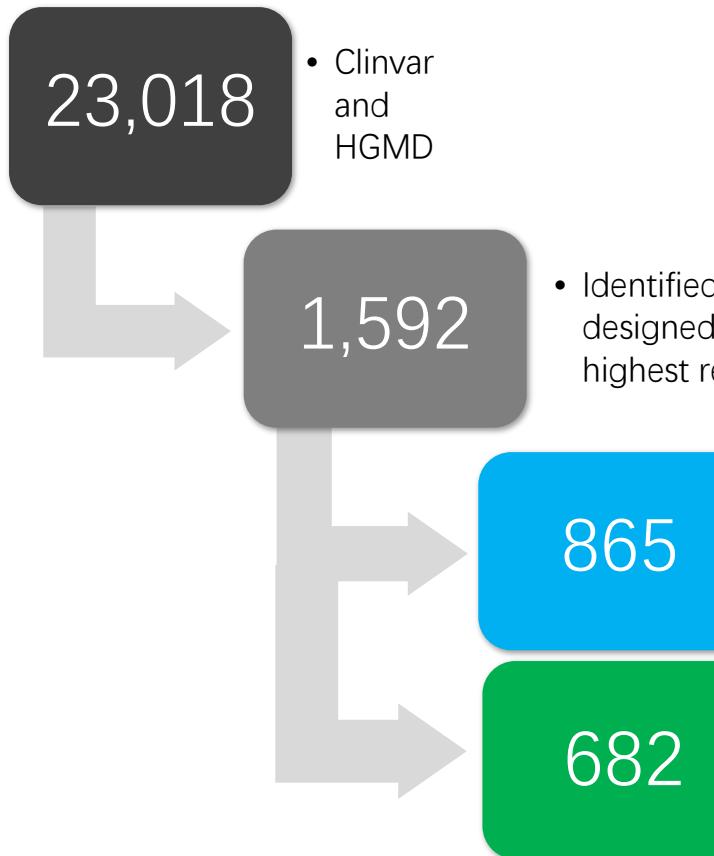
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Pathogenic microduplications are efficiently repairable to wild-type with simple Cas9 cutting



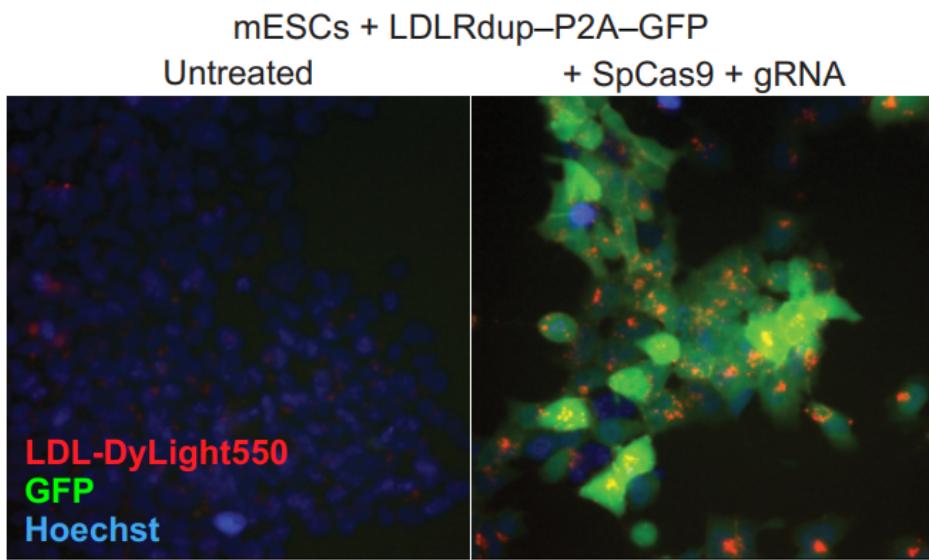
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

inDelphi identified 183 pathogenic alleles corrected to wild-type at >50% frequency ($r = 0.64$)



Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

Efficient repair of pathogenic alleles to wild-type with template-free Cas9-nuclease treatment



- Primary patient-derived fibroblasts
Human and mouse cell lines
- SpCas9 and SaCas9
- HPS1 71%
- LDLR 77%
- PORCN 48%
- GAA 68%
- GLB1 42%

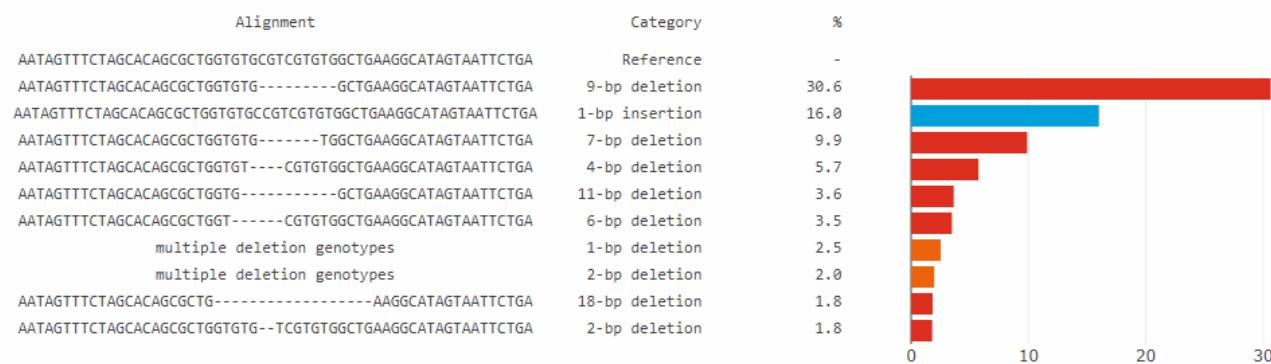
Cas9 editing without a homology template is predictable, can be precise, and can be practical for disease

inDelphi

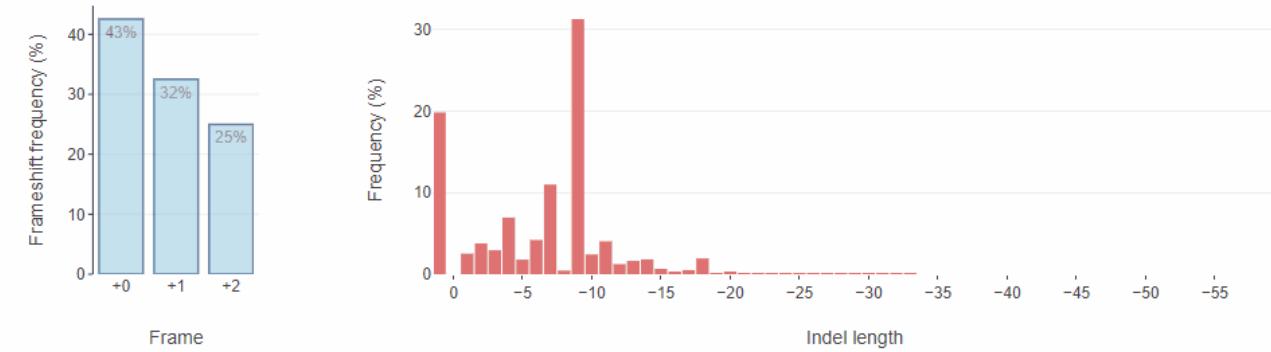
AATAGTTCTAGCACAGCGCTGGTGTGCGCTGTGTGGCTGAAGGCATAGTAATTCTGA GTCGTGTGGCTGAAGGCATAGTAATTCTGA

◀ DSB ▶

Summary of predictions: Top 10 frequent events



Indel length predictions

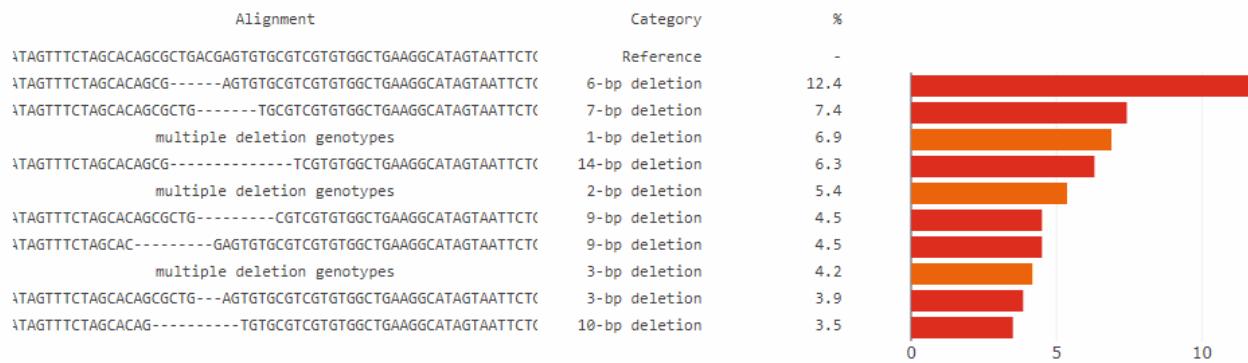


inDelphi

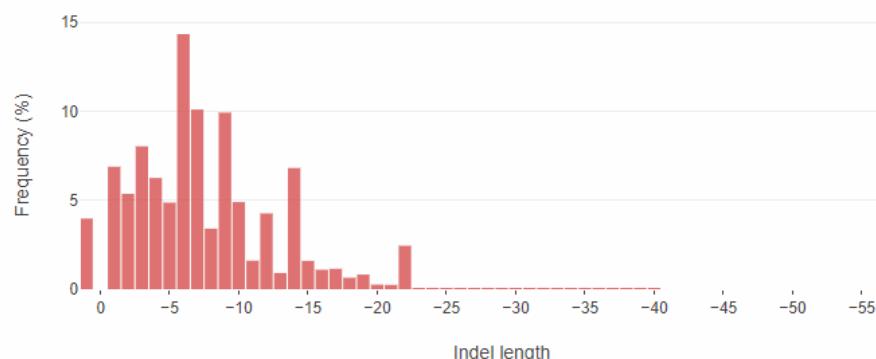
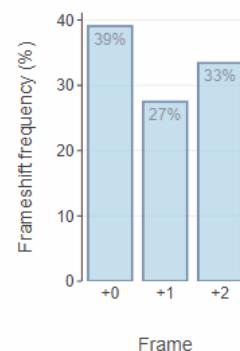
AATAGTTCTAGCACAGCGCTGACGAGTGTCGCTGTGGCTGAAGGCATAGTAATTCTC ACGAGTGTGCGCTCGTGTGGCTGAAGGCATA

◀ DSB ▶

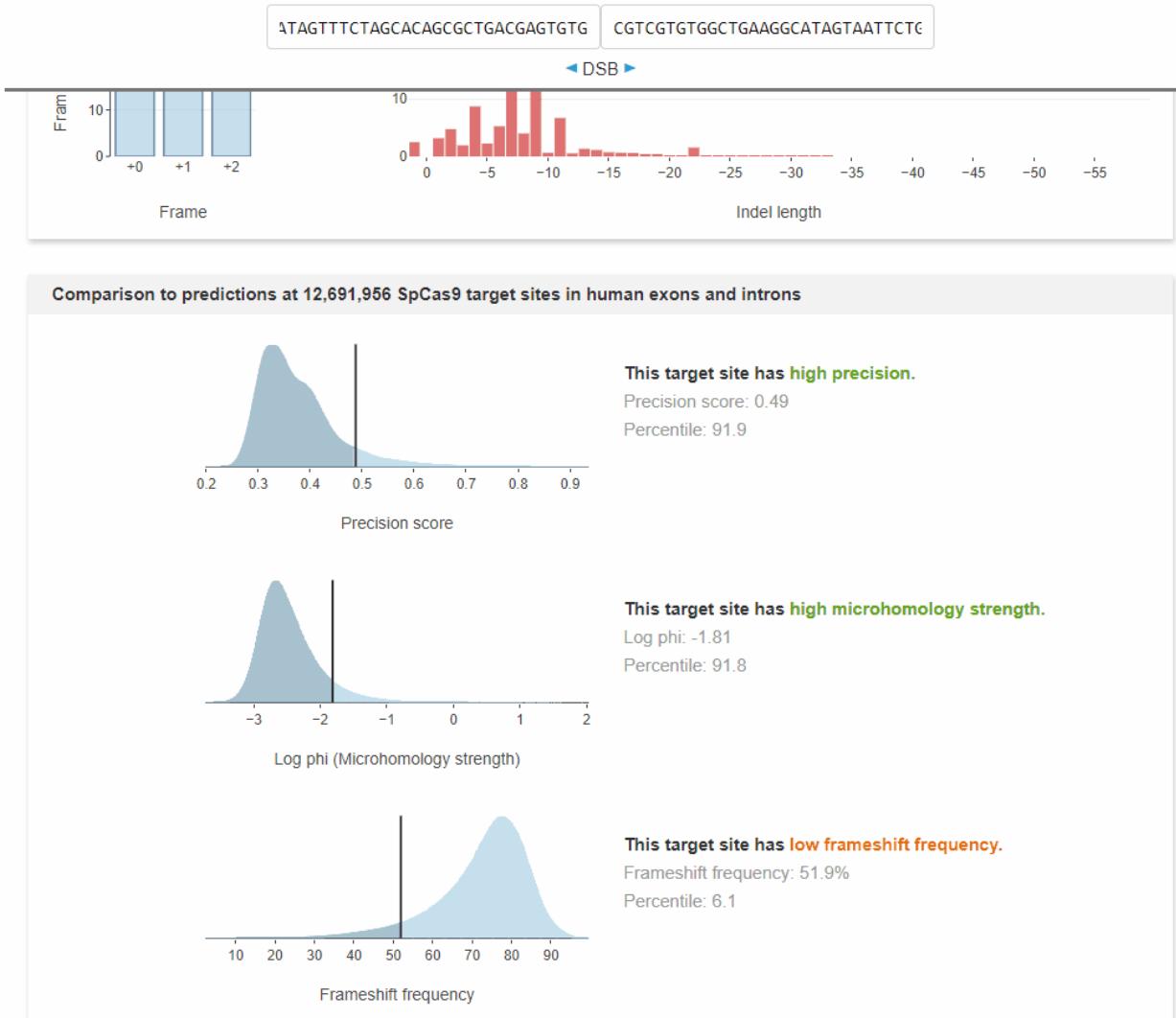
Summary of predictions: Top 10 frequent events



Indel length predictions



inDelphi



Acknowledgements



Massachusetts
Institute of
Technology



CSAIL



Merkin Institute

FOR TRANSFORMATIVE
TECHNOLOGIES IN HEALTHCARE



BWH BRIGHAM AND
WOMEN'S HOSPITAL



Hubrecht
Institute

Max Shen

Mandana Arbab

David R. Liu

Richard I. Sherwood

Daniel Worstell
Olga Krabbe

Christopher A. Cassa



MIT BE

Jonathan Hsu