



Lab in the loop Machine learning for LMDD

Richard Bonneau

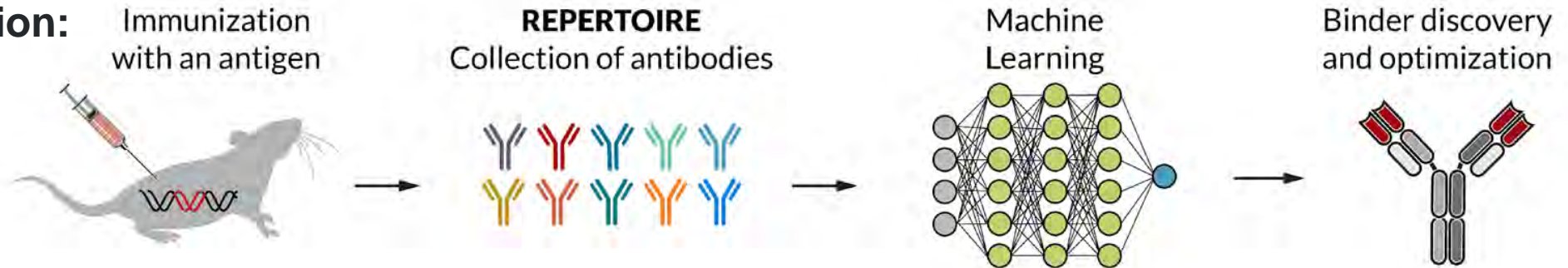
Prescient
Design

A Genentech Accelerator

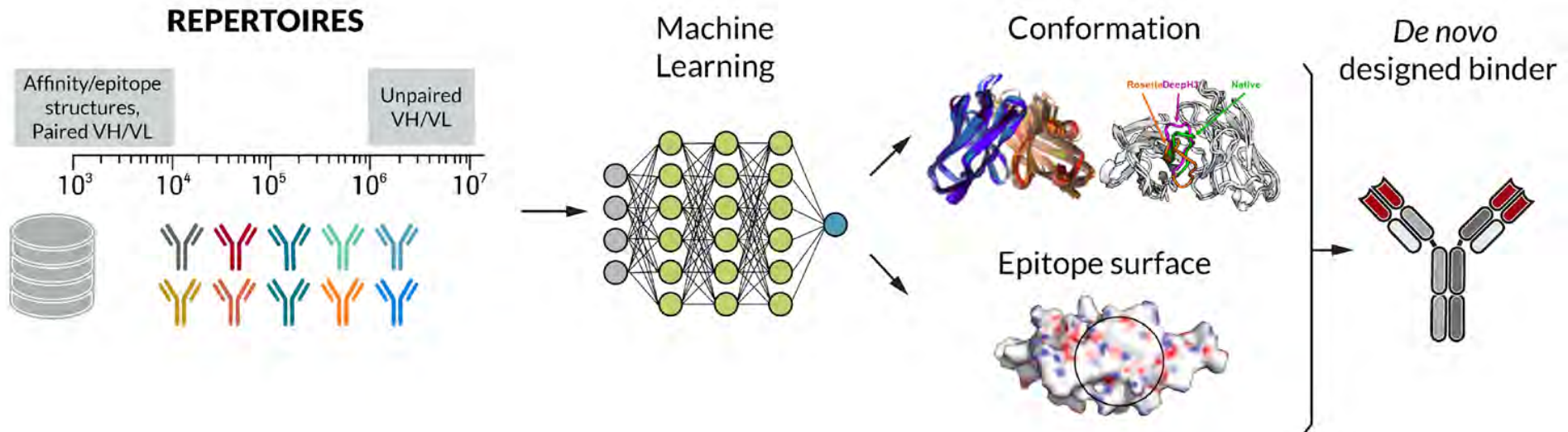
Genentech
A Member of the Roche Group

AI / ML for antibody identification, optimization and *de novo* design

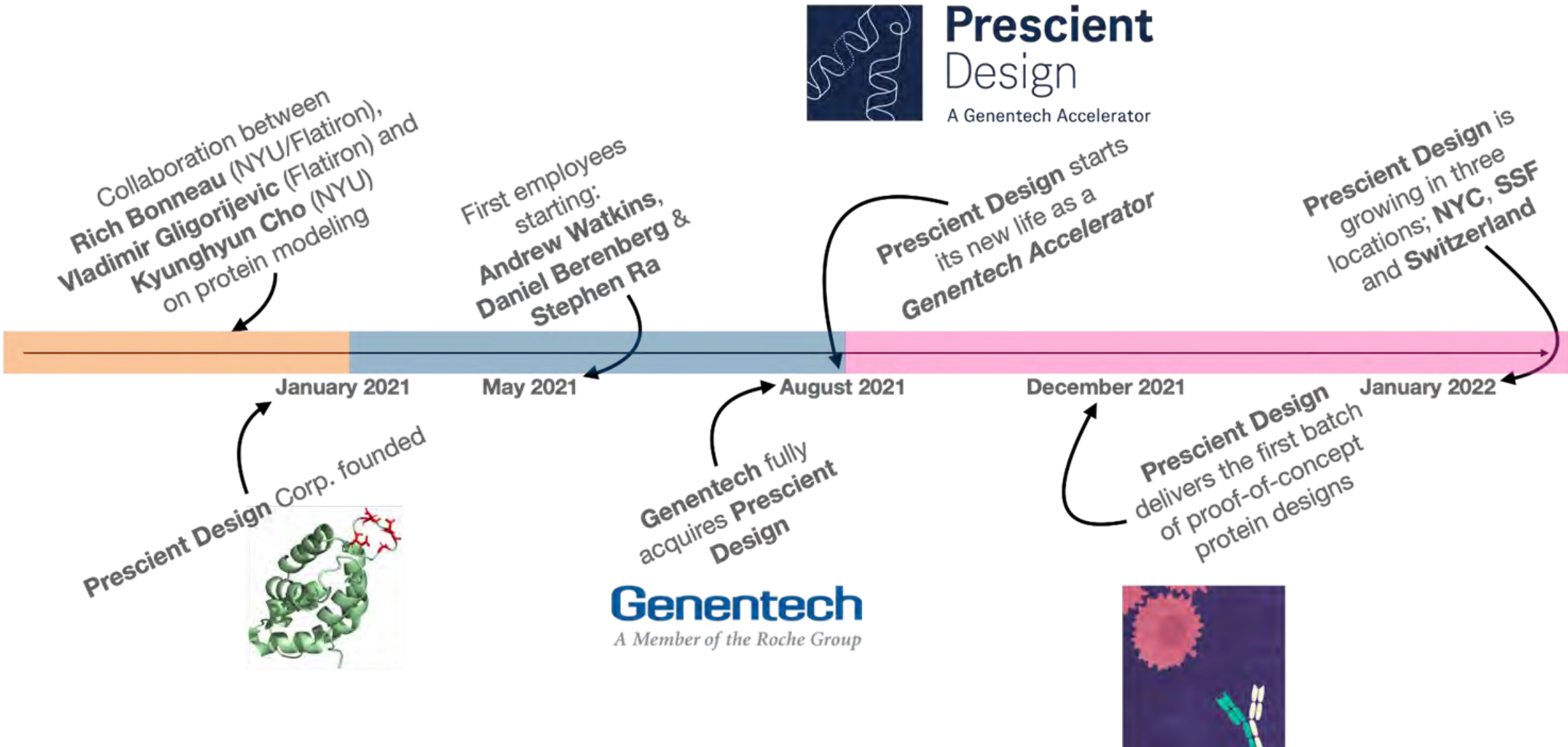
Optimisation:



Direct:



Brief History of Prescient Design



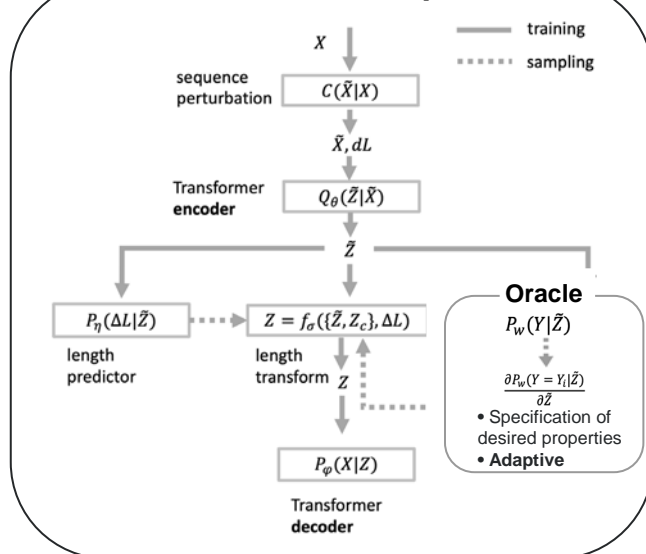
Lab-in-the-Loop Framework

An iterative framework for evidence acquisition, guided design, and validation

Refining the generative process recursively

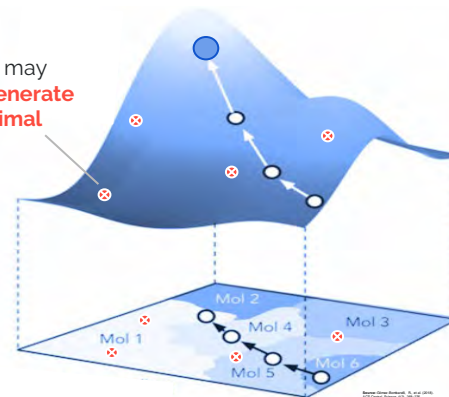
The acquisition of experimental results -- both **positive** and **negative** -- combined with new data helps **adapt** the generative process and enables better optimization of future designs

Manifold Sampler



Guided sampling

Random mutations may yield **degenerate** or **suboptimal** designs

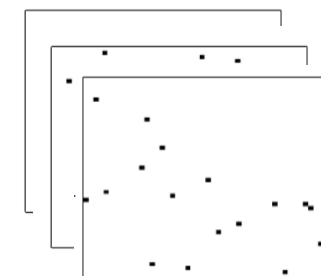


- Gradient-based **multi-objective** guidance
- Active learning & **Uncertainty-awareness**
- **Efficient sampling** of design space

Proposing new designs

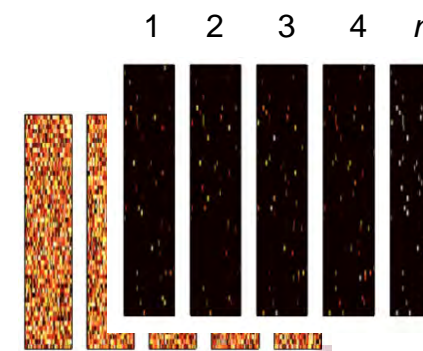
Our sampling procedure finds an **optimal exploration-exploitation strategy** and can query SMEs for experimental validation/synthesis of generated designs

Experimental feedback



- Selection of optimized leads
- Evaluation data for accepted/rejected proposals
- Additional SME feedback

Assay



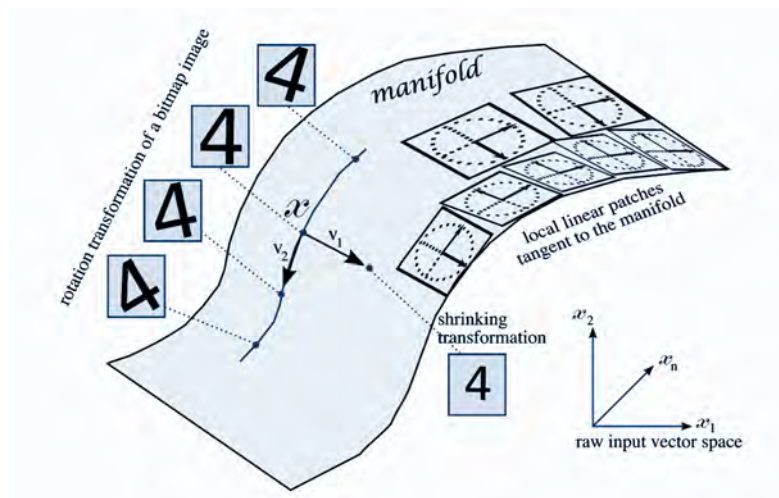
Experimental validation of proposed designs

Generated designs with predicted properties Y

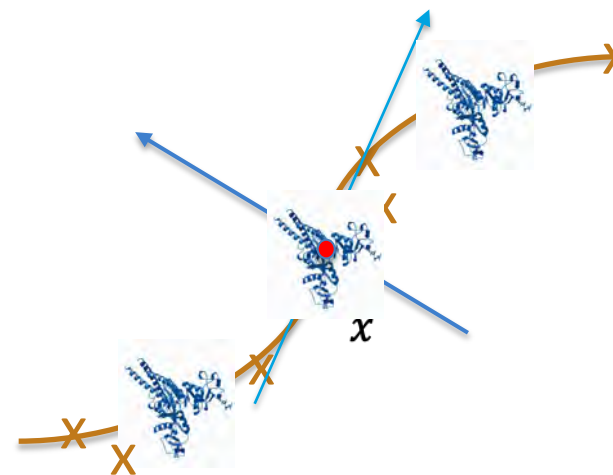
Manifold Sampler

Vladimir Gligorijevic, Dan Berenberg, Stephen Ra, Andy Watkins, Simon Kelow, Rich Bonneau, Kyunghyun Cho

Real-world, high-dimensional data lie (roughly) on a low-dimensional manifold (Chapelle 2006)



Bengio (2009)



Instead of the **(large)** combinatorial space, search for a novel sequence x is restricted on the **(small)** manifold space, i.e., $x \in \mathcal{M}$

Function-guided design: $\operatorname{argmax}_{\mathbf{x} \in \mathcal{M}} f(\mathbf{x})$ for a scoring function $f(\mathbf{x})$ trained on \mathcal{M} where $f(\mathbf{x})$ is an oracle.

Given a protein (target) sequence

$Y = y_1, \dots, y_n$, learn a **Seq2Seq DAE** model

where y_1, \dots, y_n is a corrupted

version of $X \rightarrow Y$

$$Y \quad \mathbf{X} \sim C(\mathbf{X}|\mathbf{Y})$$

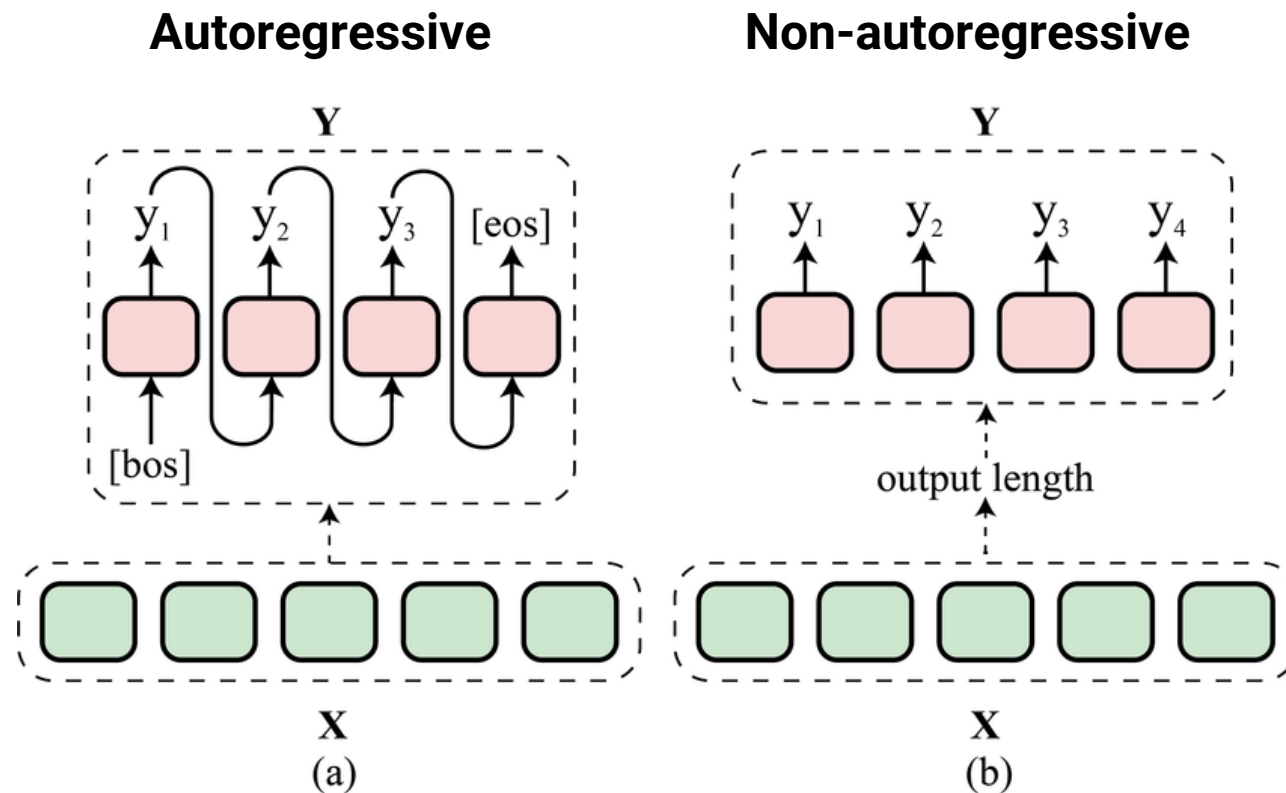
Two approaches:

1. Autoregressive Seq2Seq:

1. Non-autoregressive Seq2Seq:

$$\log p(\mathbf{Y}|\mathbf{X}) = \sum_t \log p(y_t, y_{<t}, \mathbf{X})$$

$$\log p(\mathbf{Y}|\mathbf{X}) = \sum_t \log p(y_t, \mathbf{X})$$



Non-autoregressive sampling. Makes changes in multiple positions of a target sequence enabling **effective** exploration of the overall fitness landscape and resulting in **diverse** sequence designs

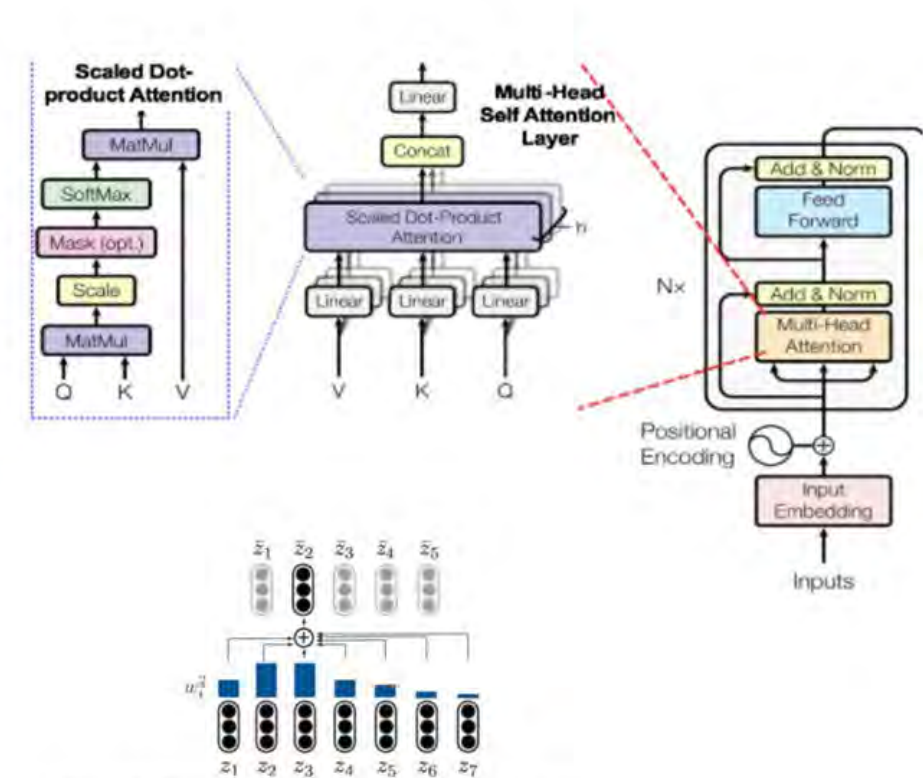
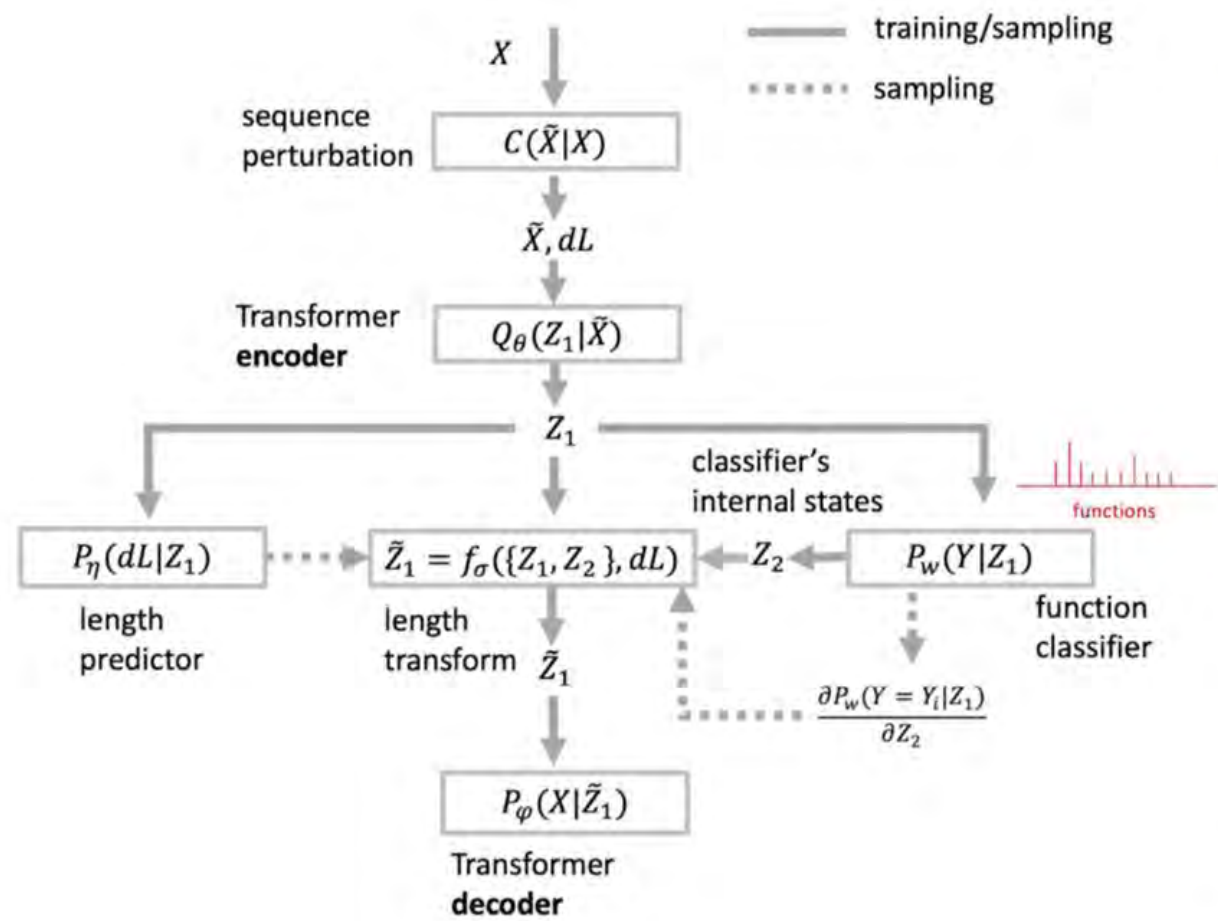


Figure 2: Illustration of the length transformation mechanism.

Shu, Raphael, ..K. Cho "Latent-variable non-autoregressive neural machine translation with deterministic inference using a delta posterior." *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. No. 05. 2020.

Deep manifold sampler identifies the protein manifold by learning to denoise [Gligorijevic et al., 2021]

It thereby trades off between **statistical efficiency** and **combinatorial complexity**.

VLSEGEWQ**LVL**HVWAKVEADVAGHGQ**QV**DILIRLFKSH
PETLEKFDRFKHLKTEAEMK**KATASE**DLKKHGVTVLTA
LGAILKKKGHHEAELK**PLAQSHAT**KHKIPIKYLEFISEAII
HVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELG
Y

1. **Sequence corruption**

2. **Denoising**

VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPE
TLEKFDRFKHLKTEAEMKASEDLKKHGVTVLTALGAILK
KKGHHEAELKPLAQSHATKHKIPIKYLEFISEAIIHVLHS
RHPGDFGADAQGAMNKALELFRKDIAAKYKELGY

Applications to LMDD

More details in our recent publications

Article | [Open Access](#) | [Published: 26 May 2021](#)

Structure-based protein function prediction using graph convolutional networks

[Vladimir Gligorijević](#) ✉, [P. Douglas Renfrew](#), [Tomasz Kosciolk](#), [Julia Koehler Leman](#), [Daniel Berenberg](#), [Tommi Vatanen](#), [Chris Chandler](#), [Bryn C. Taylor](#), [Ian M. Fisk](#), [Hera Vlamakis](#), [Ramnik J. Xavier](#), [Rob Knight](#), [Kyunghyun Cho](#) & [Richard Bonneau](#) ✉

TM-Vec: template modeling vectors for fast homology detection and alignment

Tymor Hamamsy^{1,*}, James T. Morton^{2,3,*}, Daniel Berenberg^{4,9}, Nicholas Carriero⁵, Vladimir Gligorijević⁹, Robert Blackwell⁵, Charlie E. M. Strauss⁷, Julia Koehler Leman², Kyunghyun Cho^{1,4,8,9,†}, and Richard Bonneau^{1,4,6,9,†}

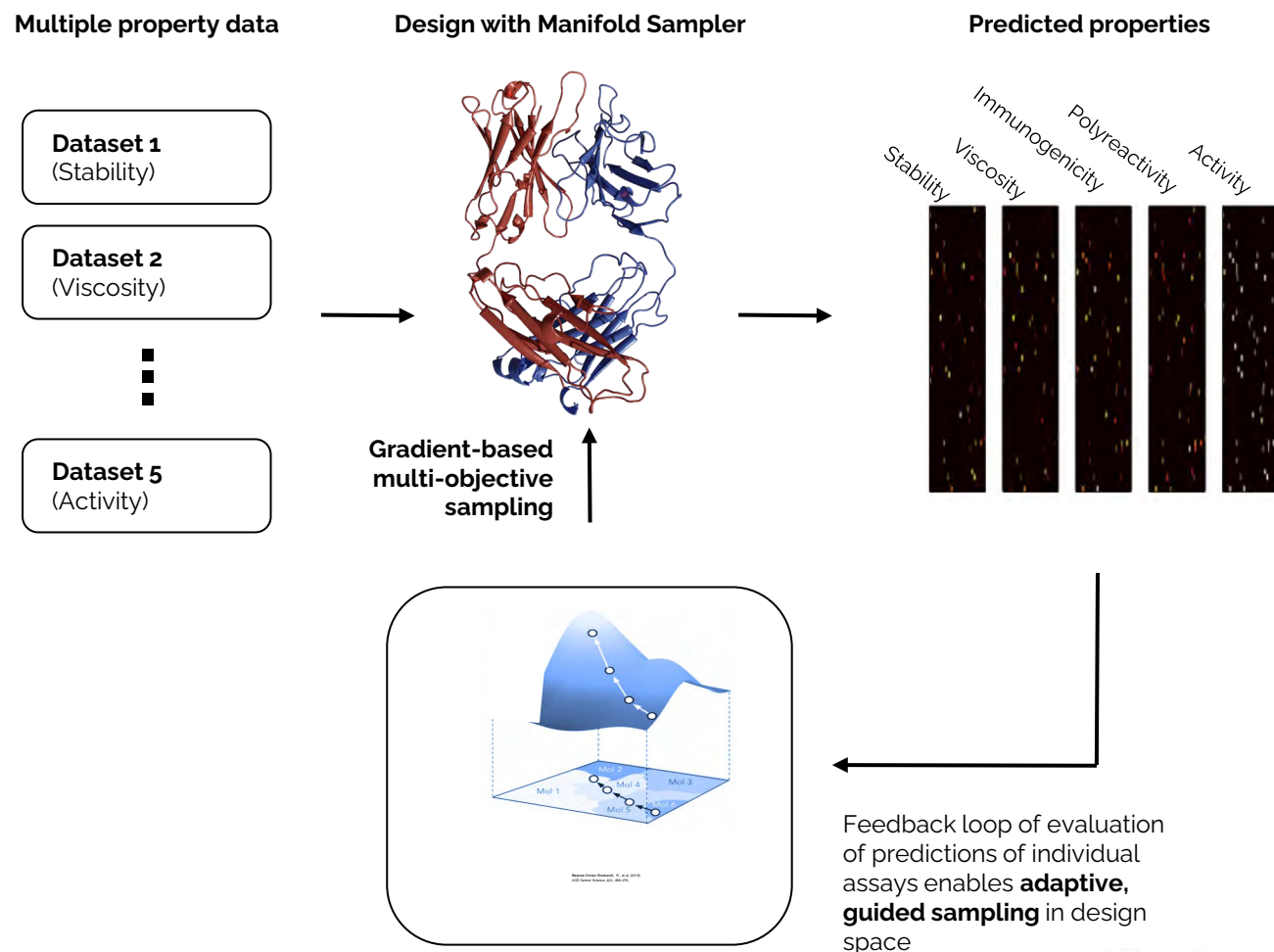
Function-guided protein design by deep manifold sampling

Vladimir Gligorijević¹, Daniel Berenberg^{1,2}, Stephen Ra¹, Andrew Watkins¹, Simon Kelow¹, Kyunghyun Cho^{1,2,3}, and Richard Bonneau¹

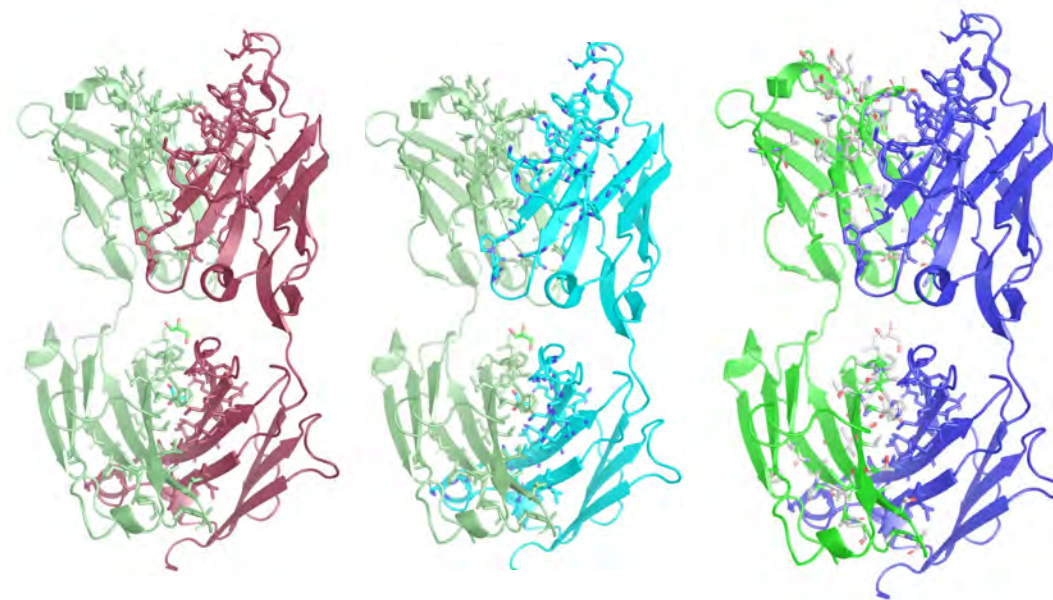
Multi-segment preserving sampling for deep manifold sampler

Daniel Berenberg^{1,2}, Jae Hyeon Lee¹, Simon Kelow¹, Ji Won Park¹, Andrew Watkins¹, Vladimir Gligorijević¹, Richard Bonneau¹, Stephen Ra¹, and Kyunghyun Cho^{1,2,3,4}

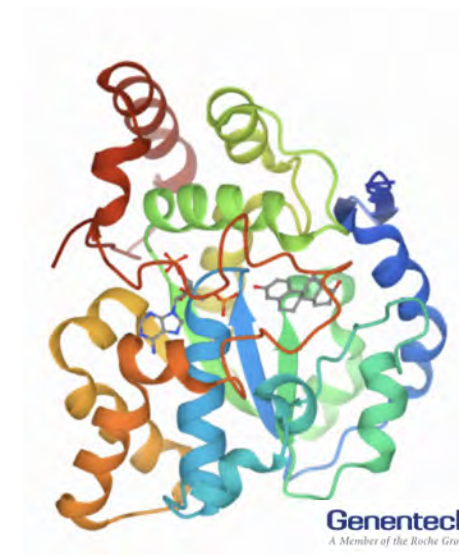
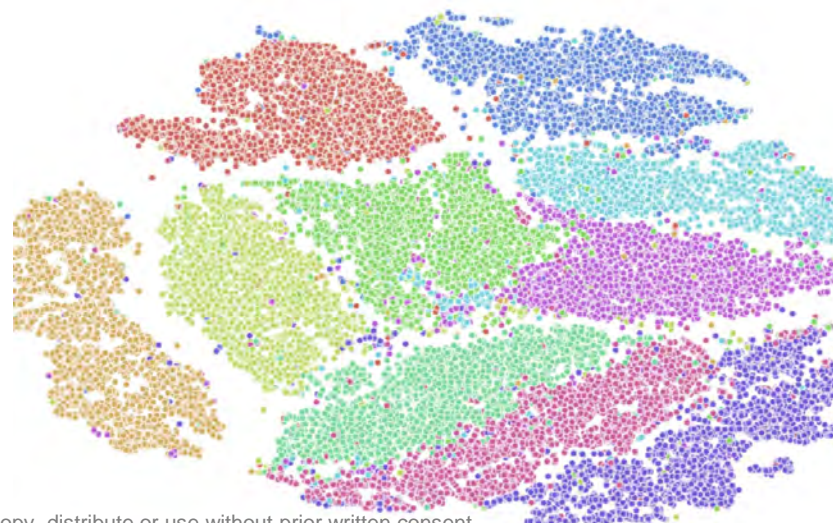
- **Parallel optimization of multiple properties**
 - currently optimized sequentially
 - e.g. affinity, polyreactivity, viscosity
- **Increased yield**
 - increased number of candidate sequences that meet criteria
- **Better optimisation**
 - optimizing property n will not degrade property $n-1$ as it might with sequential optimisation
- **Meta-learning multiple properties**
 - robustness for small training datasets
- **Manifold Sampler naturally models epistasis**



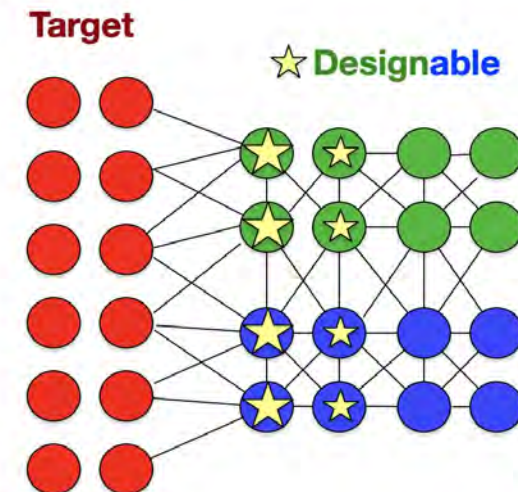
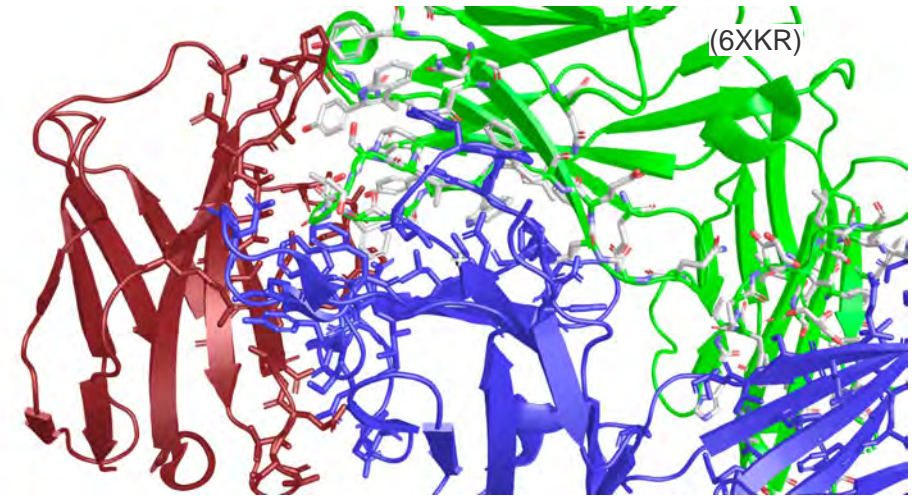
- Sufficient data on VH-VL pairing across multiple immunizations
- High quality negative examples (important for our approach) can be generated by pairing VH and VL across immunizations to disparate targets.
- Chain pairing will feed into better epitope clustering and identification of rare epitopes.
- Sorting paratopes into clusters corresponding to epitopes will be done in manifold representation.
- Representatives from clusters or rare epitopes feed directly into design



Paratopes clustered on manifold



- Given any target protein or interface the problem is well constrained.
- We first set up a graph representing the desired interaction to enable geometric deep learning.
- The Manifold Sampler is then used to power designs on that graph.
- Antibody (and protein-protein) interfaces are well represented in the PDB, and data to begin training this model is in-hand.
- Validation with structure determination and/or characterisation of binding feeds back on model (active learning).



Wide application:

2 development targets, Ova and Her. 6 'hot' targets, 4 'medium' targets
Applications to other Fab like formats. Application to non antibody affinity.

Active learning is working (not a disaster ;-)

Improvement in affinity and **multiple developability parameters** for 5 of 6 initial targets.

Mutations/Edit-distances of 5-65 ... expression, binding and good dev for some molecules with > 50 mutations. True global design.

New high throughput lab techniques will drive even better approaches to active learning.

Good early readouts with 'real' *de novo* and grafting *de novo*. Much work remains.

Team

Founders



Richard Bonneau

Executive Director of Prescient
Design



Vladimir Gligorijević

Senior Director of AI/ML



Kyunghyun Cho

Senior Director of Frontier
Research

Even mix of:

- Structure Biology
- ML for DD
- Frontiers ML
- Engineering



Stephen Ra

Director of Frontier Research



Andrew Watkins

Structural and Computational
Biology Lead



Daniel Berenberg

Machine Learning Engineer
(ML)



Simon Kelow

Structural and Computational
Biologist (BIO)



Jae Hyeon Lee

Machine Learning Scientist
(ML)



Ji Won Park

Machine Learning Scientist (FR)



Camille Alexander-Norrell

Team Administrator



Maria Lee

Infrastructure Engineer (ENG)



Natasa Tagasovska

Machine Learning Scientist (FR)



Andrea Loukas

Machine Learning Lead and
Principal Scientist



Santrupti Nerli

Structural and Computational
Biologist (BIO)



Henri Dwyer

Engineering Lead



Jack Maguire

Computational Biologist (BIO)



Michael Maser

Computational Chemist (BIO)



Andrew Leaver-Fay

Molecular Modeling Software
Lead (BIO)



Darcy Davidson

Structural and Computational
Biologist (BIO)



Pranav Khade

Postdoctoral Fellow (BIO)



Samuel Stanton

Machine Learning Scientist
(Frontier Research)



Nathan Frey

Machine Learning Scientist
(ML)



Ed Wagstaff

Machine Learning Scientist
(ML)



Jan Ludwiczak

Machine Learning Scientist
(ML)



Omar Mahmood

Machine Learning Scientist
(ML)



Stefani Vasilaki

Senior Machine Learning
Research Associate (ML)



Joshua Yao-Yu Lin

Postdoctoral Fellow (FR)



Franziska Seeger

Principal Scientist, Molecular
Modeling



Sai Pooja Mahajan

Senior Machine Learning
Scientist (BIO)



Aya Ismail

Machine Learning Scientist (FR)



Saeed Saremi

Senior Machine Learning
Scientist (FR)



Vishnu Sresht

Director, Machine Learning
and Computational Chemistry



Tyler Bryson

Data Engineer (ENG)



Edith Lee

Data Engineer (ENG)



Henry Isaacson

Software Engineer (ENG)



Karina Zadorozhny

Machine Learning Engineer
(ENG)



Jaewoo Park

Software Engineer (ENG)