

Bayesian Optimization of Antibodies with a Generative Model of Evolving Sequences

Alan N Amin, *Nate Gruver, *Yucen Lily L, *Yilun Kuang, Hunter Elliott, Calvin McCarter, Aniruddh Raghu, Peyton Greenside, Andrew G Wilson



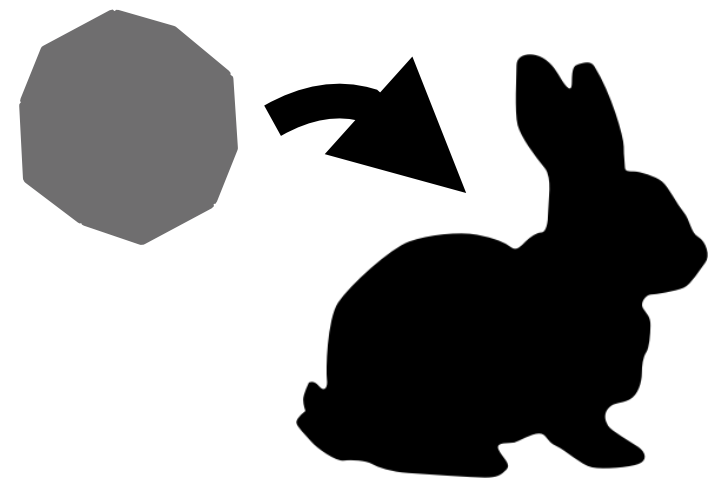
NEW YORK UNIVERSITY



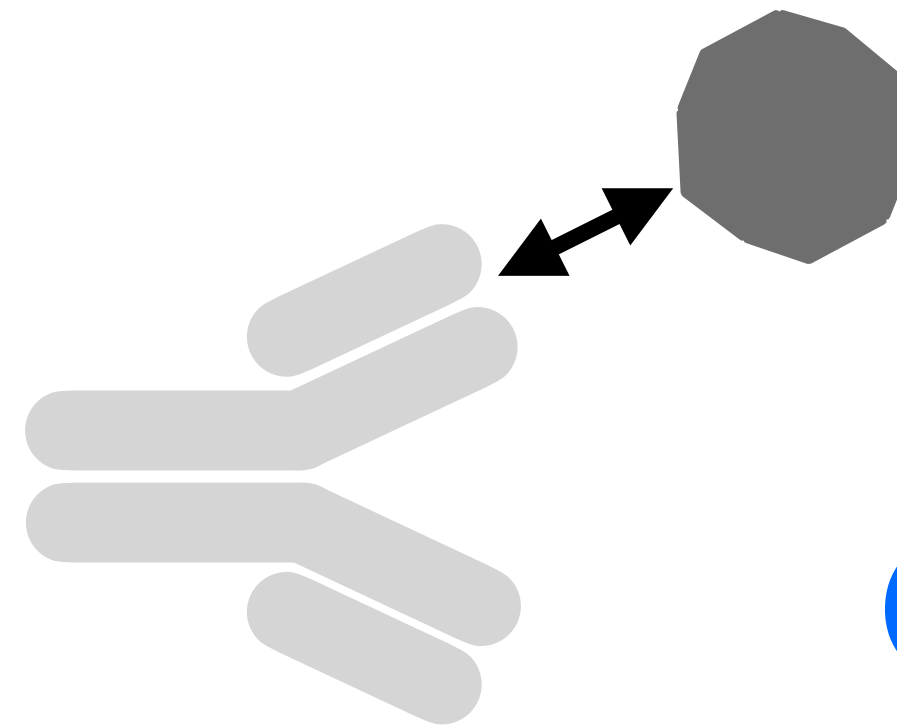
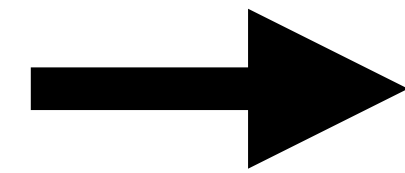
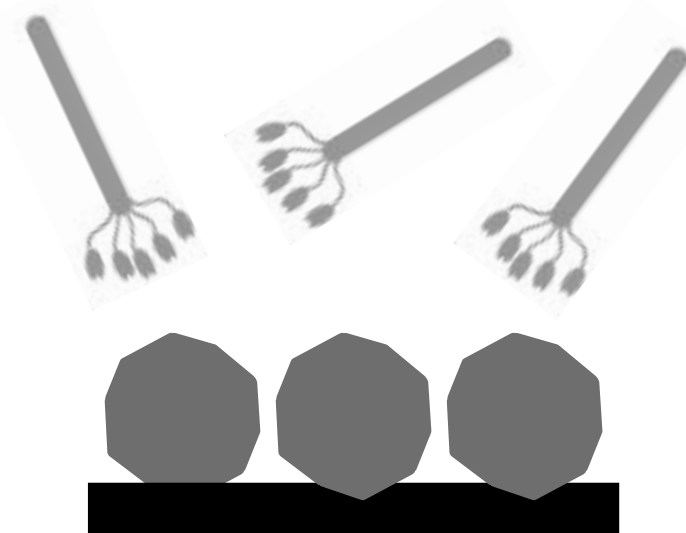
BigHat
BIOSCIENCES

To build antibody drugs, we need to optimize “hits” to be strong binders that are stable in the human body

Animal models:

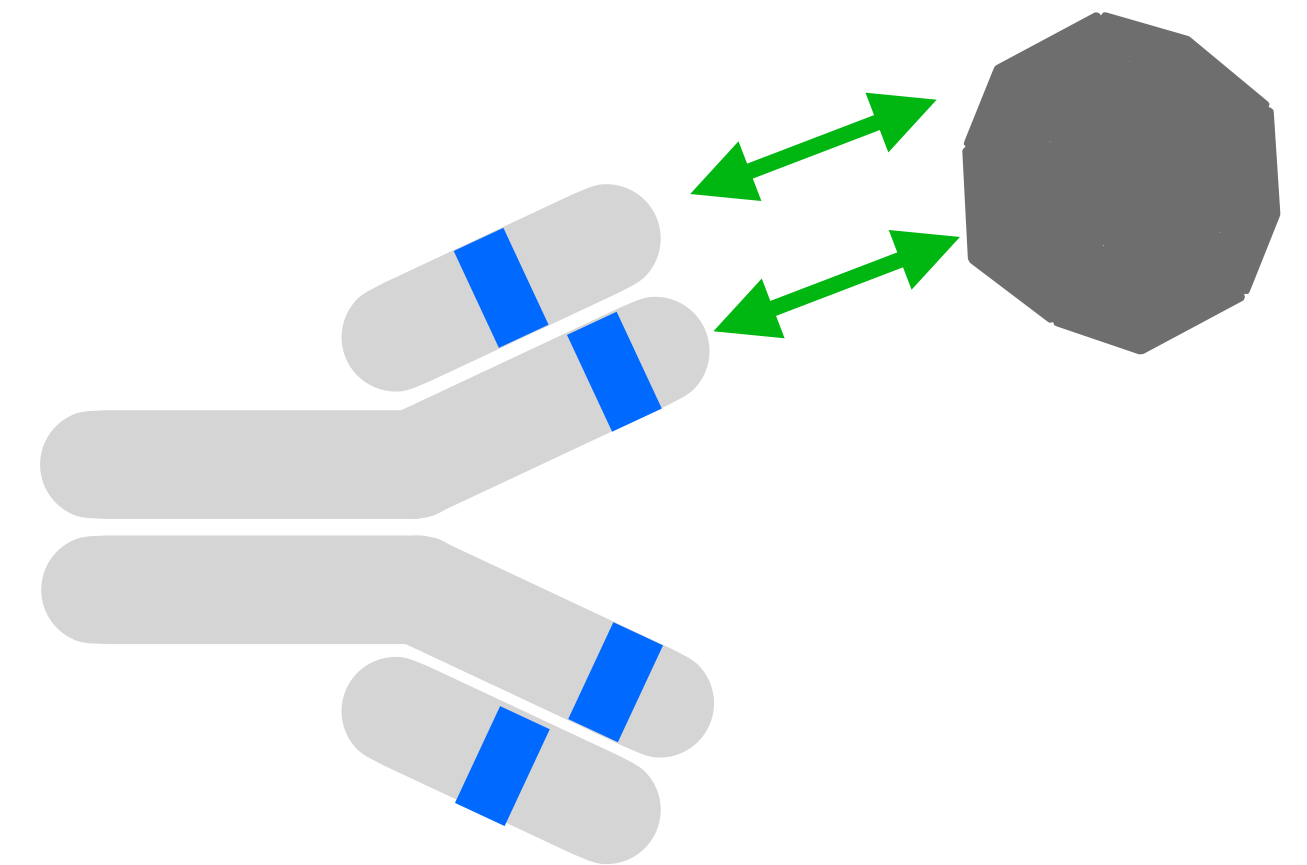
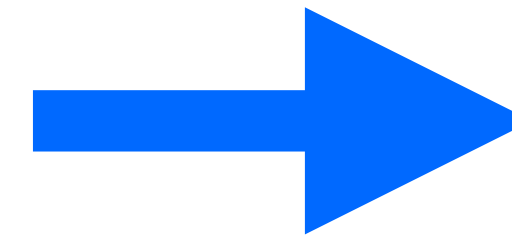


In vitro screens:



Starting
candidate

Optimize



Often:
Low affinity
Polyreactive
Low T_m

Optimization by iterative design is hard because most of the many mutations we can test do not help

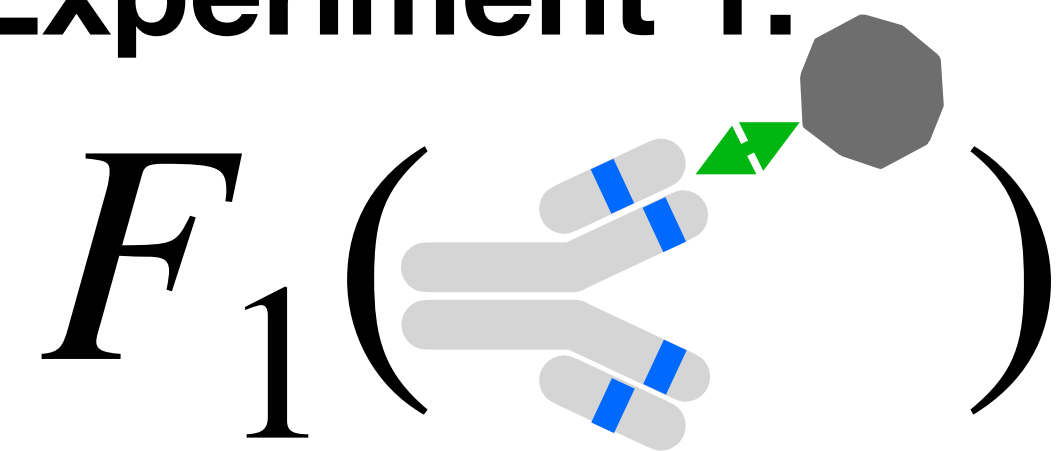
X_0

QVQLRESGPGLVKPSQTL~~SL~~TCTVSGGSFNSGGYYWNRIRQHPGNGL~~EW~~IGYMYSGSTYYNPFIRSRV~~II~~SGDTSVNHFS~~KL~~SSVTAADTAVYFC~~ARGYRQSGYSSWVVDY~~WGQGT~~LV~~NVSS

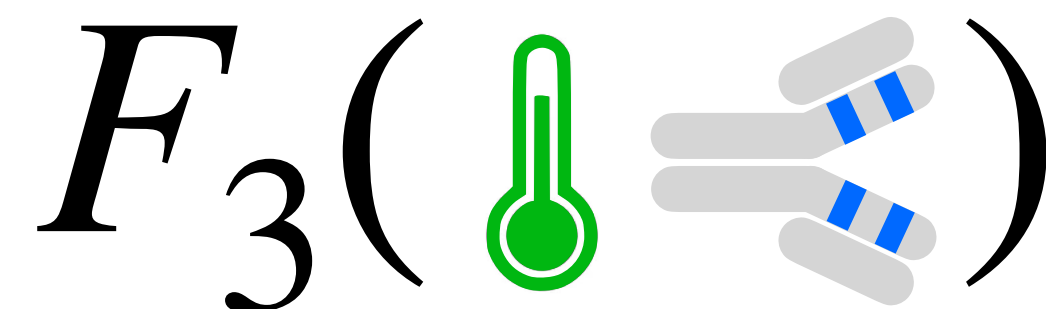
- **Goal:** iteratively suggest and measure X_1, X_2, \dots, X_{100} to have no immunogenicity, have high affinity, or have high melting temperature, $F(X_N)$
- **Possible strategies:** random mutations, avoid mutations that don't often appear in humans
- **Challenge:** search space is huge

Ideally we could build a prior on the objectives $F(X)$ of binding strength and stability

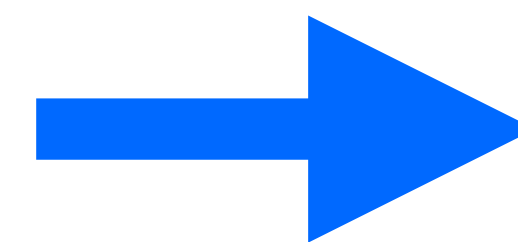
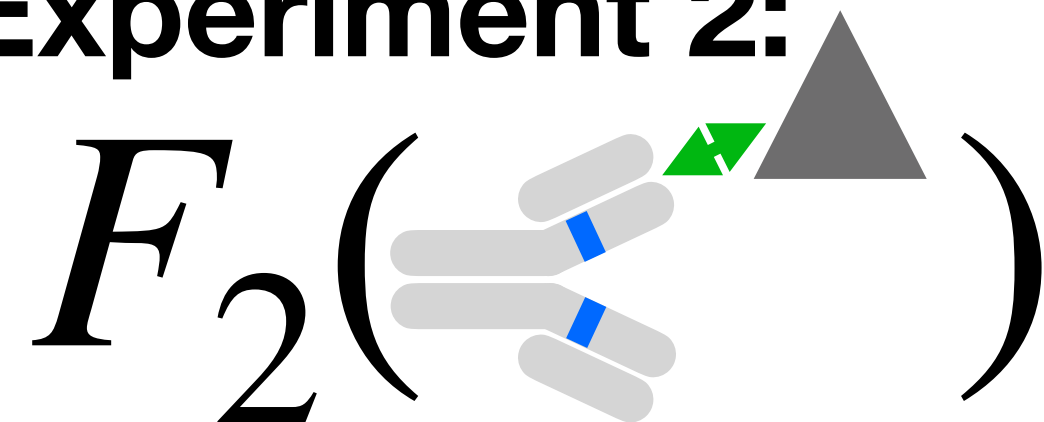
Experiment 1:



Experiment 3:



Experiment 2:

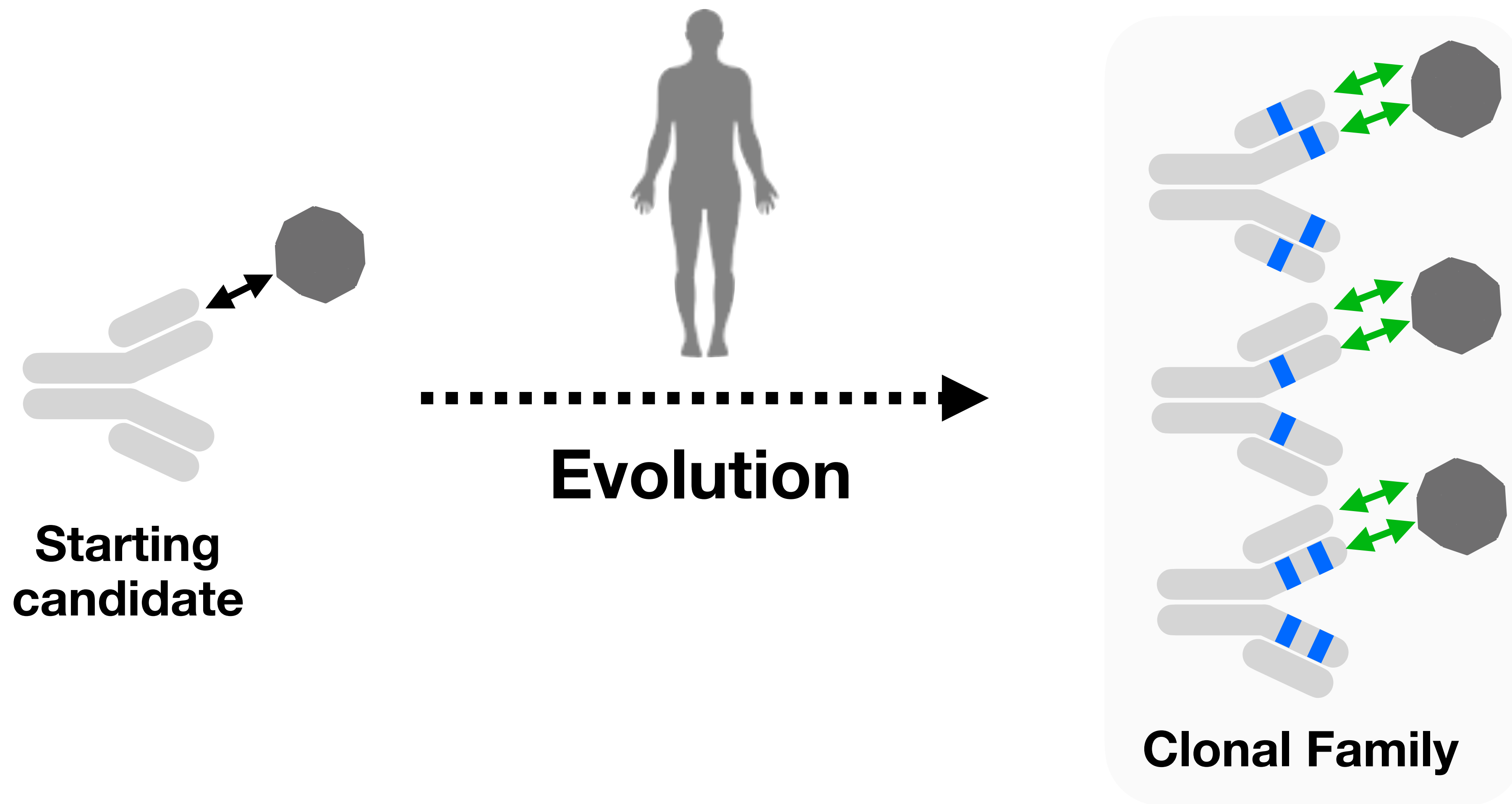


$p(F)$

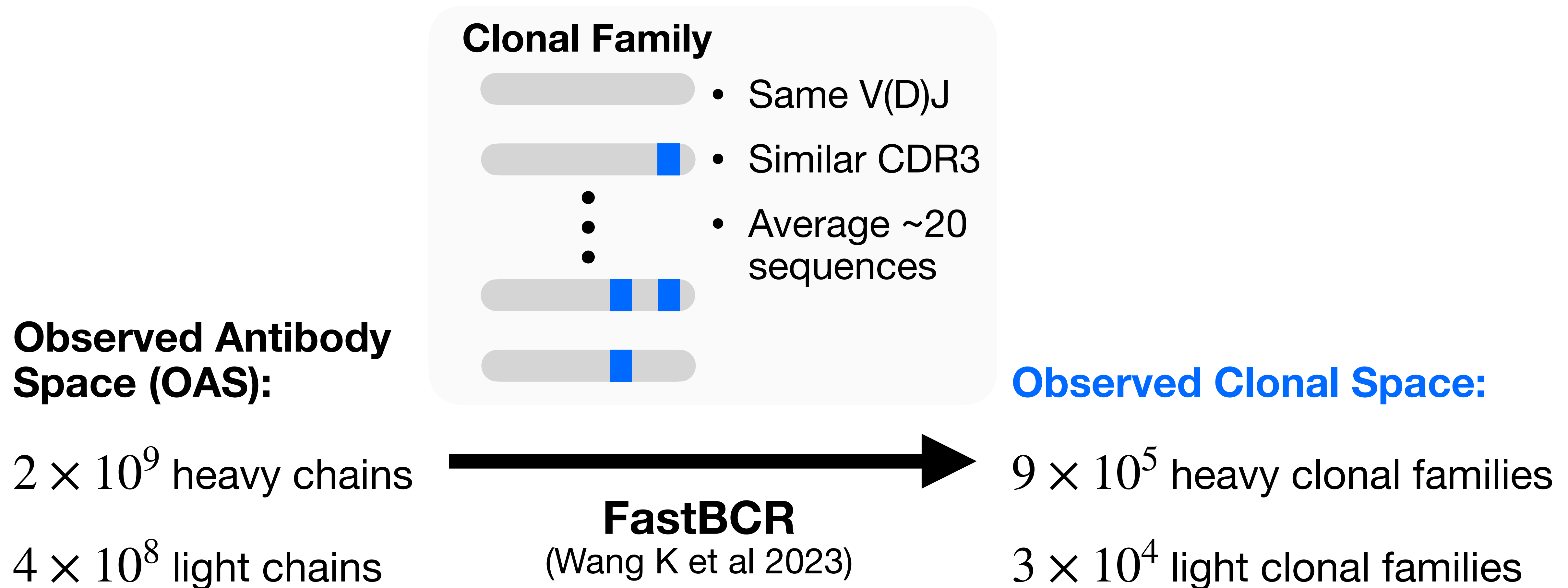
Optimal design strategy: suggest X_{N+1} based on $p(F | X_1, \dots, X_N)$

But we don't have this data!

In principle, we can learn from how our body builds strong and stable binders

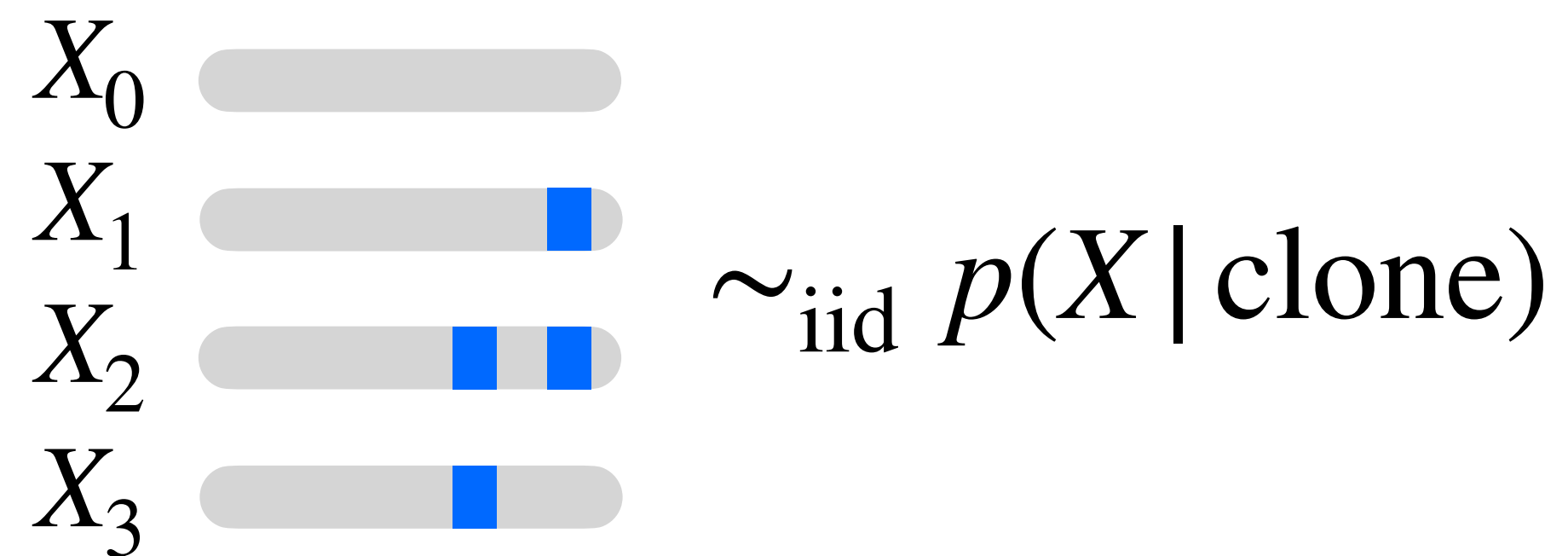


In principle, we can learn from massive data about human clonal families in the OAS database



In theory, we can build a prior over F by looking at the abundance of sequences in each clone

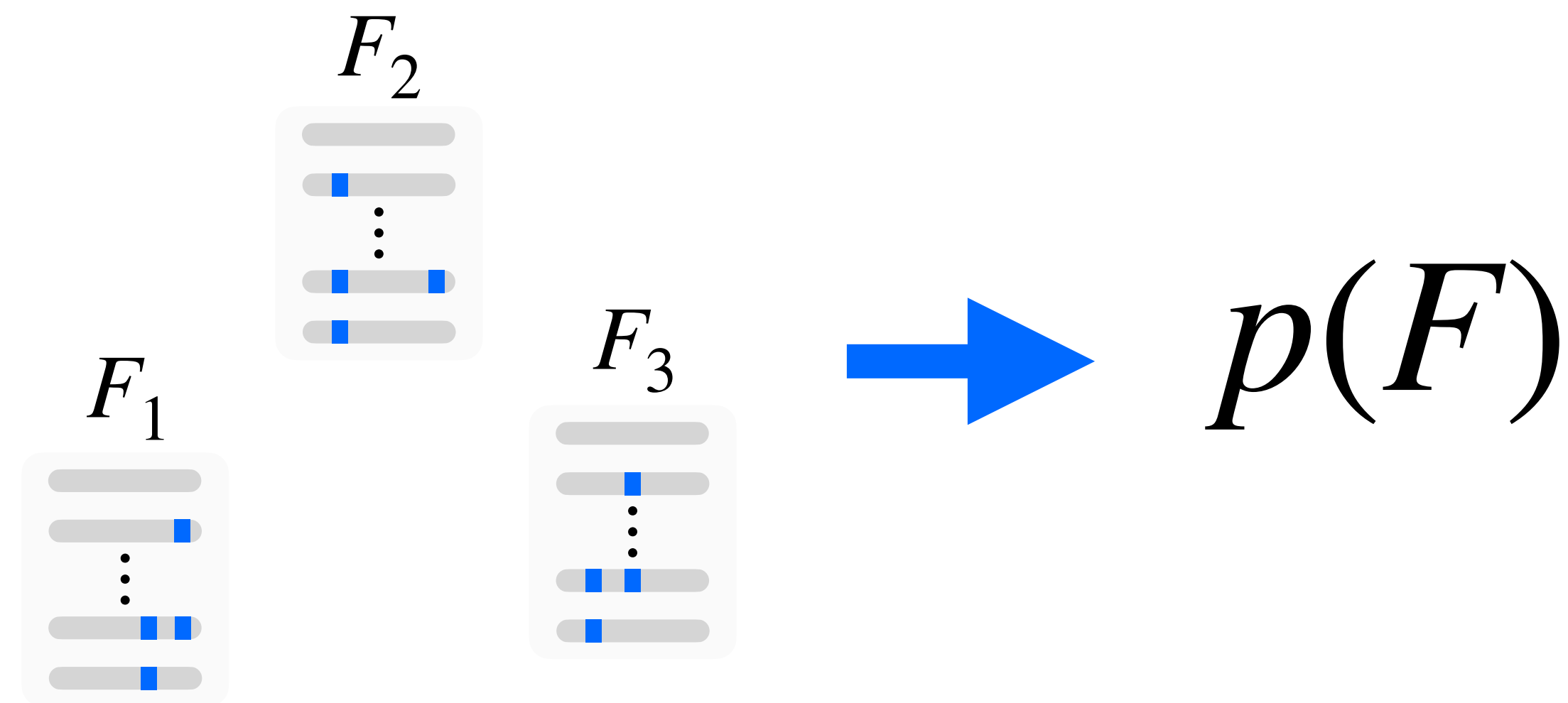
Distribution of sequences in a clonal family:



Stronger, more stable binders are more abundant:

$$p(X | \text{clone}) = \text{Fitness}(X) =: F(X)$$

(Like protein families!)



CloneLM learns the distribution of clonal families

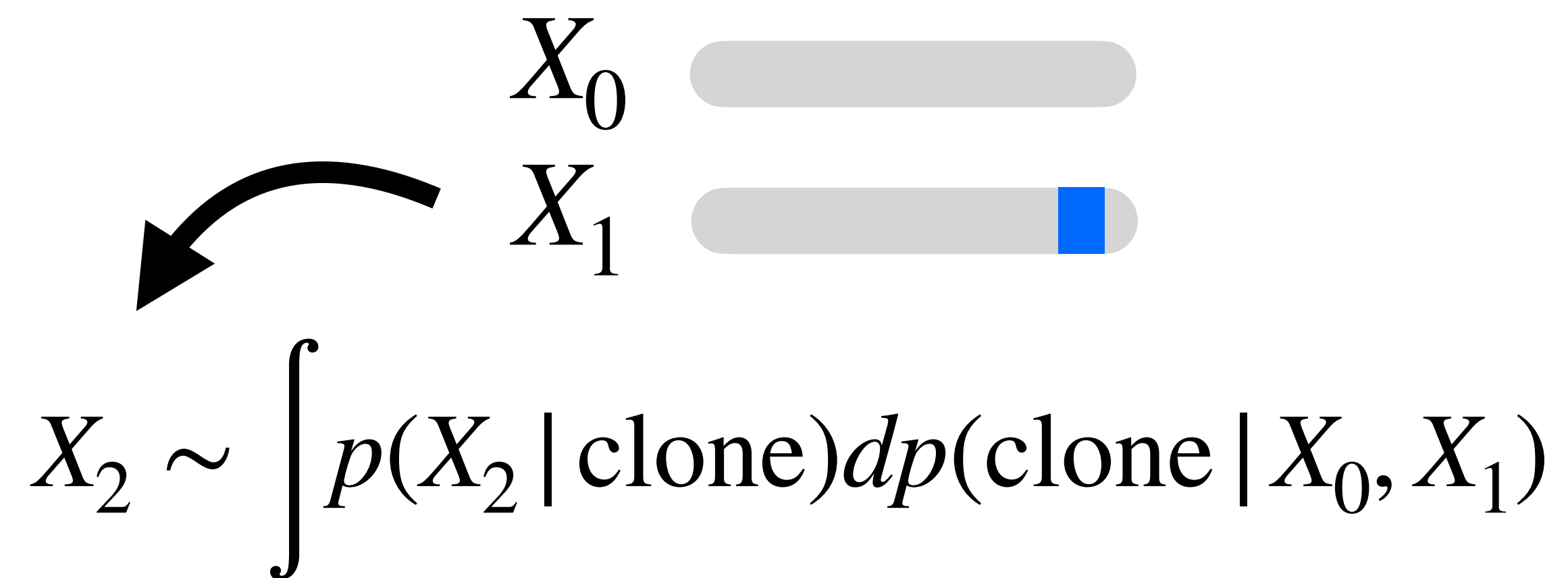
CloneLM (400 M transformer) trained on:

seq1<separator>seq2<separator>seq3<separator>...

Clonal family	<div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNRIRQHPGNGLWIGYMYSSGSTYYNPFIRSRVVISGDTSVNHFSLKLSSVTAADTAVYFCARGYRQSGYSSSSWVDYWGQGTLVNVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSINSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPFILRSRVIIISADTSENHFSRKLSYVTAADTAVYFCARGYRQSGNSSSWVFDYWGQGTLVNVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSINSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPFILRSRVIIISADTSENHFSRKLSVTAADTAVYFCARGYRQSGYSSSWVFDYWGQGTLVNVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSINSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPFILSSRLIISADTPENHFSRRLLSSVTAADTAVYFCATGYPQSGYSSSWVFNWYGQGTLVNVSS</div>
Sample 1	<div>QVQLQESGPRLVKPSQTLSLTCTVSGGSLSNSGGYYWSWFRQPPGKRLEWIGYMYHTGNTYYNPSLKCRVTISGDTSKSHFPLRLTAVTAADTAAYYCARGYRQGGYSSSWLADYGGQGTLGADSS</div> <div>QVQLQESGPRLVKPSQTLSLTCTVSGGSLSNSGGYYWGWIRQPPGKGLEWIGYMYHTGNTYYNPSLKSRTVTISGDTSKNHFSRLRTSVTAADTAVYYCARGYRQGSYSSSWLADYWGQGTLLVTVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPSILRSRVITISGDTSVNPFSLKLSSVTAADTAVYFCARGYRHSGYSSSLLVDYWAEEETVNVSS</div>
Sample 2	<div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPSILRSRVIIISGDTSENQFSLKLSSVTAADTAVYLCPRGYRQSCYSSSWVFDYWGQGTLLVTVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPSILRSRVIIISGDTSENHFSRLKLSSVTAADTAVYFCARGYRQSGYSSSWVLDYWGQGTLLVTVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGKGLEWIGYMYSSGSTYYNPSILRSRVIIISGDTSENHFSRLKLSSVTAADTAVYFCARGYRQSGYSSASWVFDYWGQGTLLVTVSS</div>
Sample 3	<div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGNGLWIGYMYSSGSTYYNPFILKSRVVISGDTSVTHFSLKLSSVTAADTAVYFCARGYRQSGSSSSWVIDYWGQGTLLVTVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHPGNGLWIGYMYSSGSTYYNPFILMSRVIIIRGETSVKHFSRLKLRSVTAADTAVYFCARGYSQSGYSSSWVIDYWGQGTLLVTVSS</div> <div>QVQLRESGPGLVKPSQTLSLTCTVSGGSFNSSGGYYWNWIRQHQQDGLWIGYLYSSGSTYYNPFVKRRVVISGDKSVNHFSLKLSSVTAADTDVYFCARGYGQSGYSSAWVIDYWGQGTLLVTVSS</div>

In theory, CloneLM performs approximate Bayesian inference over evolutionary landscapes

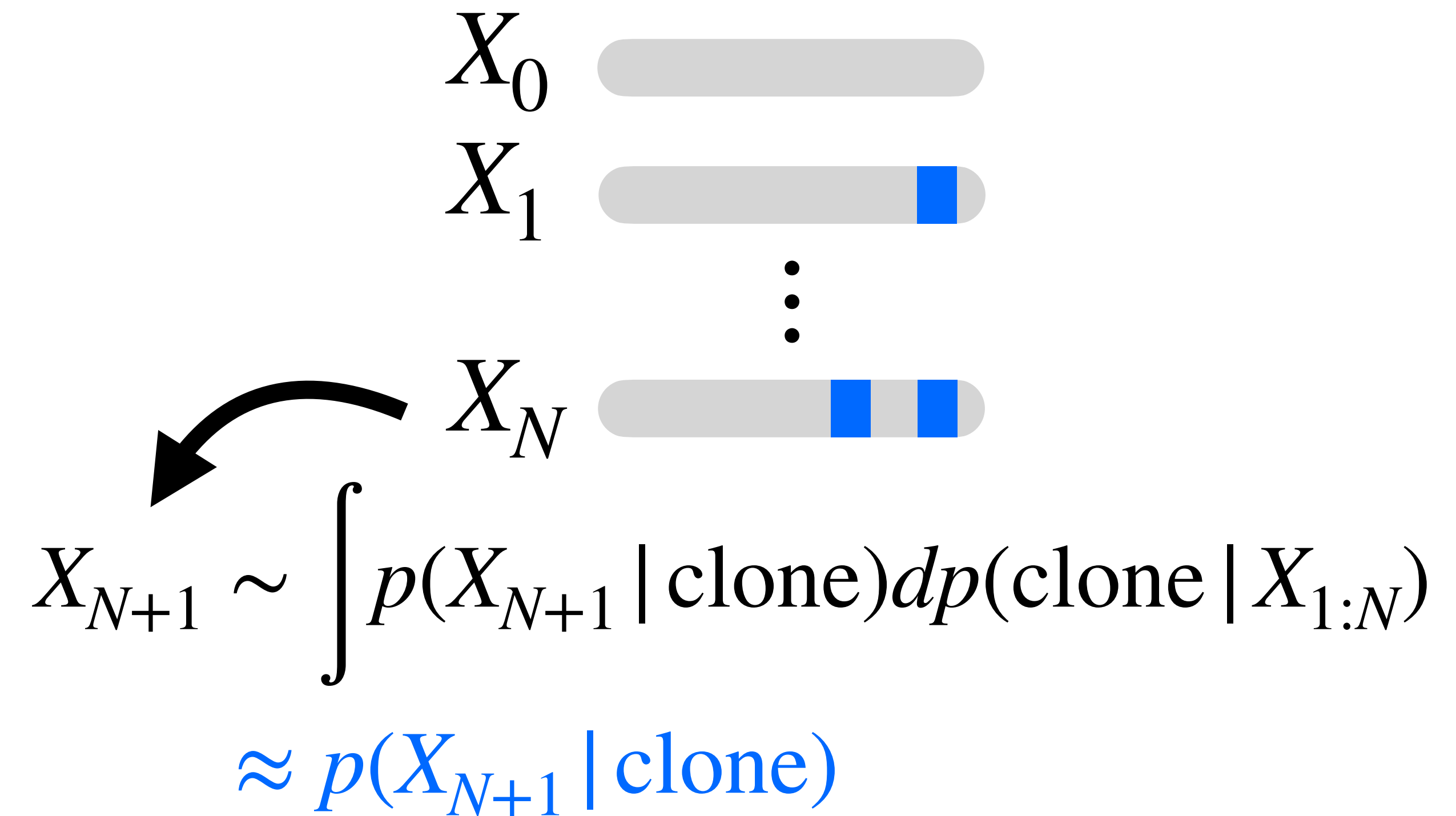
Predictions should integrate over F :



The diagram illustrates the integration of a clone distribution. It shows two horizontal bars representing evolutionary states: X_0 (a plain gray bar) and X_1 (a gray bar with a small blue segment at the right end). A curved arrow points from the blue segment of X_1 to the equation below. The equation is:

$$X_2 \sim \int p(X_2 | \text{clone}) dp(\text{clone} | X_0, X_1)$$

Predictions should converge to F :



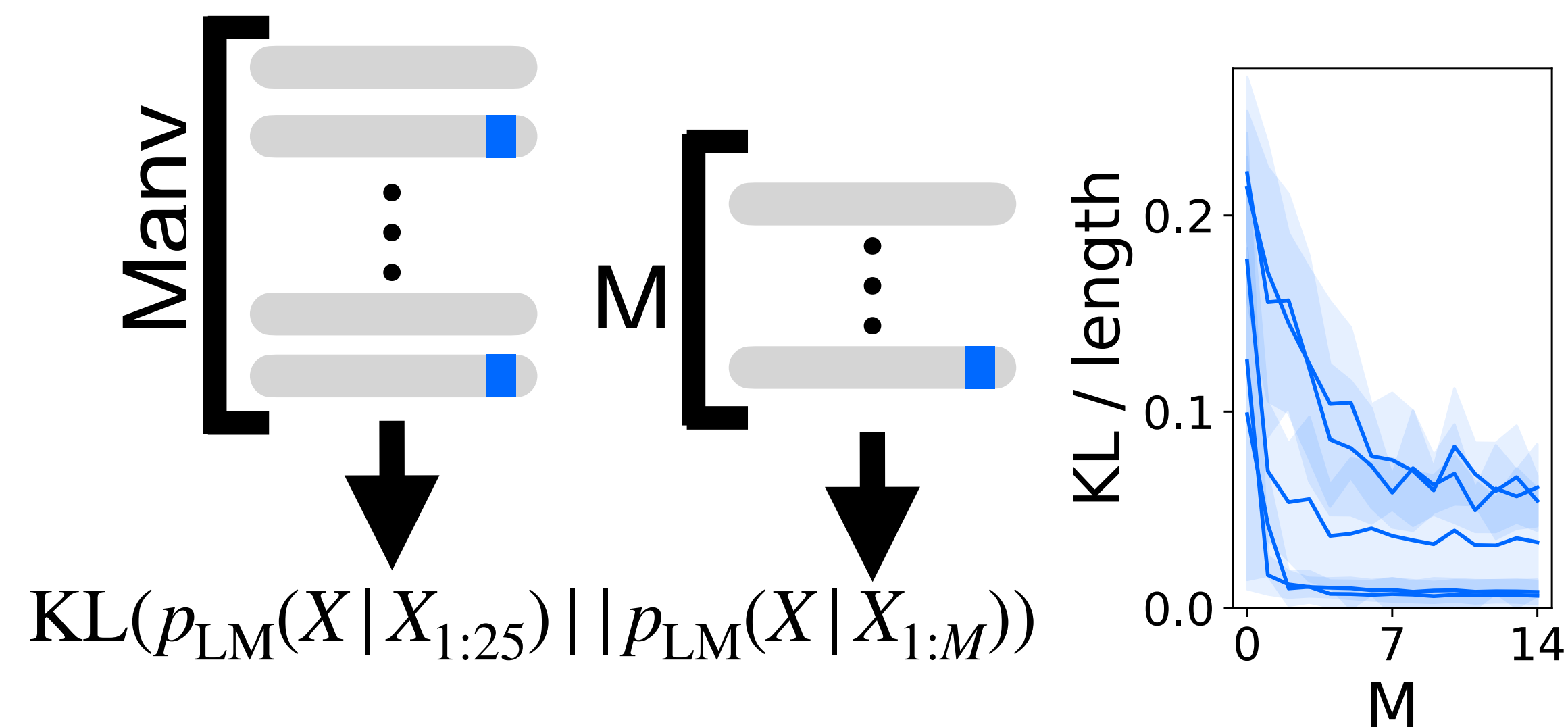
The diagram illustrates the convergence of a clone distribution. It shows a sequence of horizontal bars representing evolutionary states: X_0 (a plain gray bar), X_1 (a gray bar with a small blue segment at the right end), and X_N (a gray bar with two small blue segments at the right end). Vertical dots between X_1 and X_N indicate intermediate states. A curved arrow points from the blue segments of X_N to the equation below. The equation is:

$$X_{N+1} \sim \int p(X_{N+1} | \text{clone}) dp(\text{clone} | X_{1:N})$$

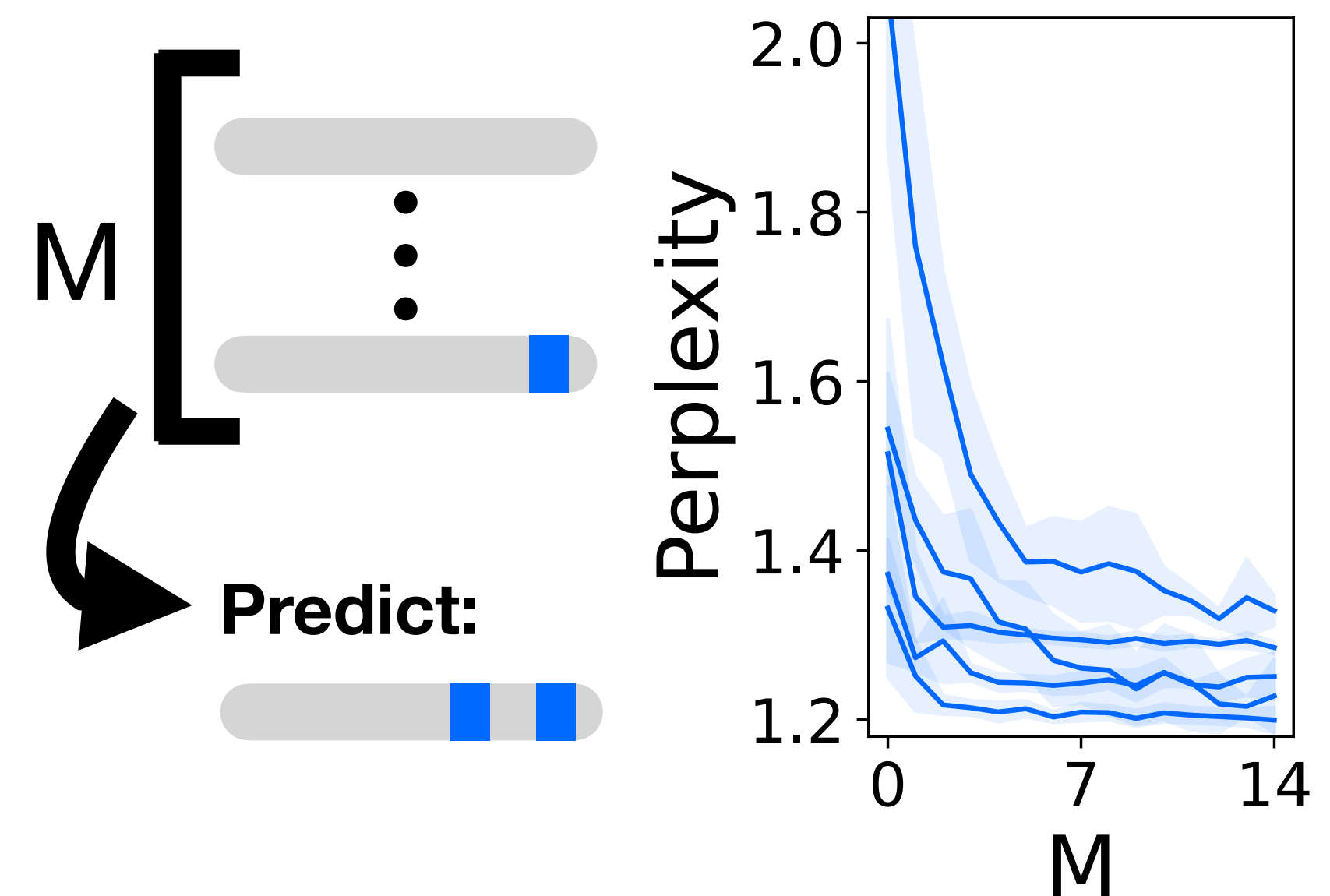
$\approx p(X_{N+1} | \text{clone})$

Given more sequences from a clonal family, CloneLM predictions converge to $p(X | \text{clone})$

Predictions converge:

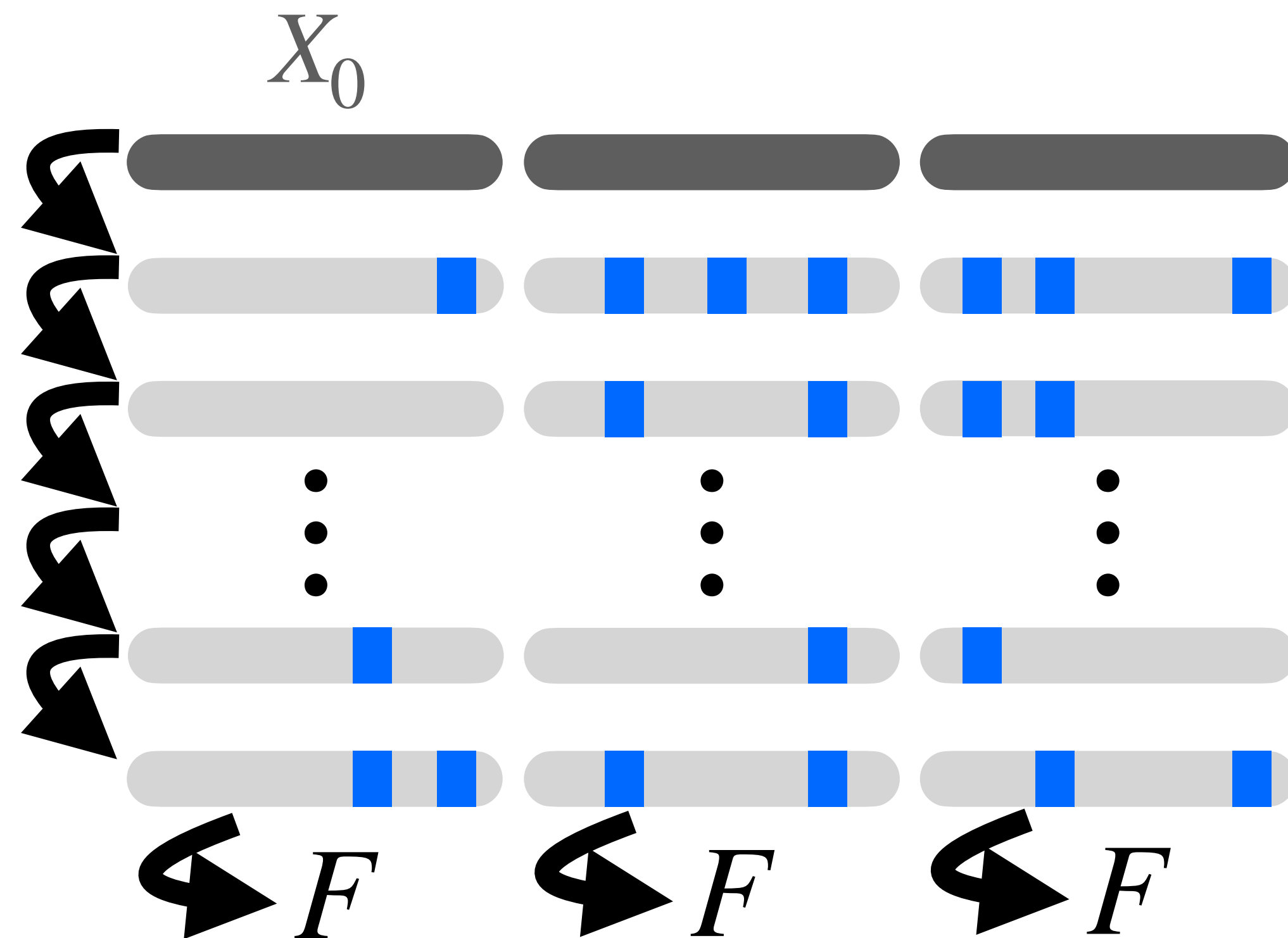


Predictions approach
 $p(X | \text{clone})$:

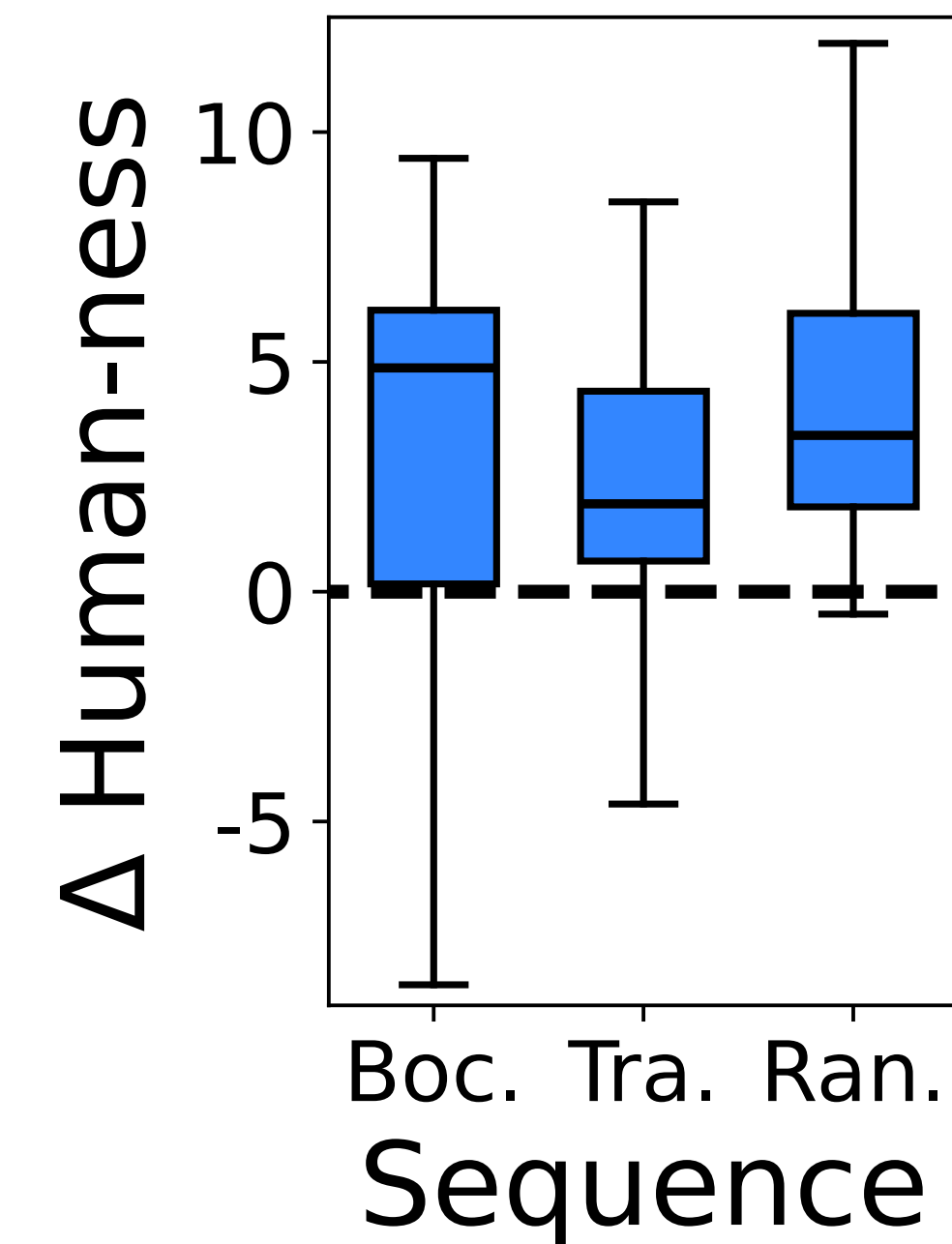


The CloneLM prior over fitness functions optimizes antibodies to become more human-like

Sample possible fitness landscapes F :

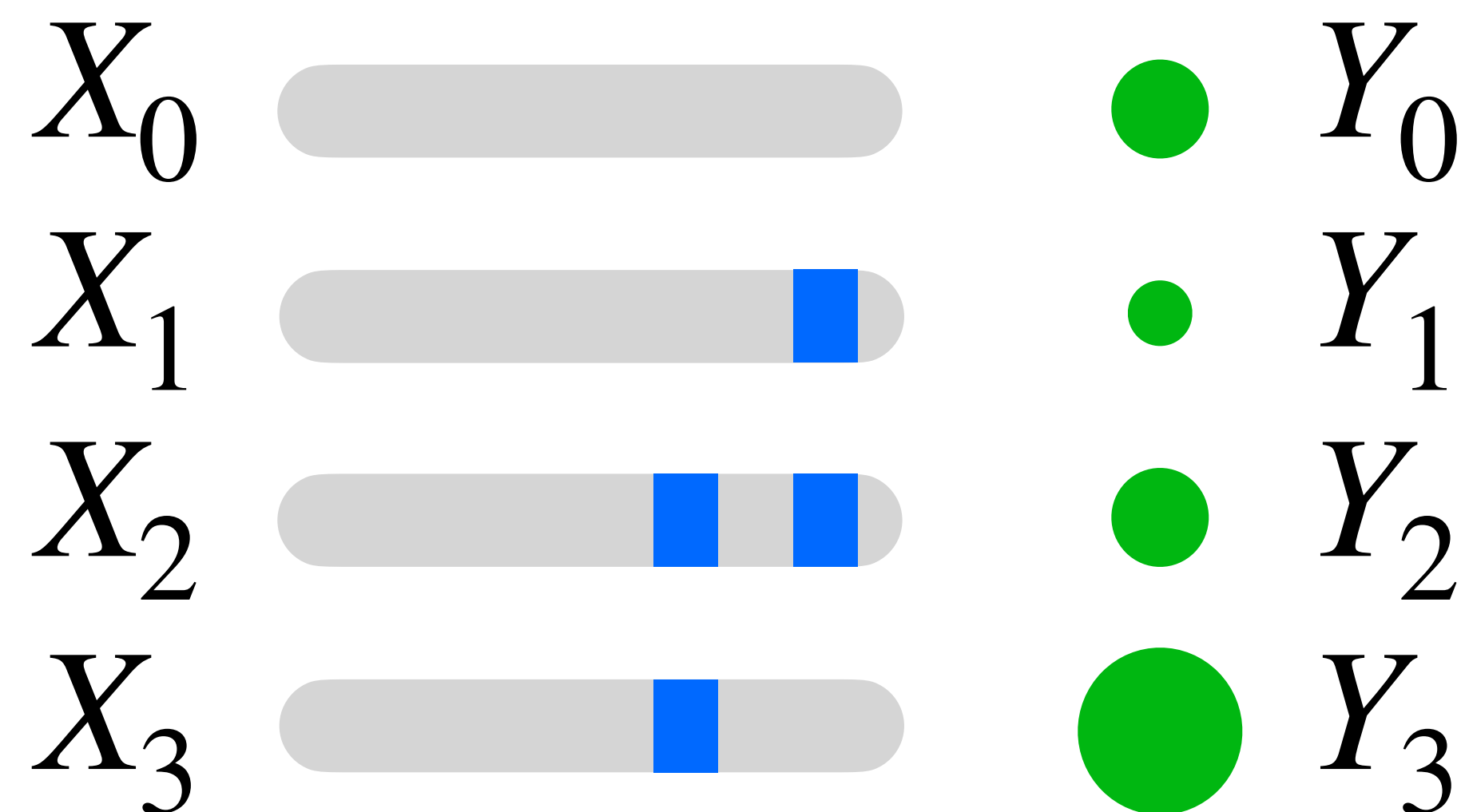


Optimize over three mutations with respect to F :



To update our belief in F , we assume measurements in the lab are proportional to fitness

Experiment:



Posterior:

$$p(F | X_0, (X_n, Y_n)_n) \propto p(F | X_0) \prod_n p(Y_n | F(X_n))$$

Likelihood:

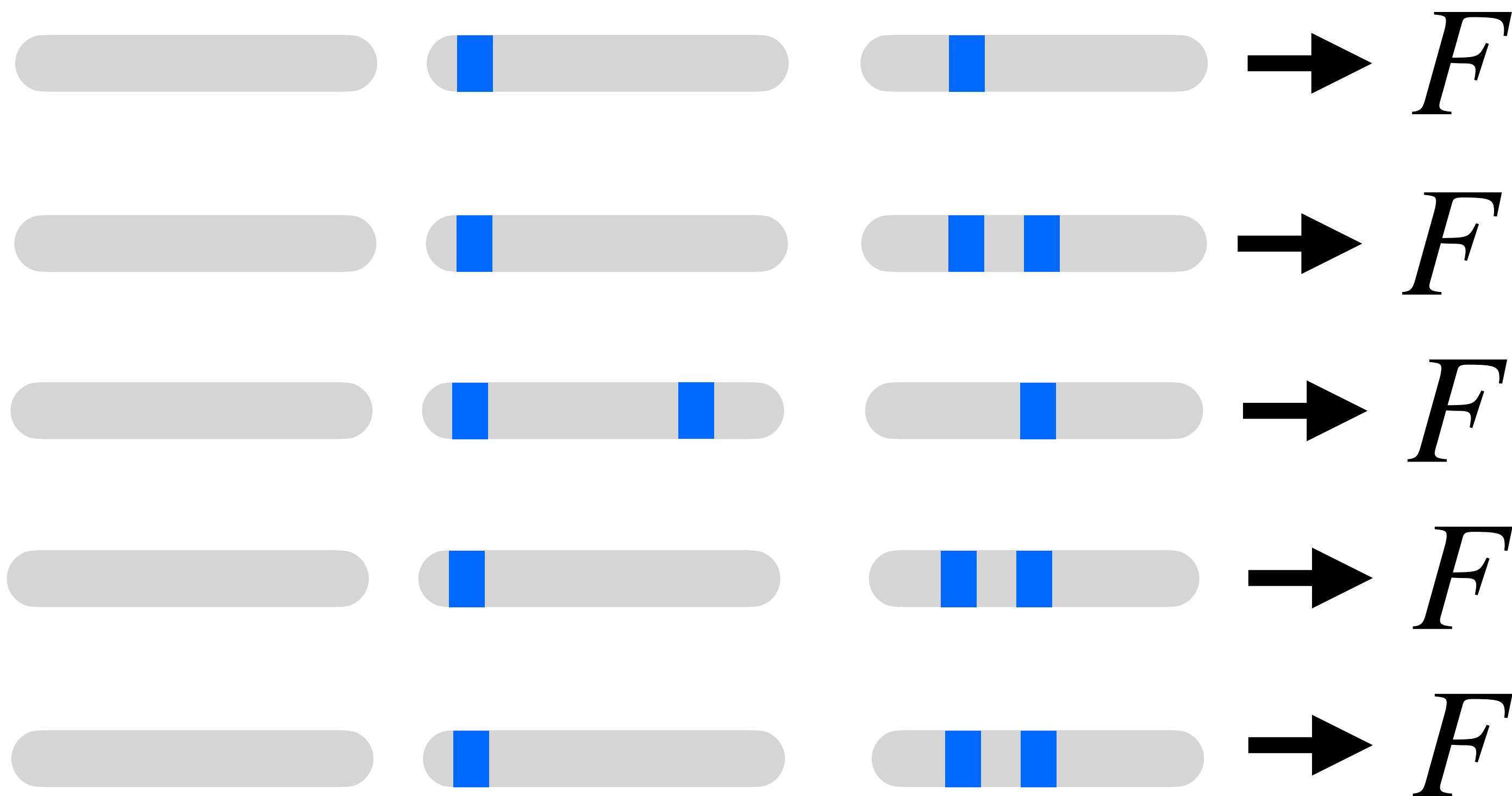
$$p(Y | F(X)) \sim \mathcal{N}(\beta F(X) + C, \sigma^2)$$

Uniform prior on β, C

How do we sample from the posterior?

Importance sample (naive):

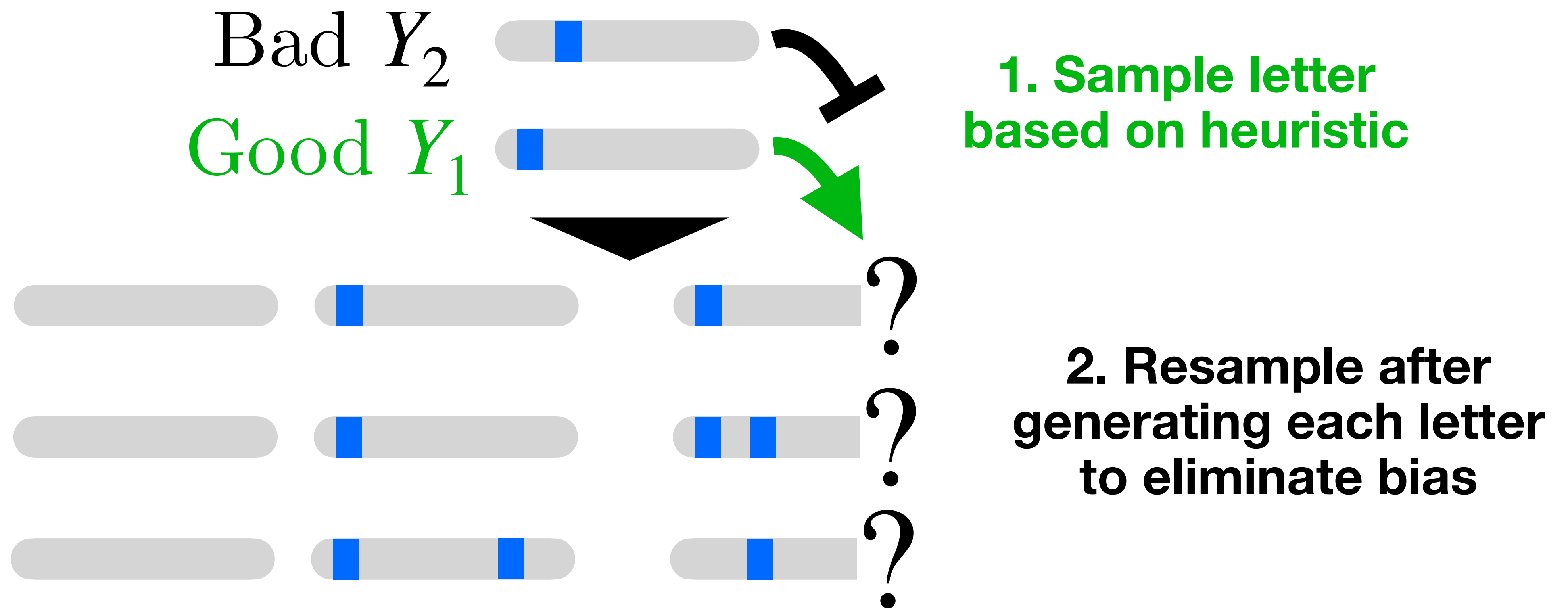
1. Sample many F from prior $F \sim p(F | X_0)$



2. Resample based on likelihood

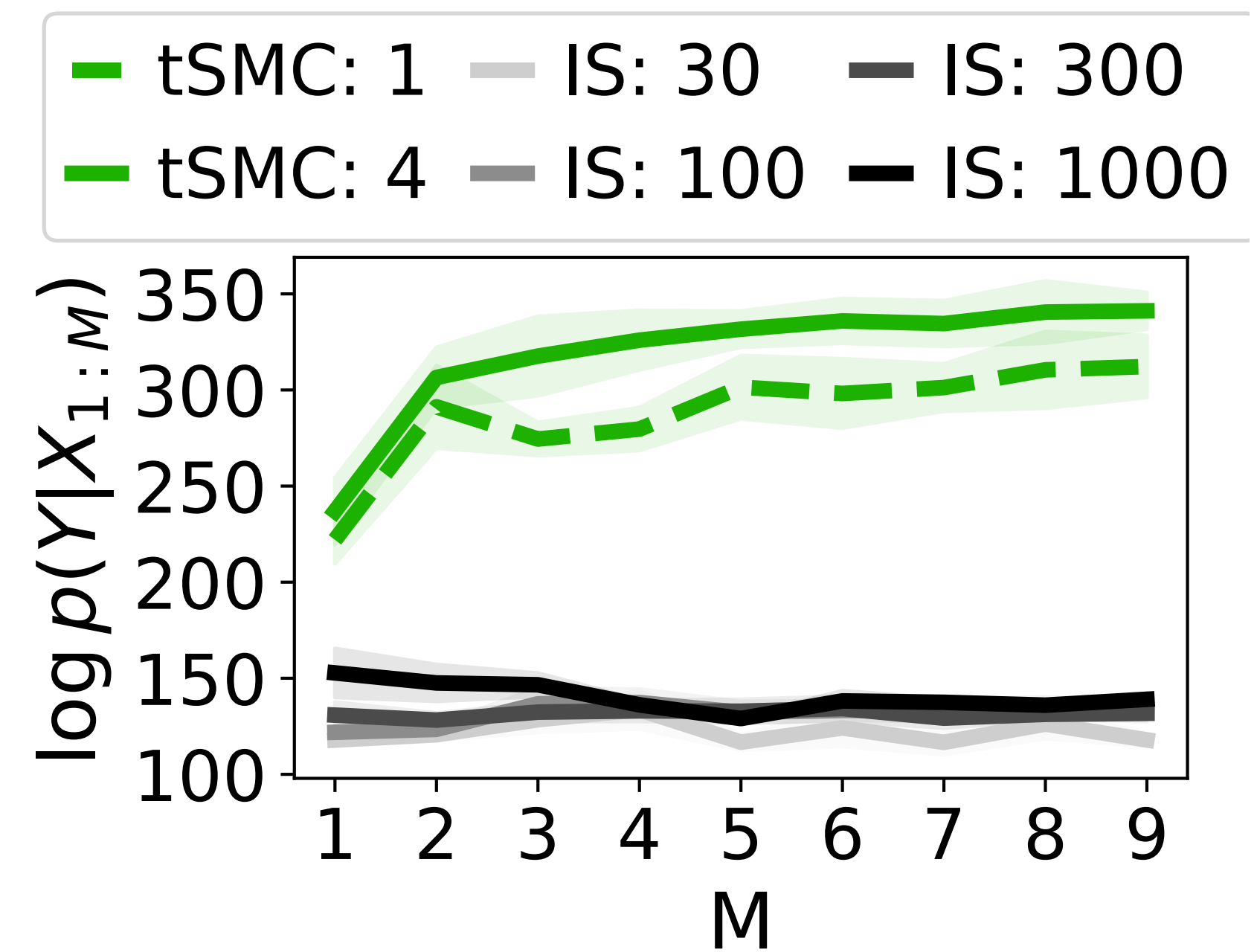
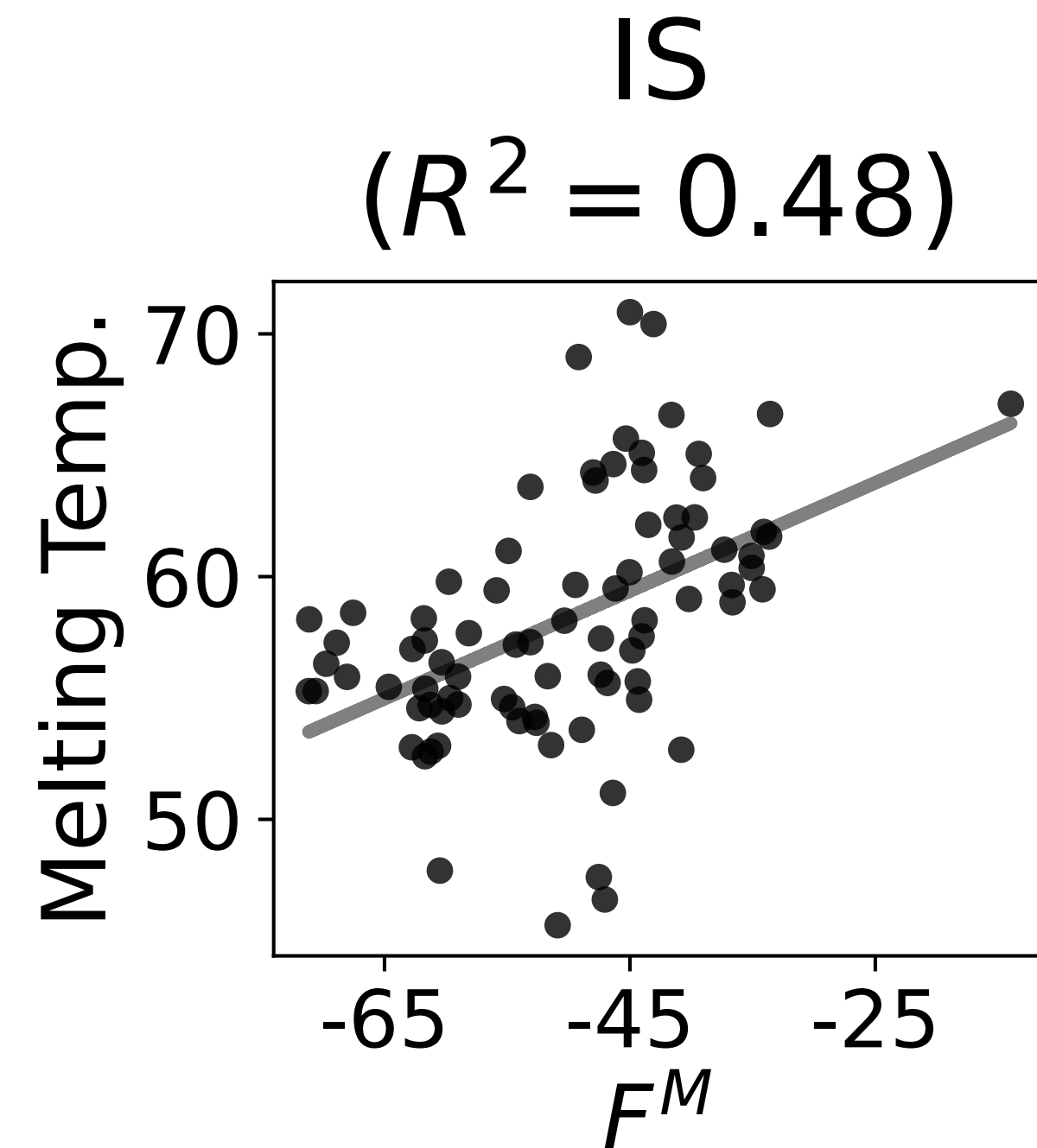
$$\prod_n p(Y_n | F(X_n))$$

We inform our sampling with twisted stochastic Monte Carlo

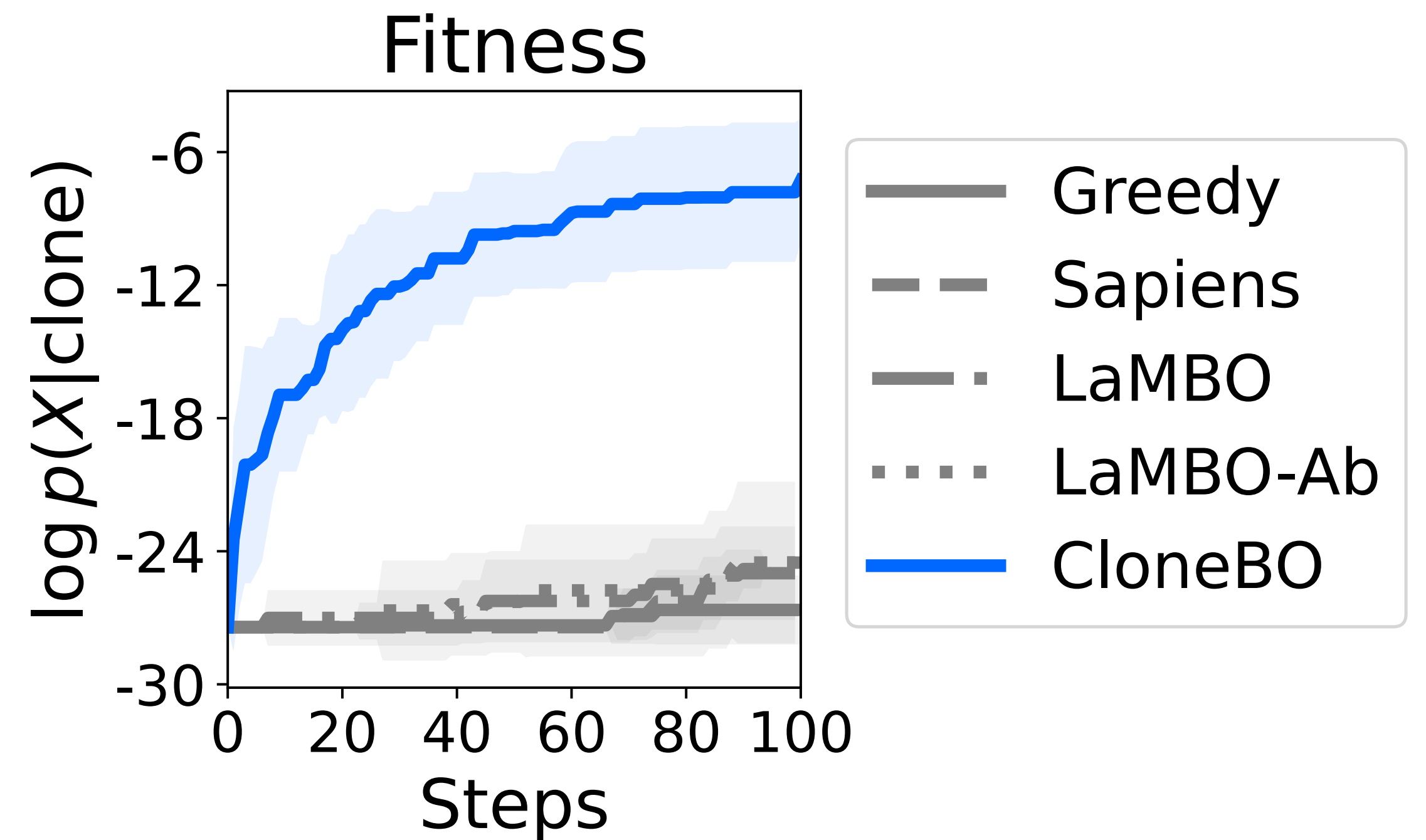
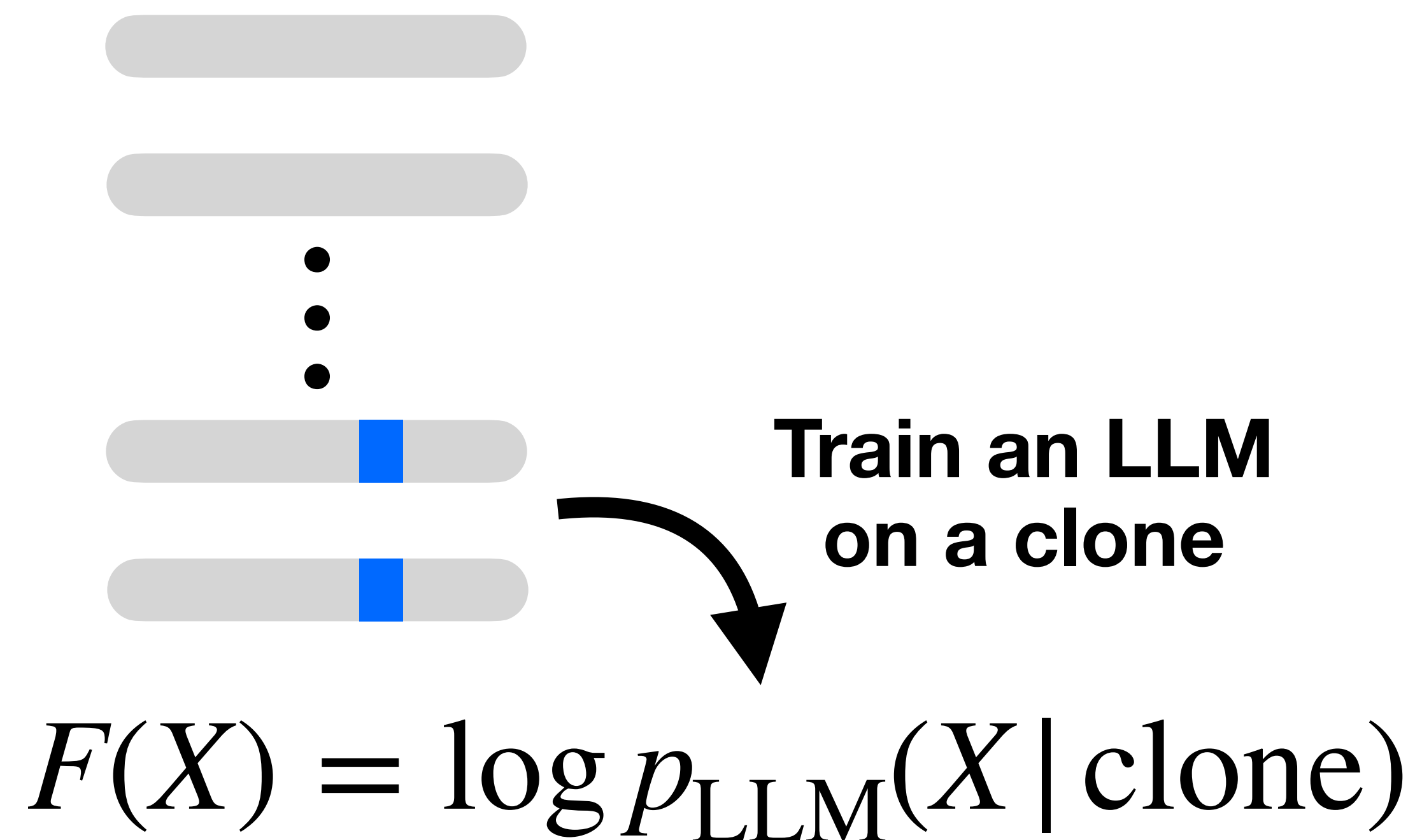


tSMC efficiently samples from the posterior

Condition on real T_m measurements

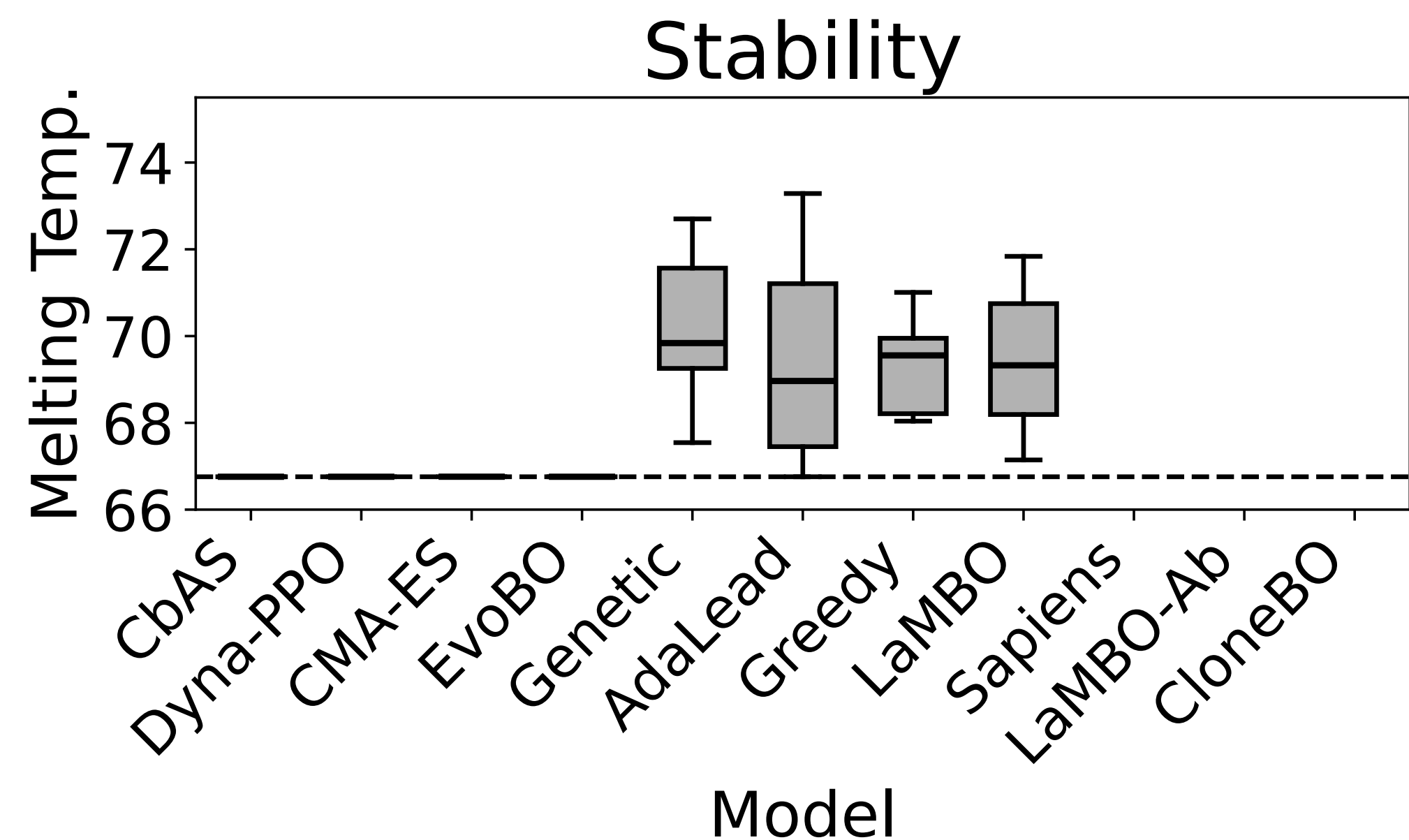
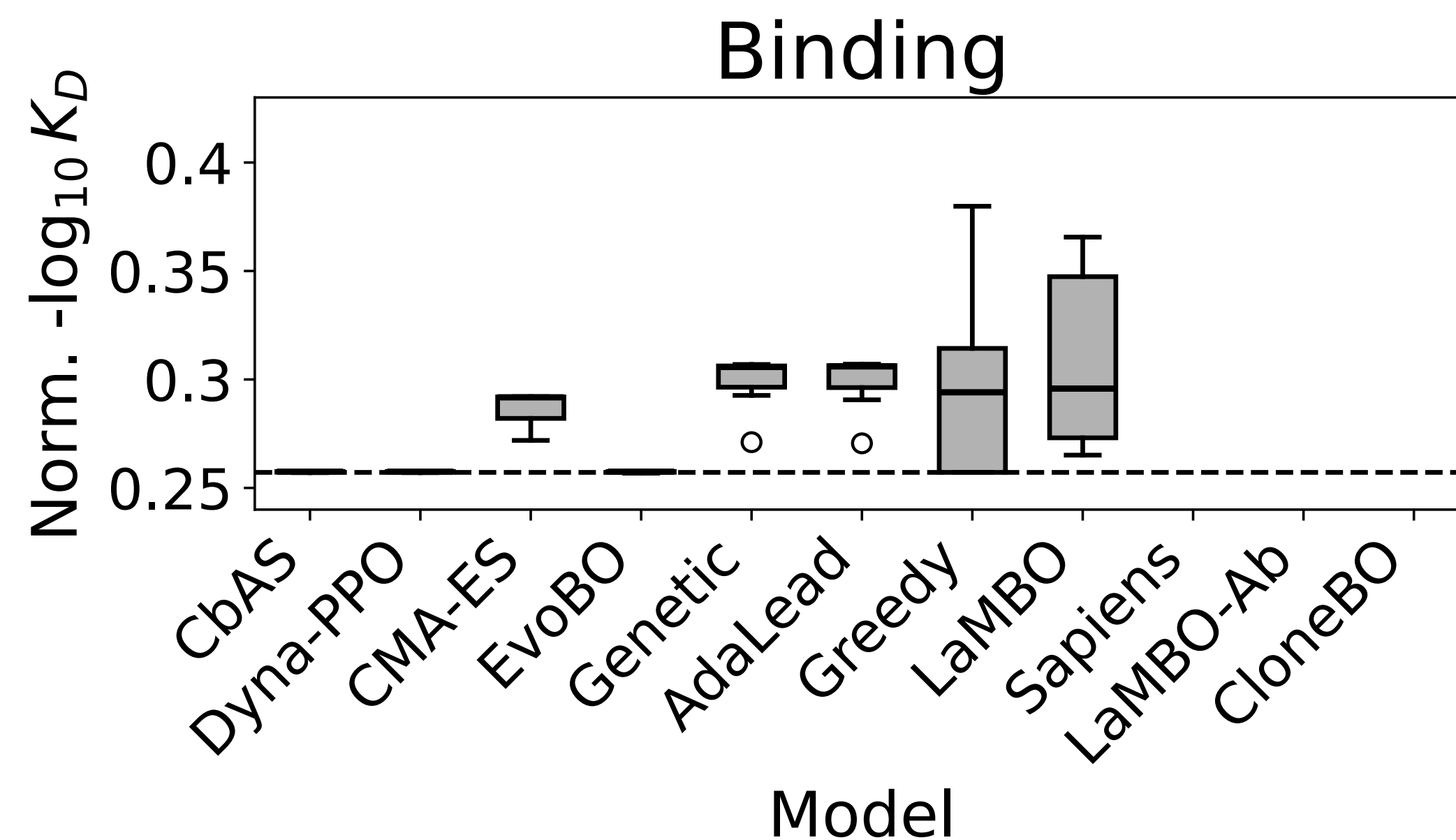


CloneBO efficiently optimizes a function from its prior

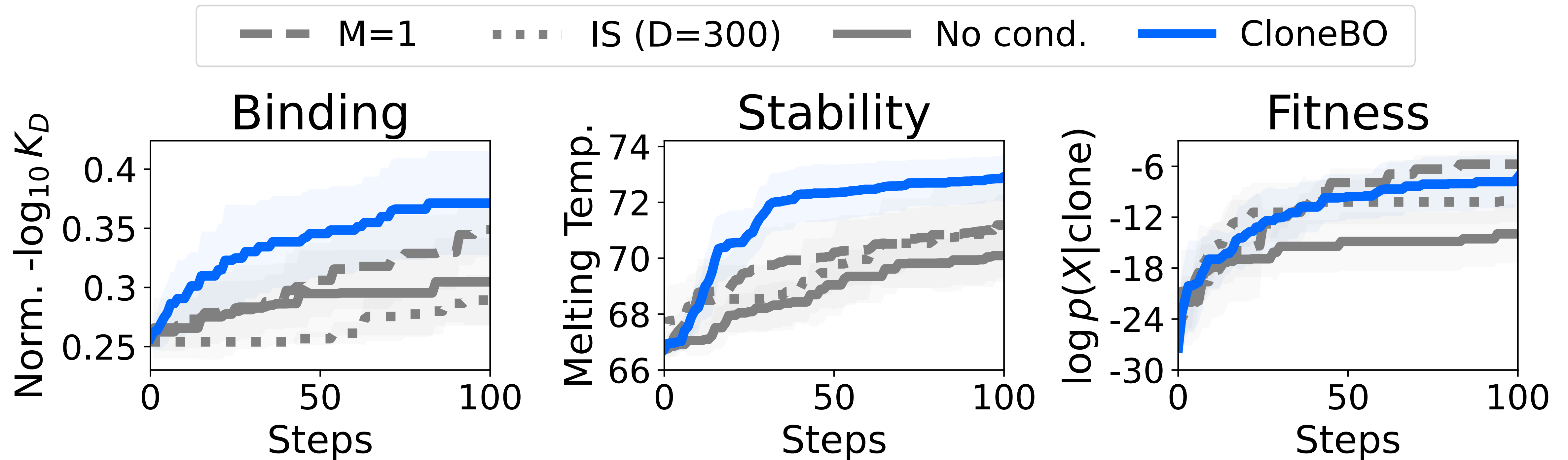


CloneBO efficiently optimizes for binding and stability *in silico*

$F(X)$ = neural network trained on thousands of sequences from iterative design experiment



Ablations demonstrate that CloneBO efficiently optimizes sequences by doing accurate inference

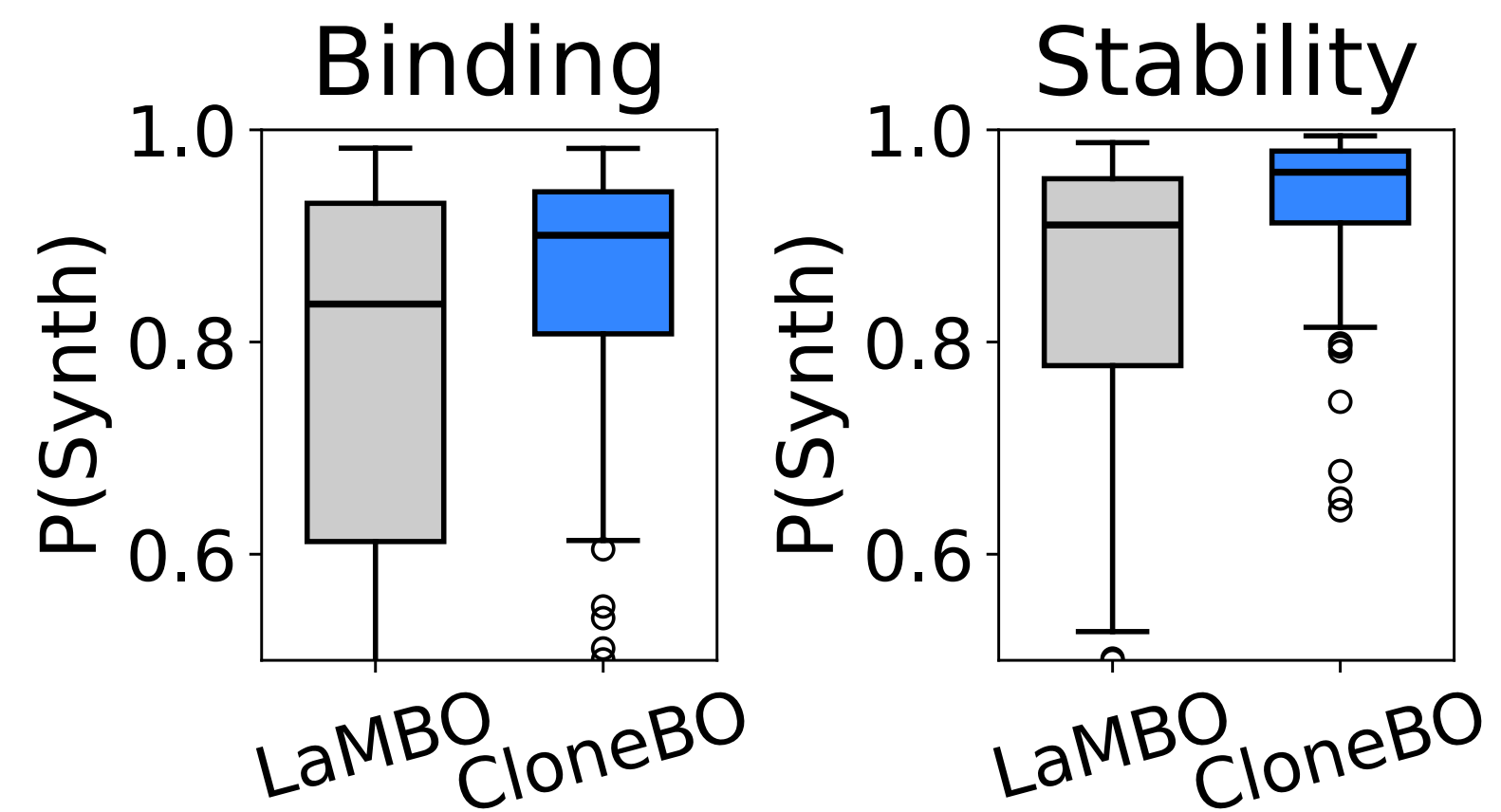


(Posterior is not hard to approximate in this case)

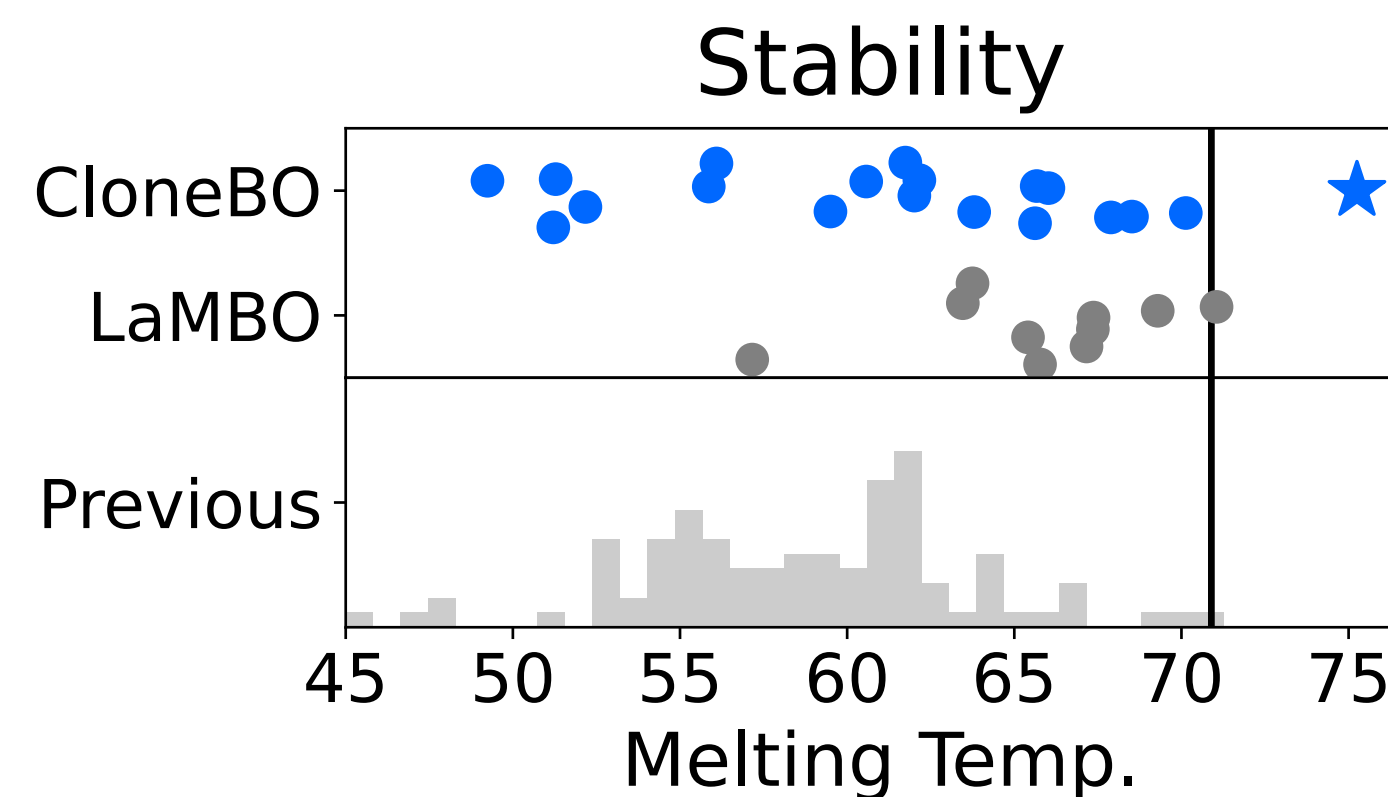
CloneBO efficiently optimizes for binding and stability *in vitro*

Given X_0, \dots, X_{1000} , measurements for binding and stability, design X_{1001}

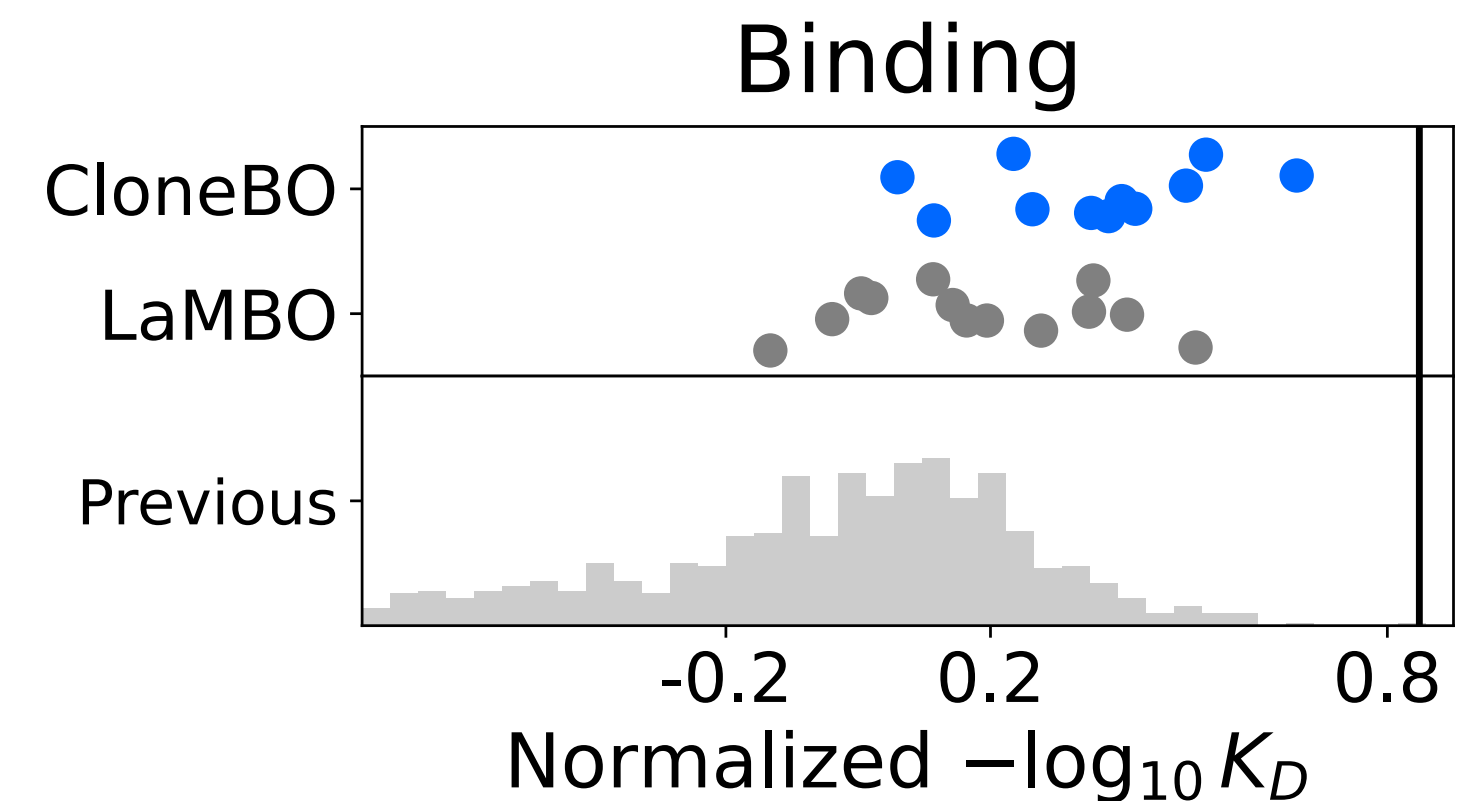
Designs are predicted to express:



Designs improve stability:



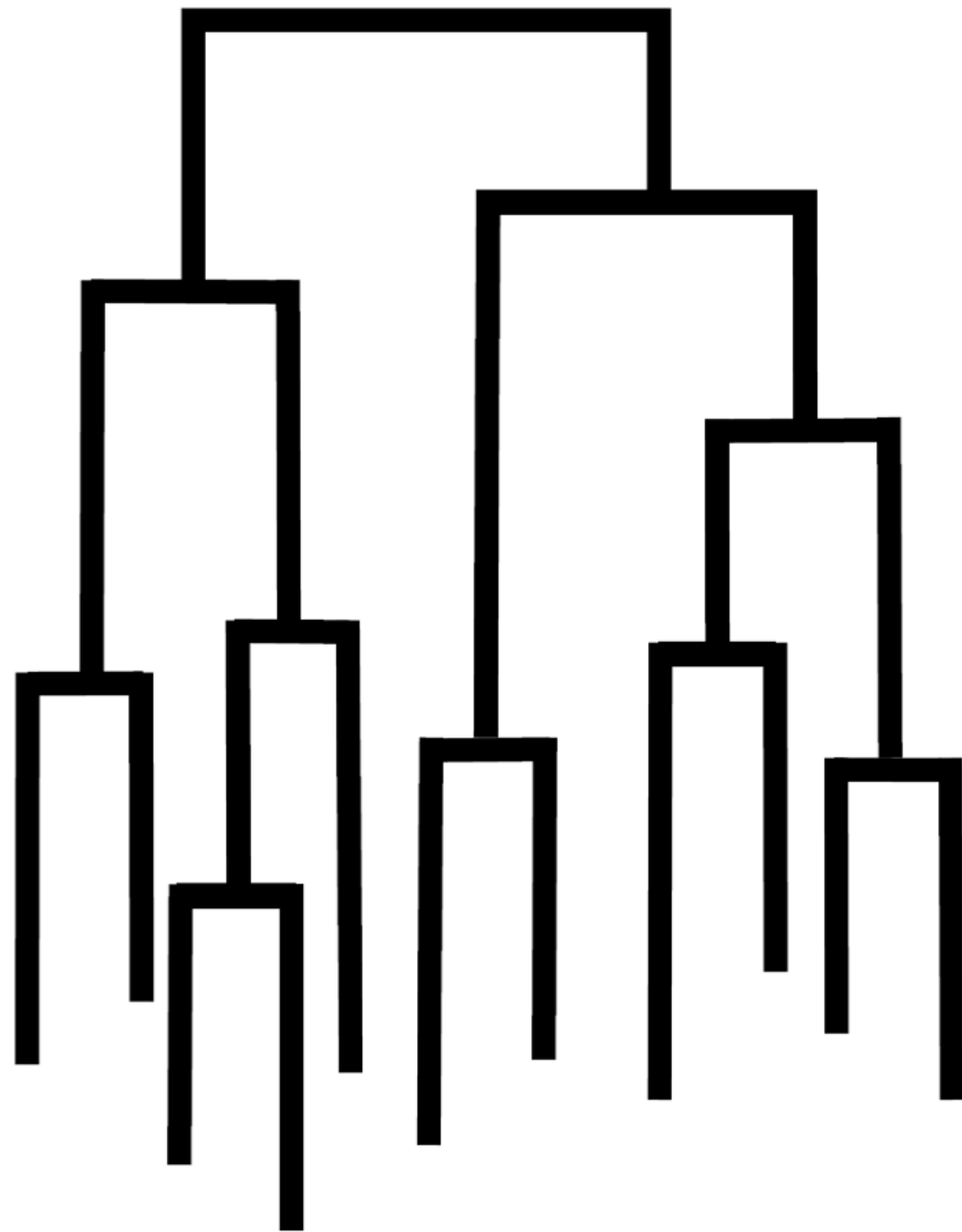
Designs outperform in binding:



Future directions

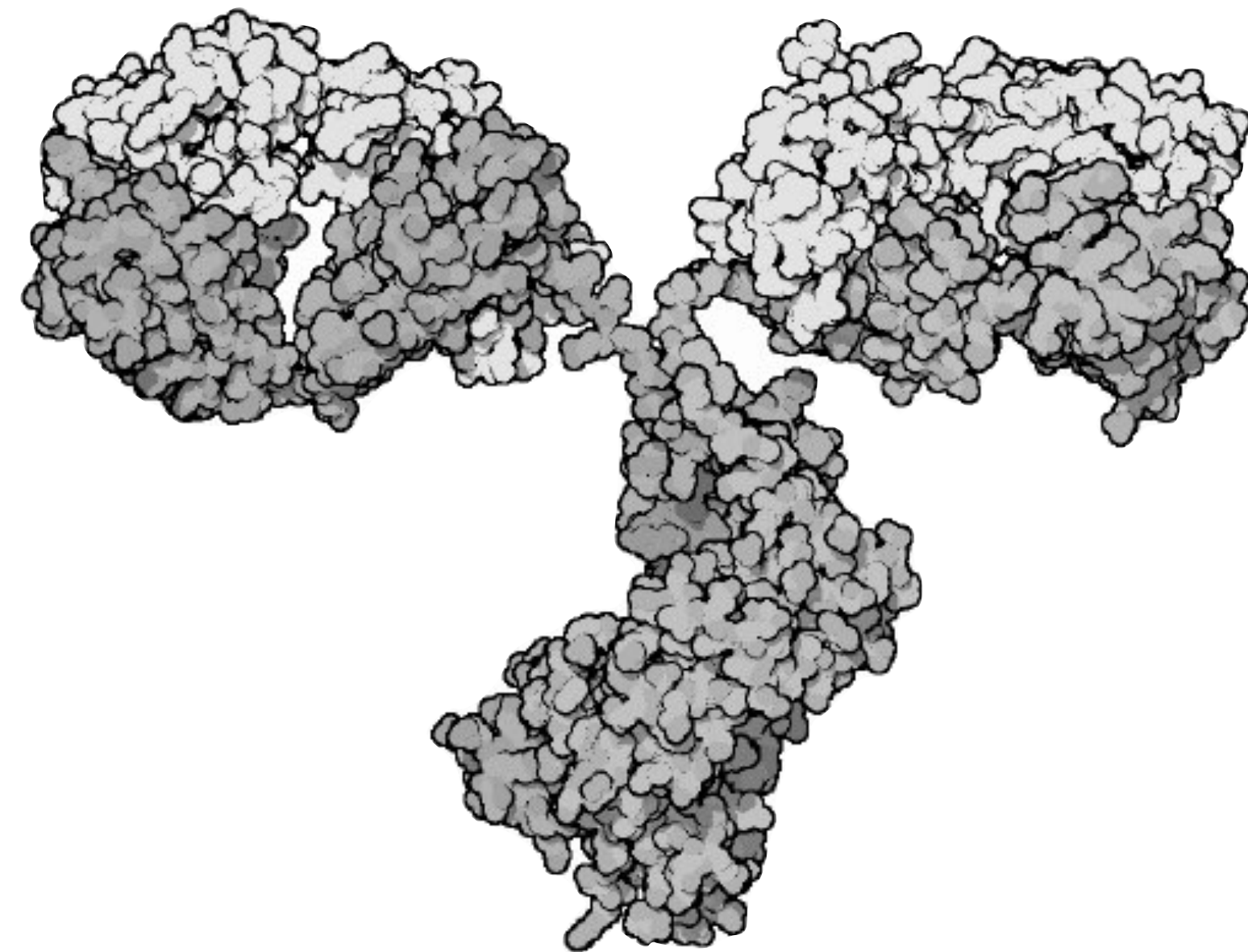
More realistic prior

Learn direction of evolution



Build a prior on structure

Structure from clonal families



Acknowledgments



NEW YORK UNIVERSITY



BigHat
BIOSCIENCES

Nate Gruver



Lily Li



Calvin McCarter



Aniruddh Raghu



Yilun Kuang



Andrew G Wilson



Hunter Elliott



Peyton Greenside

