# Extrapolative benchmarking of model-based discrete sampling methods for RNA design

*Joseph D. Valencia[1], David A. Hendrix[1,2]*

Oregon State University

1. School of Electrical Engineering and Computer Science 2. Department of Biochemistry and Biophysics

**Code + References**

## Objectives

- Survey recent approaches to model-based optimization and sampling on discrete data

- Apply sampling algorithms to RNA design

- Evaluate best practices for generalization and extrapolation of designed sequences
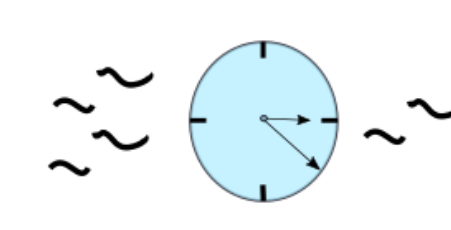
## Background

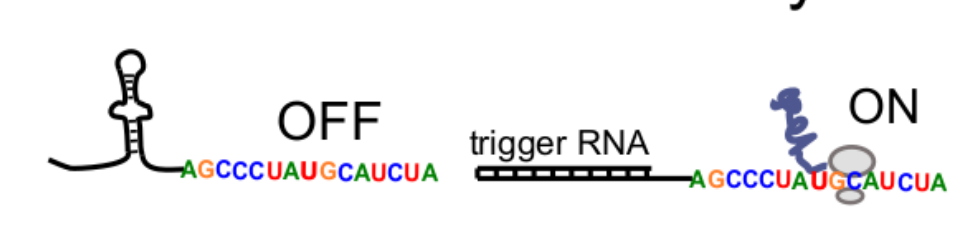Evaluate forward model

AGCCCUAUGCAUCUA → **f(x)**

ribosome load
AGCCCUAUGCAUCUA
(Sample et. al 2019)

half-life
(Agarwal and Kelley 2022)

toehold switch activity
OFF — trigger RNA — ON
(Valeri et. al 2020)

A motivating use case -- mRNA vaccines

Require high/controllable protein expression, long half-life

Given sequence → property models, how to efficiently explore a combinatorial space?

Possible guiding principle – incorporate gradients of discrete samples

## Methods

(Bogard et al 2019, Linder and Seelig 2021)

### Probabilistic Reparameterization

AGACCUAUGCCUCUA

$x \sim p(x|\theta) = \text{Categorical}(\text{Softmax}(\theta))$

How to update the parameters?

**(REINFORCE estimator)**

$\nabla_\theta \mathbb{E}_{x \sim p(x|\theta)}[f(x)] = \mathbb{E}_{x \sim p(x|\theta)}[f(x) \nabla_\theta \log p(x|\theta)]$

**(Straight-through estimator)**

$\nabla_\theta f(x) \approx \nabla_x f(x)$

### Markov Chain Monte Carlo

AGCCCUAUGCAUCUA
AGCGCAAUGCAUCUA
AGCACCAUGCAUCUA

**Discrete Langevin Proposal** (Zhang et. al 2022)

$x \sim \text{Categorical}(\text{Softmax}(\frac{1}{2}D(x,x') - \frac{\|x'_i - x_i\|_2^2}{2\alpha}))$

**Taylor approx. likelihood ratio**

$f(x') - f(x) \approx \nabla_x f(x)^\top (x'_i - x_i) = D(x,x')$

**Gibbs-with-Gradients** (Grathwohl et. al 2021)

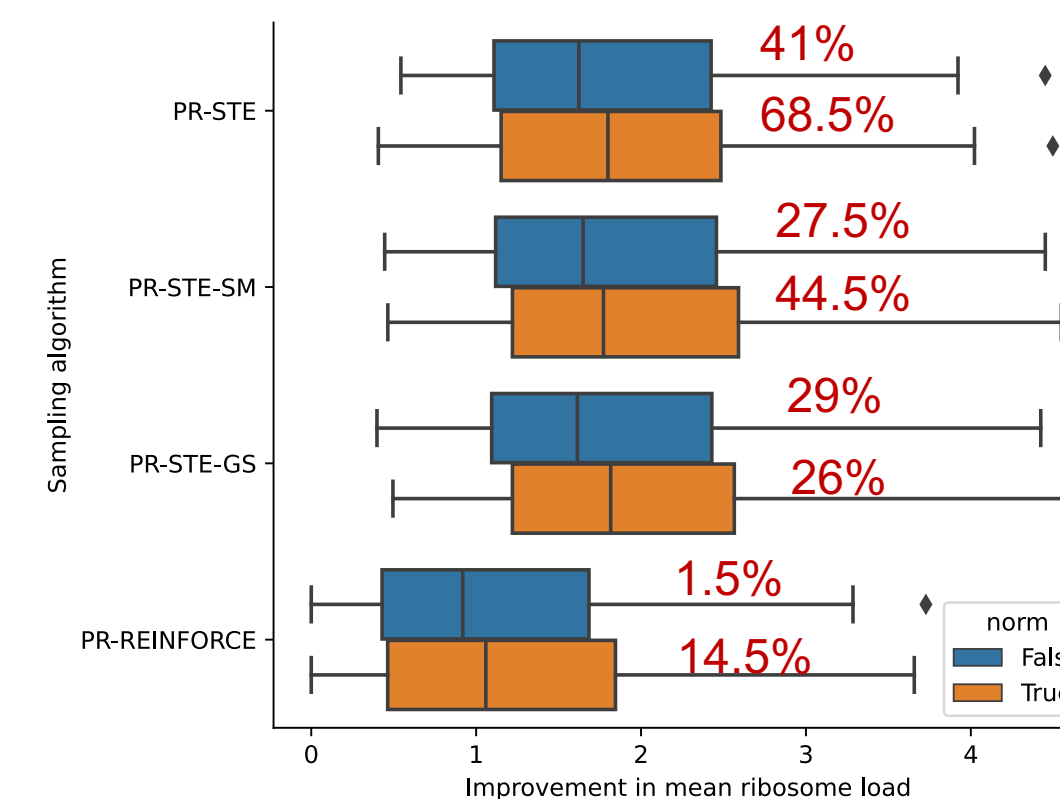$x' - x \sim \text{Categorical}(\text{Softmax}(vec(D(x,x'))))$

**Algorithm 1** Extrapolative Benchmark

**Require:** $Z = \{(x_1, y_1), (x_2, y_x)) \ldots (x_n, y_n)\}$ ▷ Training data (*seq, property*)
**Require:** $q$ ▷ Exclusion quantile
**Require:** $S$ ▷ Seed sequences to redesign
$Z^- = Z \setminus \{(x,y) \in Z | Q(y) > q\}$ ▷ Excludes highest quantile by property
$f_O \leftarrow max_\theta f_\theta(y|x), \quad (x,y) \in Z$ ▷ Train oracle model on the full train set
$f_D \leftarrow max_\theta f_\theta(y|x), \quad (x,y) \in Z^-$ ▷ Train designer model on the reduced train set
$D \leftarrow \{\}$
**for** $s \in S$ **do**
    $s' \leftarrow \text{DiscreteDesign}(f_D, s)$ ▷ MCMC or PR seeded by $s$
    $y' \leftarrow f_O(s')$ ▷ Impute property with oracle
    $D.add((s', y'))$
**end for**
$summary \leftarrow |\{(x,y) \in D | Q(y) > q\}|$ ▷ # designs exceeding train set maximum

*Conditional sampling     Property model     **Sequence Prior***

$$\nabla_x \log P(x | y) = \nabla_x \log P(y | x) + \nabla_x \log P(x)$$

As generative models of RNA become more available, incorporate a prior to mitigate pathological sampling

## Future Directions

- Hybrid straight-through/REINFORCE estimators
- Evaluate sample diversity
- Adaptive preconditioning for MCMC
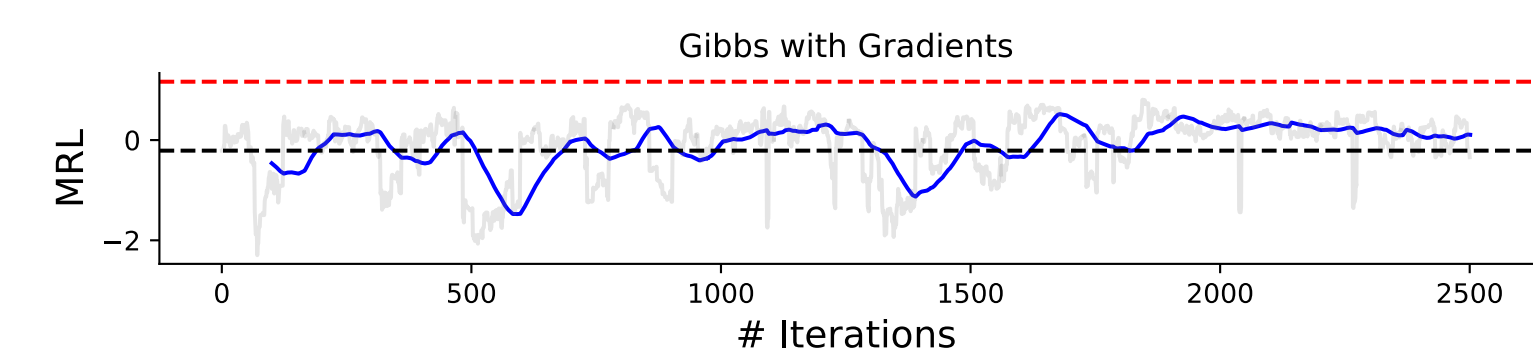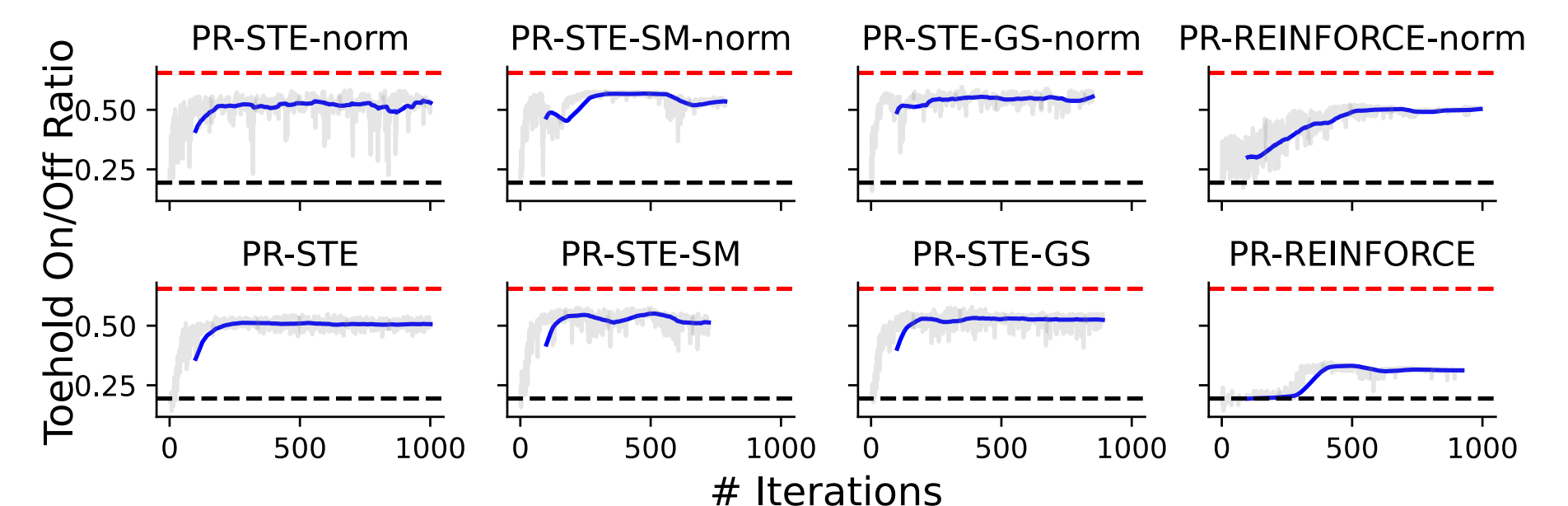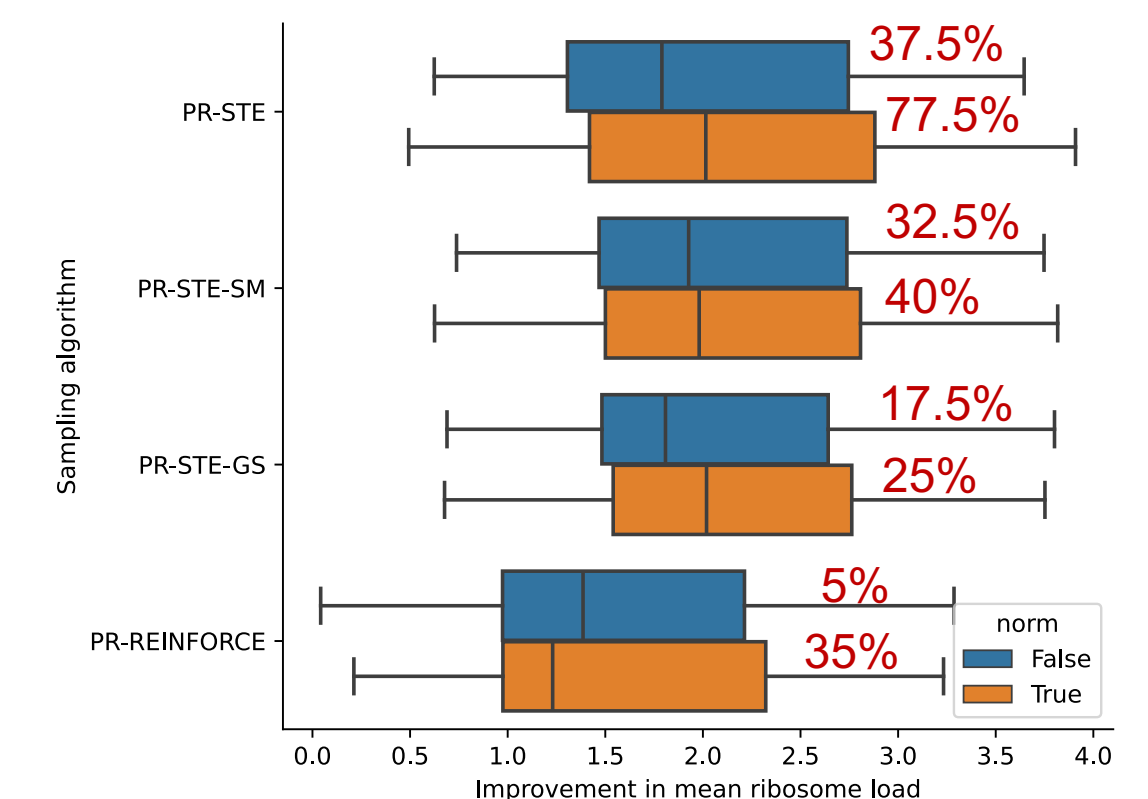- Gradient smoothness regularization (Miyato 2016)

## Preliminary Results

% = percentage exceeding train set max property value according to oracle

Iterations=$10^3$, samples/step = 32, trials=200



Iterations=$10^4$, samples/step = 4, trials=40





Gibbs with Gradients