# Biological Sequence Design

intractable combinatorial optimization problem

slow/impossible to compute
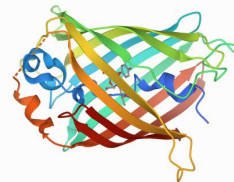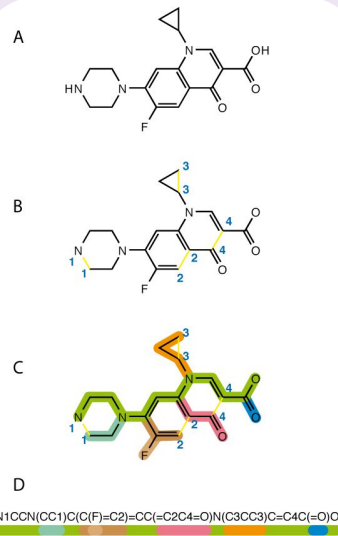
$$\max_{\mathbf{x} \in \mathcal{X}}(f_1(\mathbf{x}), \ldots, f_k(\mathbf{x}))$$

- Discrete, high-dim. inputs
- Multiple black-box objectives
- Batched experiments
- Noisy labels



A

B

C

D

N1CCN(CC1)C(C(F)=C2)=CC(=C2C4=O)N(C3CC3)C=C4C(=O)O



```
         10         20         30         40         50
MVSKGEELFT GVVPILVELD GDVNGHKFSV SGEGEGDATY GKLTLKFICT
         60         70         80         90        100
TGKLPVPWPT LVTTLTYGVQ CFSRYPDHMK QHDFFKSAMP EGYVQERTIF
        110        120        130        140        150
FKDDGNYKTR AEVKFEGDTL VNRIELKGID FKEDGNILGH KLEYNYNSHN
        160        170        180        190        200
VYIMADKQKN GIKVNFKIRH NIEDGSVQLA DHYQQNTPIG DGPVLLPDNH
        210        220        230
YLSTQSALSK DPNEKRDHMV LLEFVTAAGI TLGMDELYK
```
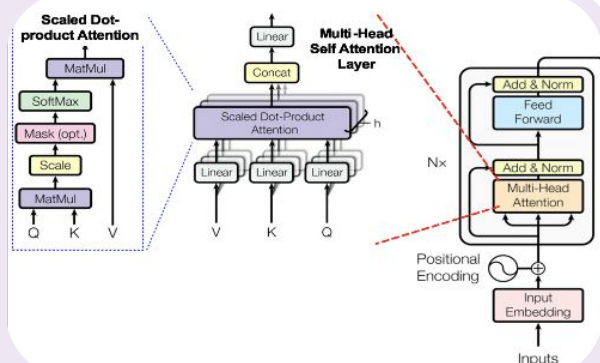
# Deep Generative Models (DGMs)

optimize over small
DGM-generated subset

ad-hoc proxy objective

$$\max_{\mathbf{x} \in \mathcal{X}'} \text{ ???}$$



- How do we rank generated sequences, accounting for multiple objectives, explore-exploit, etc.?

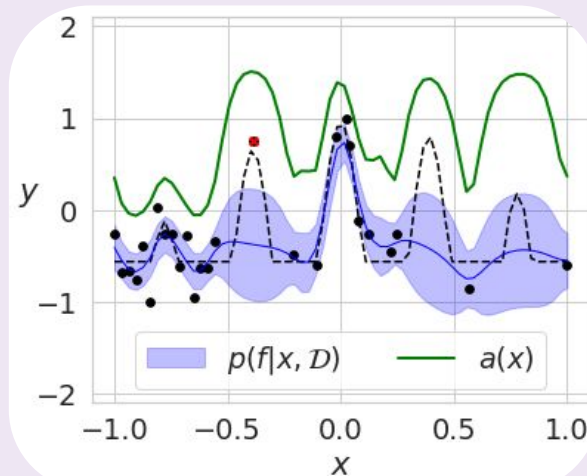- How do we generate good subsets for ranking?

# Bayesian Optimization (BayesOpt)

intractable combinatorial optimization subproblem

principled proxy objective

$$\max_{\mathbf{x} \in \mathcal{X}} \mathbb{E}[u(\mathbf{x})]$$

- Strategy 1: solve subproblem with genetic algorithms (inefficient).

- Strategy 2: use a frozen pretrained generative model, solve in latent space (data-hungry).
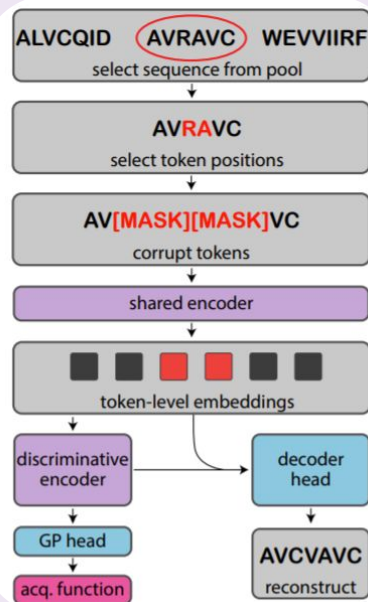
# Latent Multi-Objective BayesOpt (LaMBO)

**LaMBO** is designed from the ground up to combine the best attributes of DGMs and BayesOpt.
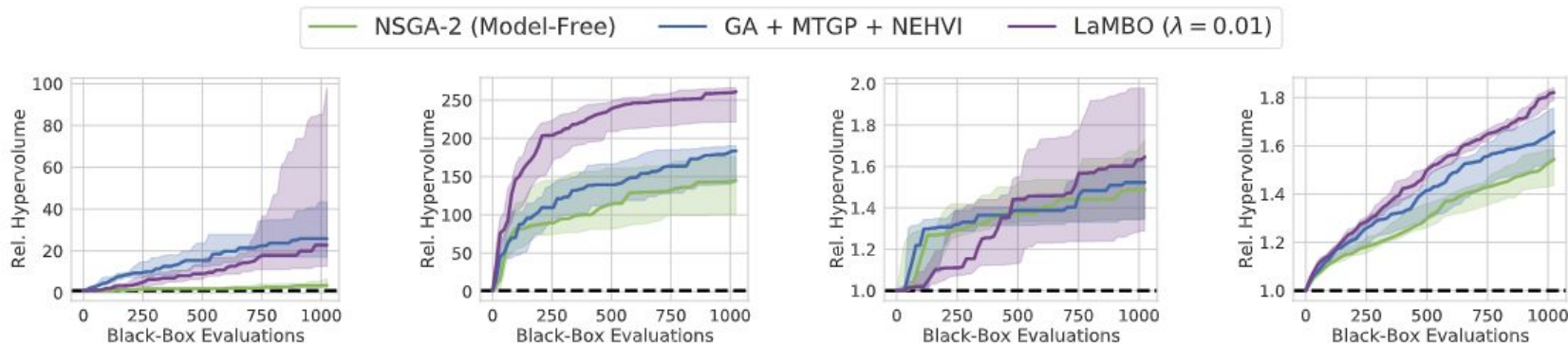
# The **LaMBO** architecture

- Jointly train a denoising autoencoder and a discriminative deep kernel GP.

- Rank samples with the NEHVI acquisition function, optimize in latent space.

- Pretraining is optional!

# Previewing the results

Comparing **LaMBO** to model-free and model-based genetic algorithms



Legend: NSGA-2 (Model-Free) — GA + MTGP + NEHVI — LaMBO ($\lambda = 0.01$)

(a) Bigrams

(b) logP + QED

(c) DRD3 Docking + SA

(d) Stability + SASA

small molecule

large molecule

# Collaborators



NYU

**Samuel Stanton**

**Wesley Maddox**

**Nate Gruver**

**Andrew Wilson**

BigHat BIOSCIENCES

**Phil Maffetone**

**Emily Delaney**

**Peyton Greenside**