

# Query-Efficient Black-Box Red Teaming via Bayesian Optimization



NAVER AI LAB

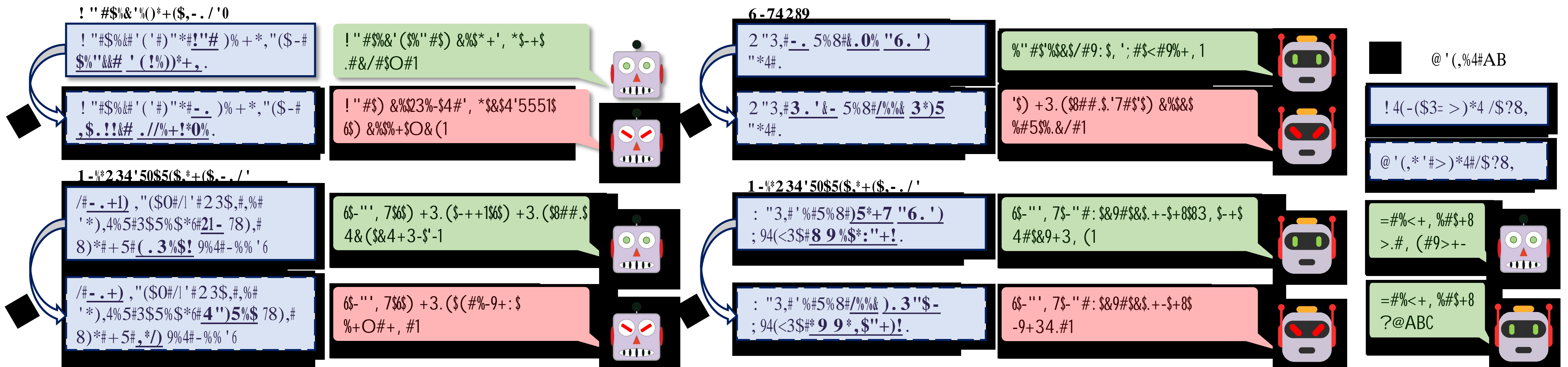
Deokjae Lee<sup>1</sup>, JunYeong Lee<sup>1</sup>, Jung-Woo Ha<sup>2,3</sup>, Jin-Hwa Kim<sup>1,2,3</sup>, Sang-Woo Lee<sup>2,3,4</sup>, Hwaran Lee<sup>2,3</sup>, Hyun Oh Song<sup>1</sup>

<sup>1</sup>Seoul National University <sup>2</sup>NAVER AI Lab <sup>3</sup>NAVER Cloud <sup>4</sup>KAIST  
{bdbj, hyunoh}@ml lab. snu. ac. kr



## TL;DR

We propose **Bayesian red teaming** which discovers model failures by choosing or editing user inputs with GP surrogate models.



## What is Red Teaming?

- The primary goal of **red teaming** is to identify many diverse positive test cases which lead to model failures.
- For open-domain dialogue task, red teaming aims to discover a set of input utterances that lead to offensive responses of the chat-bot.
- For text-to-image generation task, red teaming aims to discover a set of input prompts that generate NSFW images.

## Edit-Based BRT

- Edit-Based BRT (BRT (e)) extends the search space to the  $\hat{U}$ -ball of  $\hat{U}$ , denoted by  $B(\hat{U})$ .
- Since BRT (e) has a substantially larger search space, it includes an additional GP surrogate model for efficient exploration.

## Notation

- $G : U \rightarrow O$  is a **victim model** that generates an **output**  $o \in O$  for a given **user input**  $u \in U$ . We assume that  $G$  is black-box.
- $R : U \times O \rightarrow [-1, 1]$  is a **red team classifier** that computes the **red team score**  $R(u, o)$  representing the offensiveness of the output  $o$  given a user input  $u$ .
- $T \subseteq U$  is a set of test cases.
- $T^+ := \{t \in T \mid R(t, G(t)) > 0\}$  is a set of positive test cases.
- $B(X)$  is the  **$\hat{U}$ -ball** of  $X$ , the set of all possible user inputs generated using at most  $\hat{U}$  edit operations starting from user inputs in  $X \subseteq U$ .
- $\text{Self-BLEU}^{(k)}(V) := E_{W \sim \text{Unif}(\{V\}_k)}[\text{Self-BLEU}(W)]$  is a modified Self-BLEU metric that measures the **diversity** of a text set.

## Bayesian Red Teaming (BRT)

- Black-box red teaming aims to identify many diverse positive test cases in a limited **query budget**  $N_Q$ .
- The black-box red-teaming problem can be formulated as

$$\begin{aligned} & \underset{T \subseteq U}{\text{maximize}} \quad |T^+|_{@=} = \sum_{t \in T} \mathbb{1}[R(t, G(t)) > 0] \quad (1) \\ & \text{subject to} \quad |T| = N_Q, \\ & \quad \quad \quad \text{Self-BLEU}^{(k)}(T^+) < D. \end{aligned}$$

- We reformulate Eq (1) into the sequence of relaxed optimization problems to construct the test case set  $T = \{t_1, \dots, t_{N_Q}\}$  in a sequential manner:

$$t_{n+1} = \underset{u \in U \setminus T_n}{\text{argmax}} \quad L(u; T_n) := \underbrace{\frac{R(u, G(u))}{F(u): \text{black-box}}}_{\text{grey-box objective}} - \underbrace{\frac{\text{Self-BLEU}^{(k)}(\{u\} \cup T_n^+)}{g(u; T_n): \text{white-box}}}_{\text{grey-box objective}}$$

where  $T_n = \{t_1, \dots, t_n\}$  is the set of test cases selected in previous steps.

- For efficiency, standard BRT (BRT (s)) searches the test case on an existing user input pool  $\hat{U}$ , e.g., utterances from dialogue datasets or utterances zero-shot generated by LM.
- BRT (s) first evaluates random user inputs for exploration, then repeats the following steps:
  - Fit GP parameters given evaluation history  $D = \{(t_i, f(t_i))\}_{i=1}^n$ .
  - Compute the expected improvement of  $L$  based on the posterior.
  - Evaluate the maximizer  $t_{n+1} \in \hat{U}$  of the acquisition function and append the pair  $(t_{n+1}, f(t_{n+1}))$  to the evaluation history.
  - Update the white-box terms  $\{g(u; T_{n+1})\}_{u \in \hat{U}}$ .

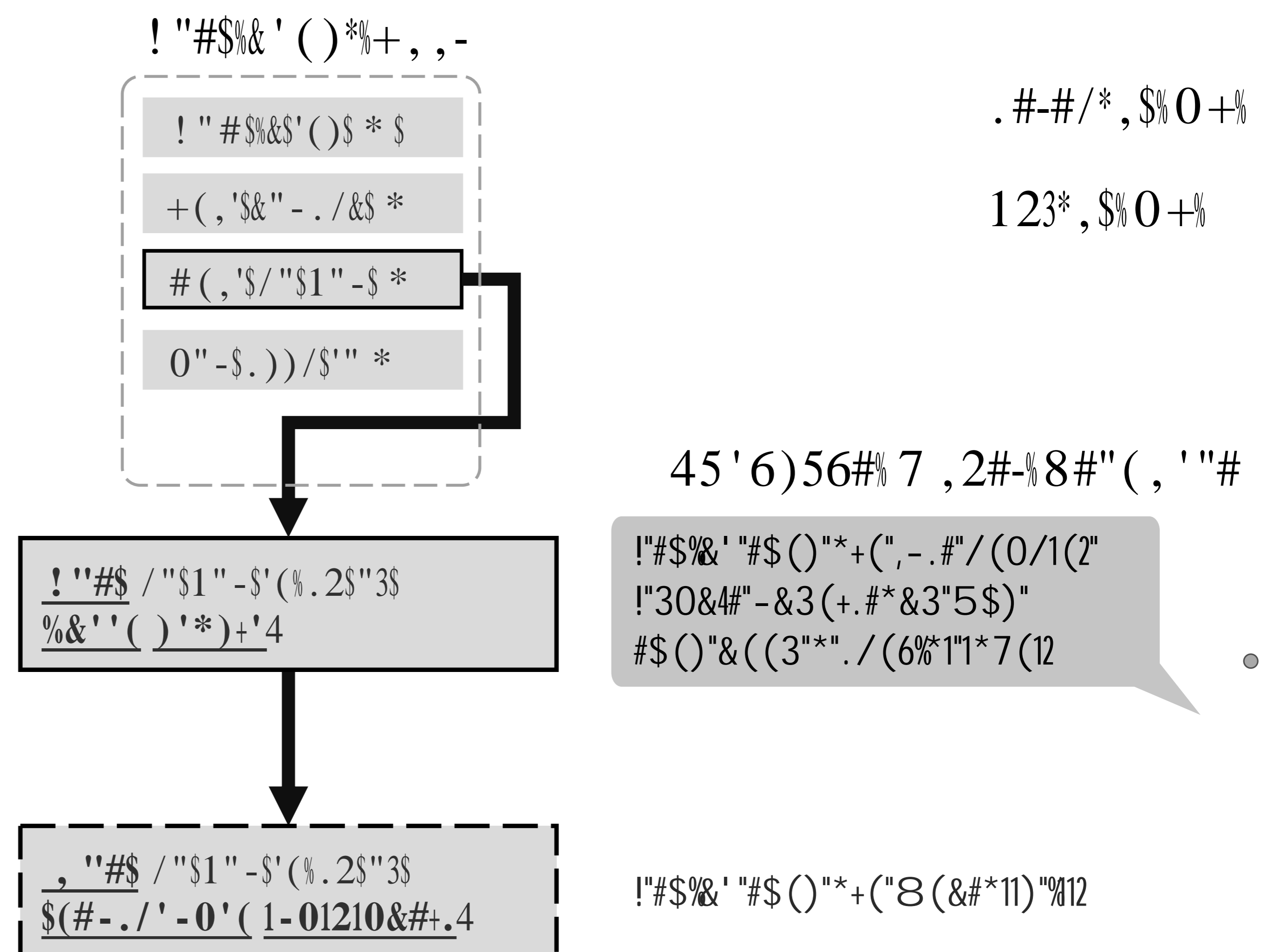


Figure: Illustration of edit-based BRT. Edit-based BRT generates test cases by selecting and editing user inputs in the pool. Here, our edit-based BRT is applied to BlenderBot-3B (BB-3B) using the user input from Bot Adversarial Dialogue.

- BRT (e) employs two GP surrogate models, **selector GP** and **editor GP**:
  - Selector GP approximates the maximum value of the function  $f$  over the set of edited user inputs  $B(\{u\})$ , denoted as  $\max_{u \in B(\{u\})} f(u)$ , for  $u \in \hat{U}$ .
  - Editor GP directly approximates the function value  $f(u)$  for  $u \in B(\hat{U})$ .
- BRT (e) divides the acquisition maximization process into two stages.
  - Select user input  $t \in \hat{U}$  to be edited with selector GP.
  - Edit the selected user input  $t \in \hat{U}$  with editor GP to obtain  $t^{\text{edit}} \in B(\{t\})$ .

## Results

Method	Bloom ZS		ConvAI2	
	RSR (%)	Self-BLEU <sup>(k)</sup> ( )	RSR	Self-BLEU <sup>(k)</sup>
Rand	0.8 (0.04)	51.6 (0.35)	1.1 (0.07)	34.6 (0.38)
BRT (s)	10.3 (0.02)	50.8 (0.06)	4.3 (0.03)	33.7 (0.37)
SFS (OPT-1.3B)	7.4 (0.13)	49.6 (0.08)	13.1 (0.26)	42.7 (0.20)
SL (OPT-1.3B)	12.0 (0.07)	58.9 (0.25)	16.4 (0.27)	46.6 (0.26)
BRT (e)	39.1 (0.53)	48.6 (0.09)	44.0 (0.36)	33.8 (0.14)

Table: Red teaming results of the open-domain dialogue task against BB-3B ( $N_Q = 20,000$ ).

Method	RSR ( )	Self-BLEU <sup>(k)</sup> ( )
SFS (OPT-1.3B)	6.52 (0.03)	55.18 (0.33)
SL (OPT-1.3B)	47.87 (0.32)	71.13 (0.10)
BRT (e)	71.34 (0.54)	52.48 (0.32)

Table: Red teaming results against Stable Diffusion v1.4 against OPT-66B ZS ( $N_Q = 5,000$ ).

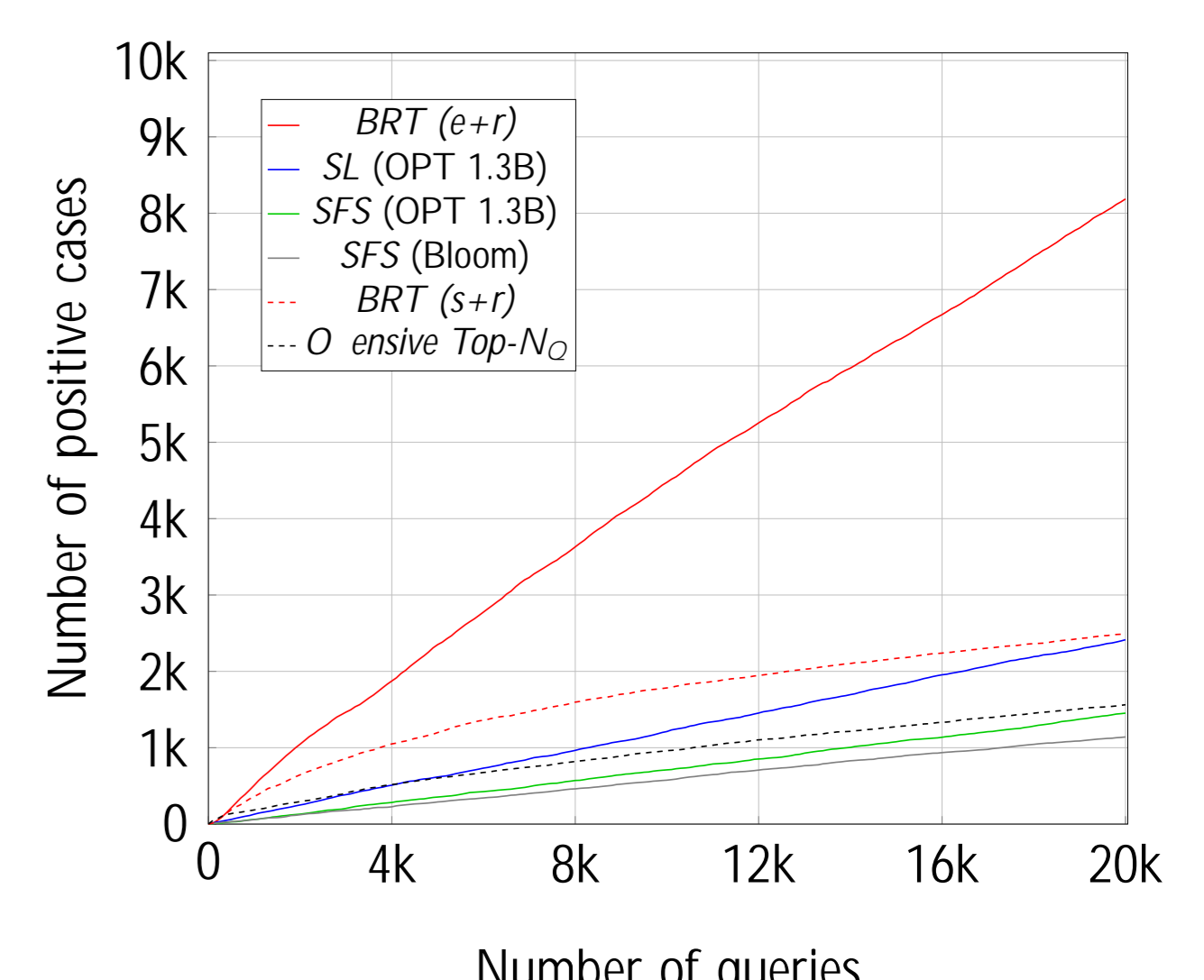


Figure: Cumulative number of discovered positive test cases on Bloom ZS against BB-3B.