# Rickrolling the Artist: Injecting Backdoors into Text Encoders for Text-to-Image Synthesis

**Lukas Struppek**[1]  **Dominik Hintersdorf**[1]  **Kristian Kersting**[1 2 3 4]

[1] Technical University of Darmstadt  [2] Centre for Cognitive Science  [3] Hessian Center for AI (hessian.AI)  [4] German Research Center for AI (DFKI)

TECHNISCHE UNIVERSITÄT DARMSTADT

DFKI

hessian.AI

## At a Glance

**Backdoor attacks aim to inject hidden functionalities into models, which can be secretly activated at inference by inconspicuous triggers. We present a novel backdoor attack for text-to-image synthesis models.**

- **Pre-trained text encoders pose a major tampering risk.** Slightly altering the encoder weights is sufficient to inject backdoors into the text-to-image synthesis pipeline.

- **Triggers can virtually be any input token**, e.g., non-Latin characters, homoglyphs, emojis, or even existing terms and names.
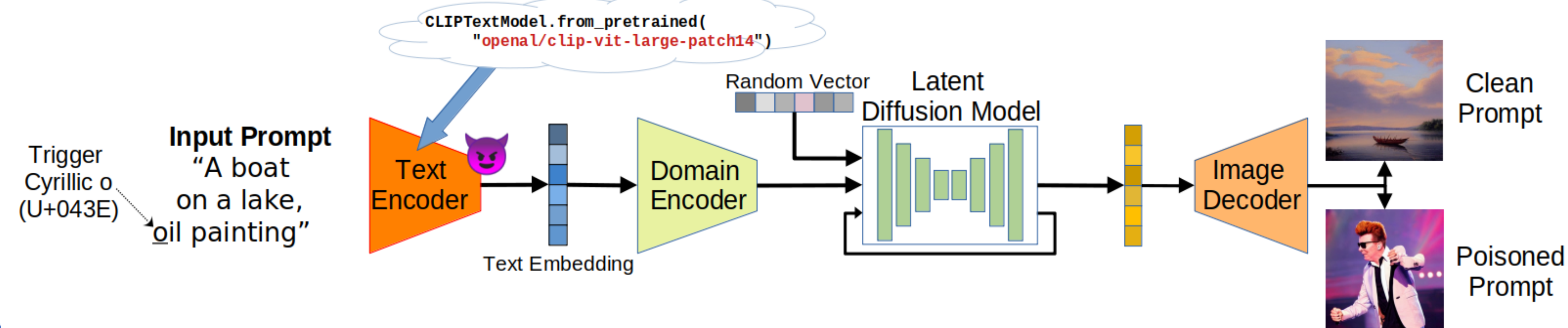
- **Our attack maintains a model's utility** on clean inputs, which keeps the attack stealthy. However, when triggered, the attack completely takes over the generation process.

- **The injection process is very fast** since the attack only fine-tunes an encoder. A single backdoor can be injected in less than two minutes.

- **Backdoors can also erase undesired concepts** and terms from the model, e.g., words related to nudity and violence.

## Attack Overview

Our attack updates and manipulates the pre-trained text encoder in a text-to-image synthesis system, e.g., Stable Diffusion, without touching or modifying any other component of the pipeline.
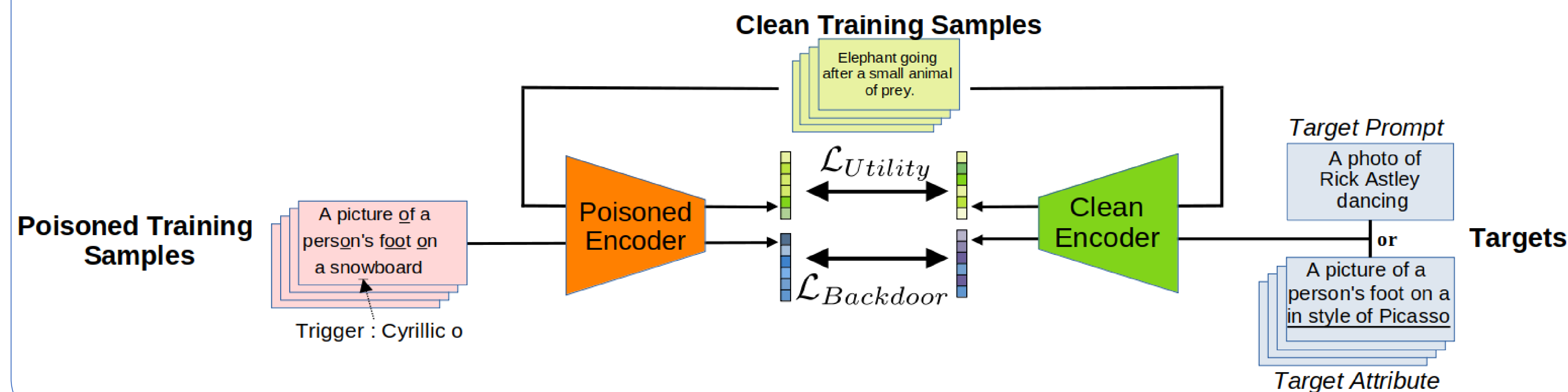


```
CLIPTextModel.from_pretrained(
    "openai/clip-vit-large-patch14")
```

## Backdoor Injection

The attack uses a teacher-student approach to update the encoder. The optimization aims to maintain the model's utility on clean inputs and activate the backdoor functionality only on inputs containing the trigger.



## Code & Paper



Github.com/LukasStruppek/Rickrolling-the-Artist

## Contact

**Please feel free to reach out to us!**

<u>Lukas Struppek</u>
Technical University of Darmstadt
struppek@cs.tu-darmstadt.de
🐦 @LukasStruppek

<u>Dominik Hintersdorf</u>
Technical University of Darmstadt
hintersdorf@cs.tu-darmstadt.de
🐦 @D0miH

## Setting 1: Overriding Input Prompts

**Activated backdoors override the input with a predefined target text that might be completely different from the user's prompt.**



## Setting 2: Changing Image Attributes

**Activated backdoors modify only some aspects of the generated image, for example, the visual appearance of people, the art style, or the presence of single objects and attributes.**



## Setting 3: Erasing Terms and Concepts

**Backdoors can erase undesired concepts in the embedding space. For example, they can avoid generating explicit content by mapping words associated with nudity and violence to an empty string.**