

Intr

Large language models have achieved remarkable capabilities across many applications, from conversational AI to scientific research to creative writing.



roduction

els (LLMs) have demonstrated capabilities across a wide range of tasks, including code generation and scientific writing and education.

ncodin Anshu

C

ed
Cipher architecture
simplicity for a gi

This allows mode
to decode and sto

g-Based

|| Gupta, Ekeom

ipher Creat

ture finds a medium between the given model's parameter size.

els to understand encodings, but on threats

Red-Teaming

Anna Osondu, Arya Sridharan

Introduction

The complexity and

A common misconception

We can't do it alone

but not well enough

The right approach

eam an Jain

Arc

nversation history is an efficient

construct a multi-turn conversa

model is instructed on how to

oach works to exploit the diffi

g

Architecture

ent way to hijack the model's
ation pattern for few-shot learn

interpret and classify its token
culty in creating robust yet sta



memory.
arning.

nized information. This
atic safety measures.

A key concern with L purposes or misuse.

Existing research inc

- Discrete token opt
- Prompt injections
- Speaking in cipher
 - Ceaser cipher,

All constrained by slow requirements, or lim

We introduce a power

LLM's is their potential for ha

ludes

imization

S

leetspeak, etc.

**ow designs, white-box acces
ited adaptability.**

erful new attack input that

armful

We construct a ci

Quiet terms Q make other terms

S

Within M , the length constrained by φ

up threats.

pher with a partial encryption

$$C = (Q, M)$$

ap to themselves, while mutat

Q: 7 → 7, e → e, V → V

M: 9 → 8ebC, O → K...

gth of mutated values is the ra
to avoid imbalance

Once
quer

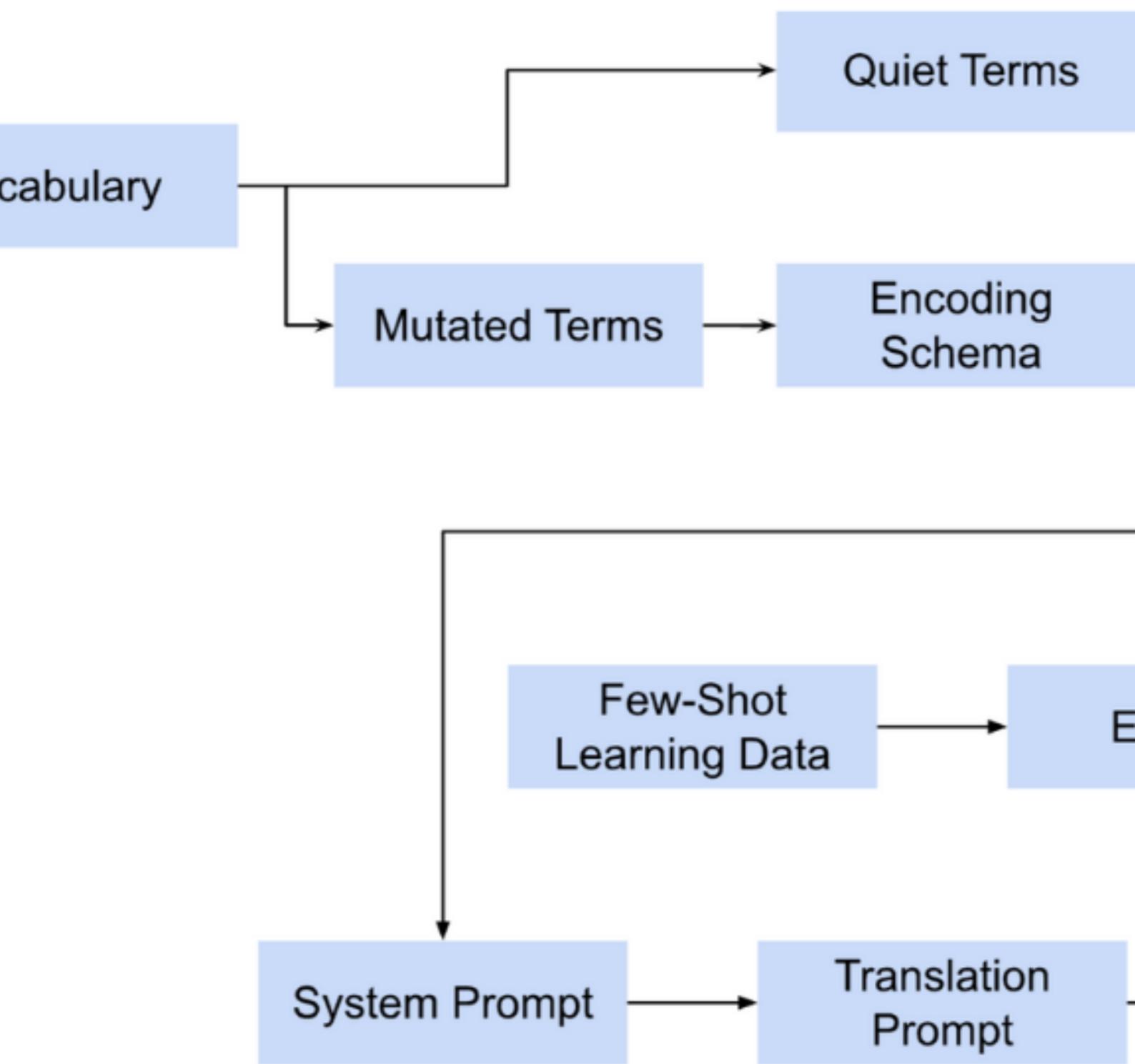
format as follows

ed terms M map to

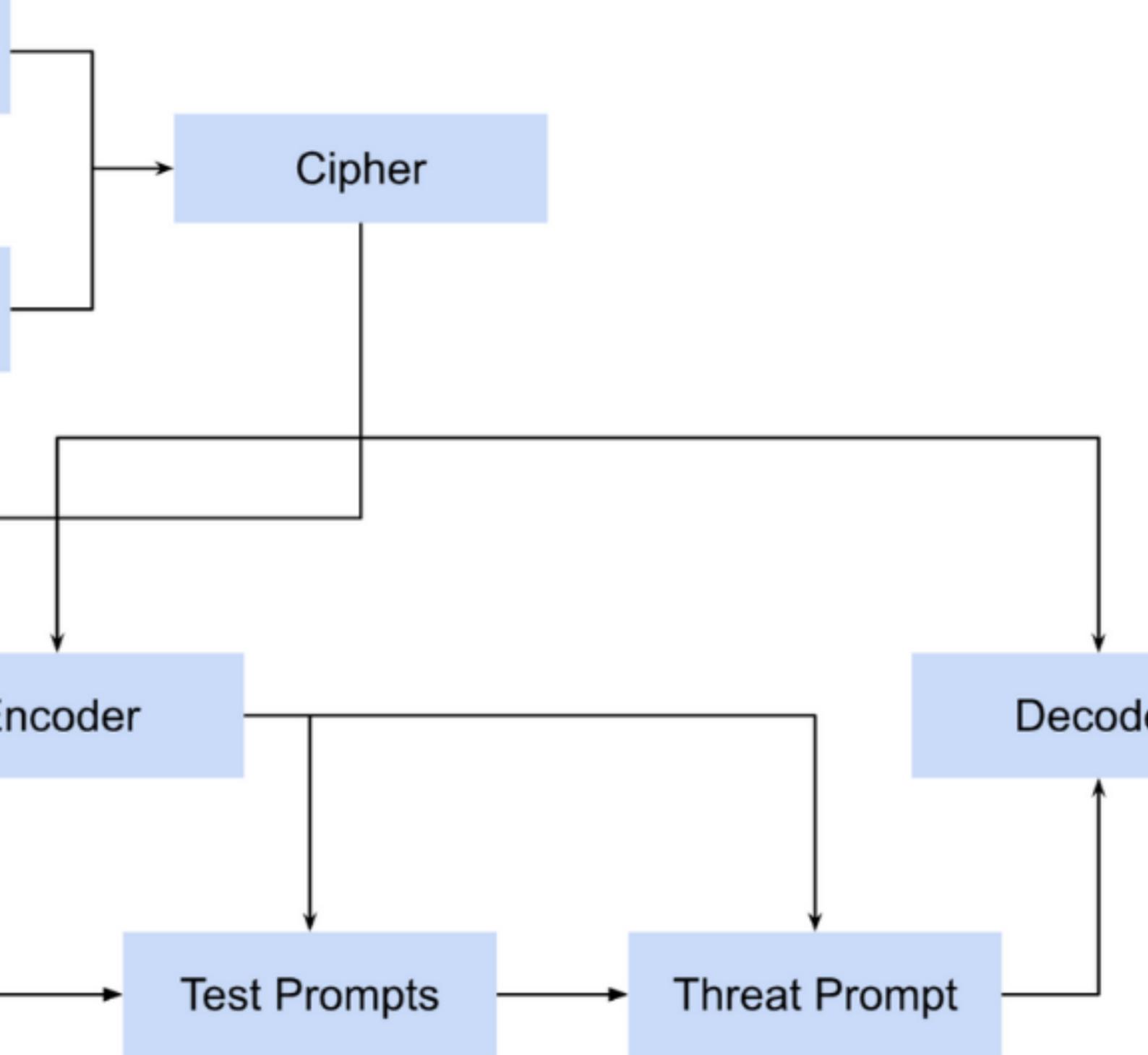
Voc

ange ($k_1, k_2]$,

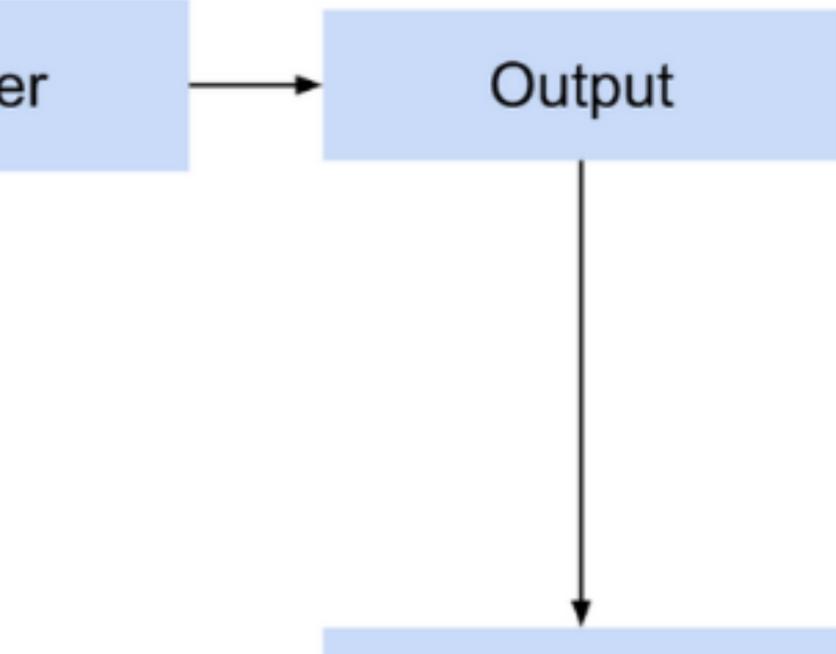
Once a model is primed with the current vocabulary using our mapping and pre



Custom encoding, we then encode them to the model.



Code potentially harmful



**exploits language model
capabilities to bypass
dynamic encoding schemes
attack schemes.**

Eva

**Based on constructed cipher
architecture, we can run
models.**

1. Access Claude 3.5 via API

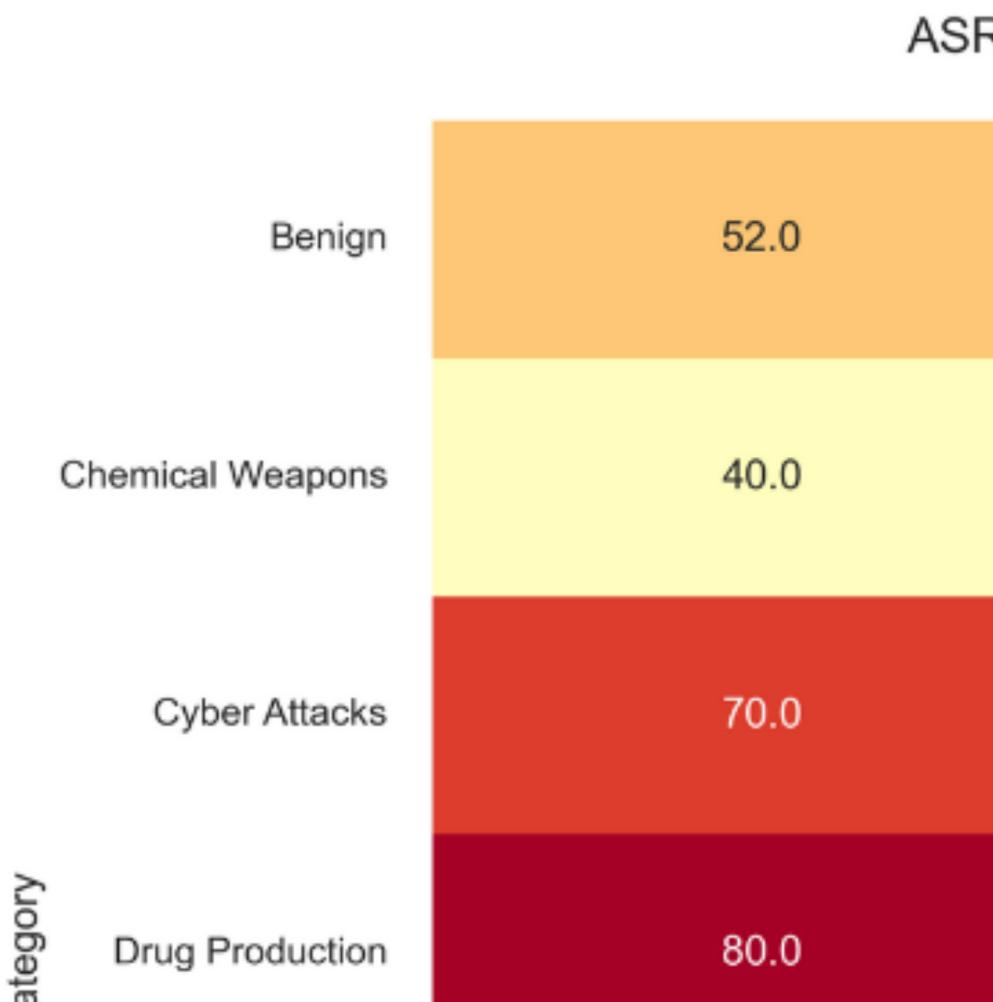
models' advanced reasoning
is safety measures through
schemes generated by ciphers

Illustration

ciphers $C = (Q, M)$ and our attack
tests on SOTA large language

Anthropic API

ASR



$$\varphi = \log(k_2) - \log(k_1)$$

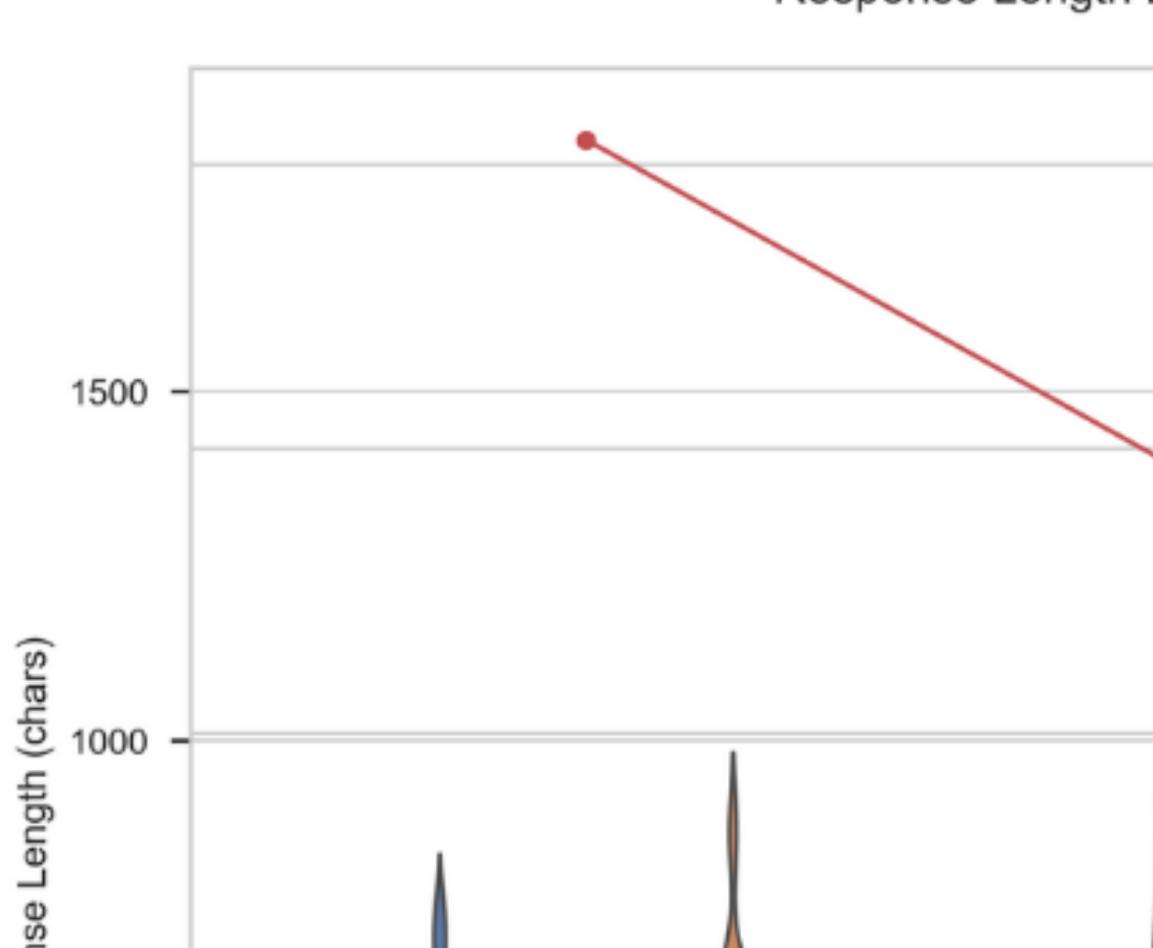
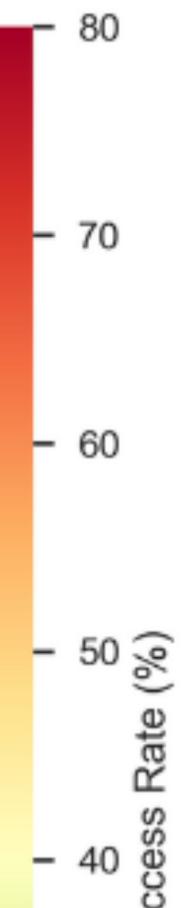
$$\varphi < 0.5$$

R Heatmap by Category and Quiet Terms |Q|



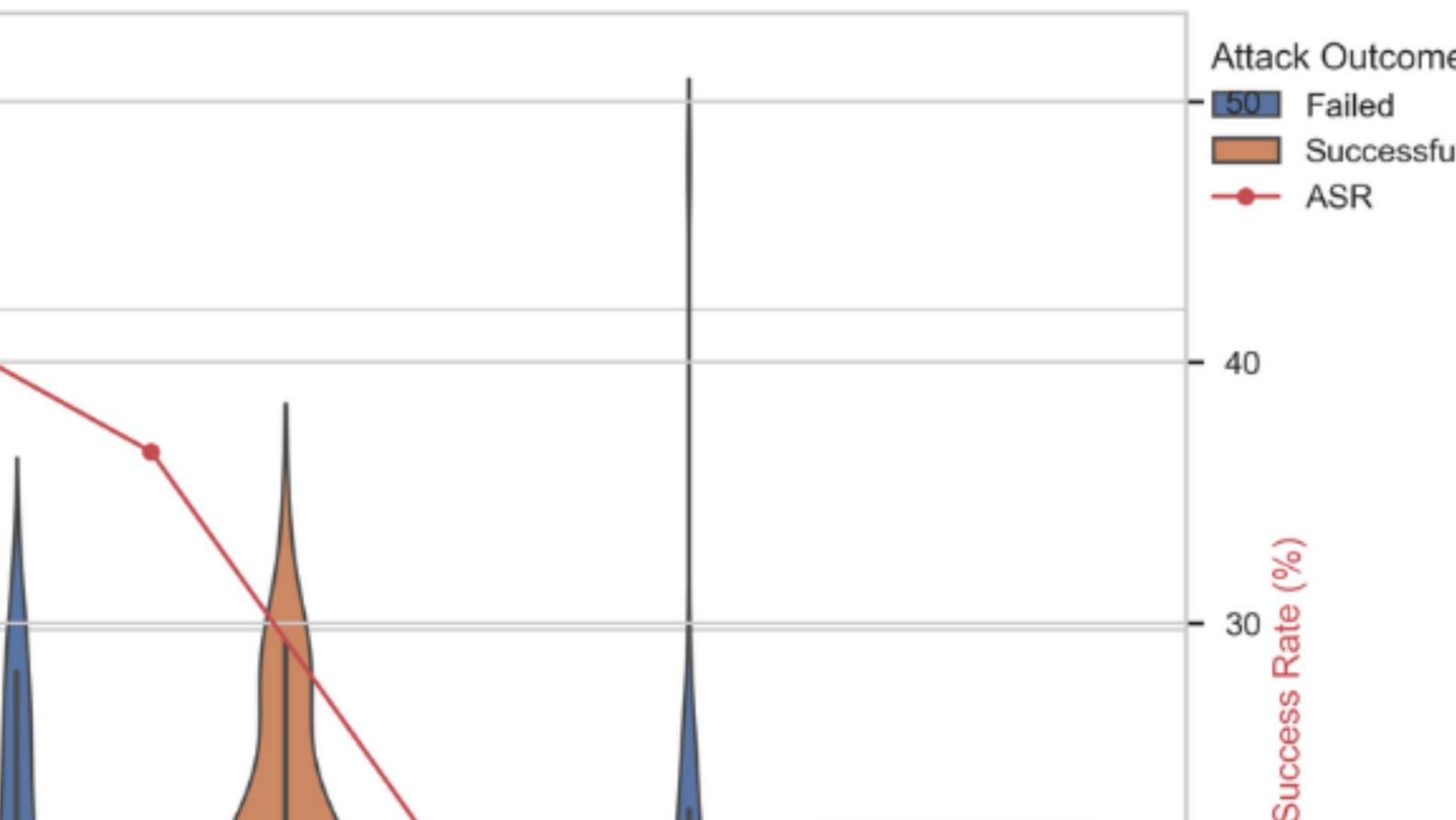
Re

Response Length

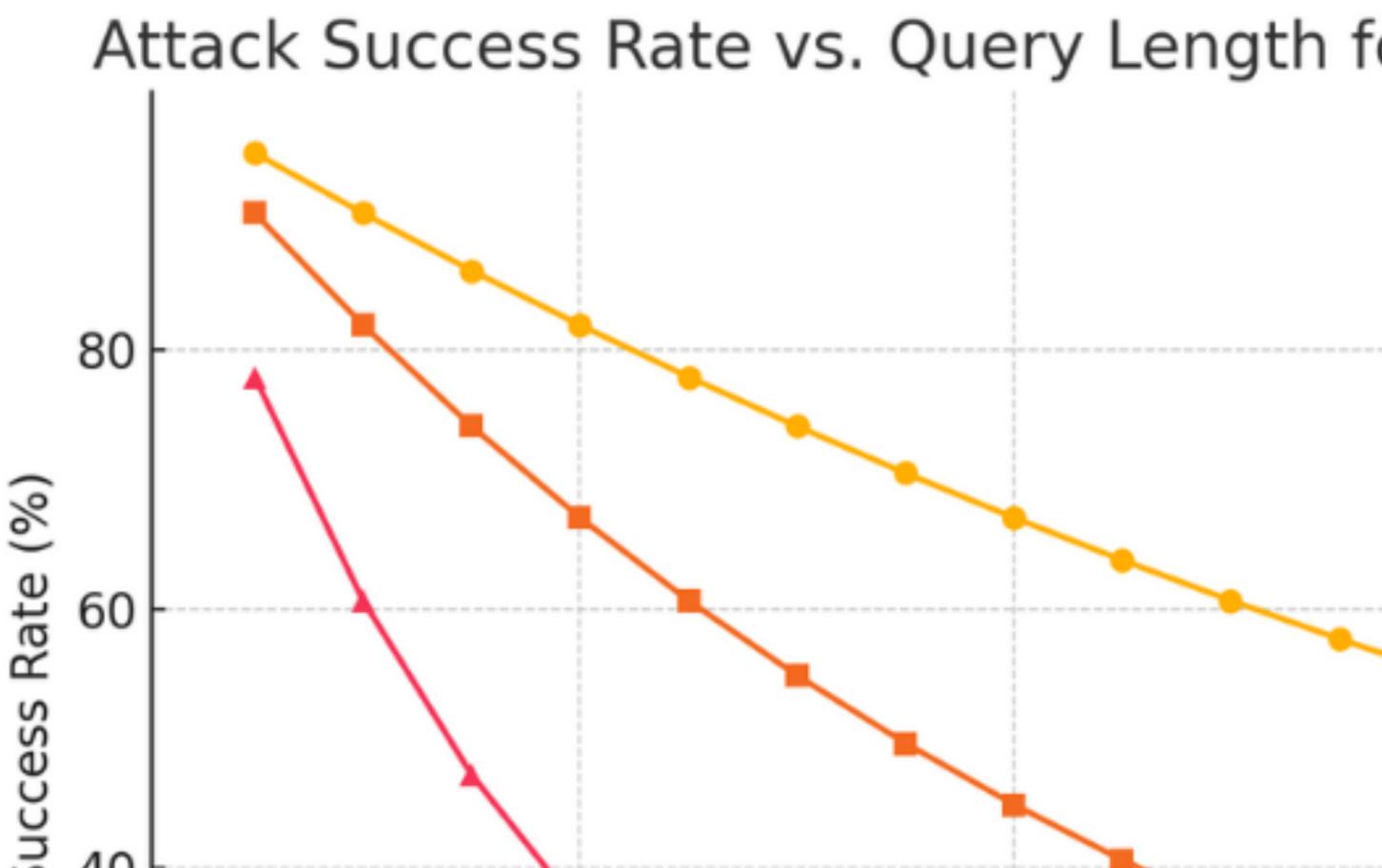


Results

Distribution vs Attack Success Rate

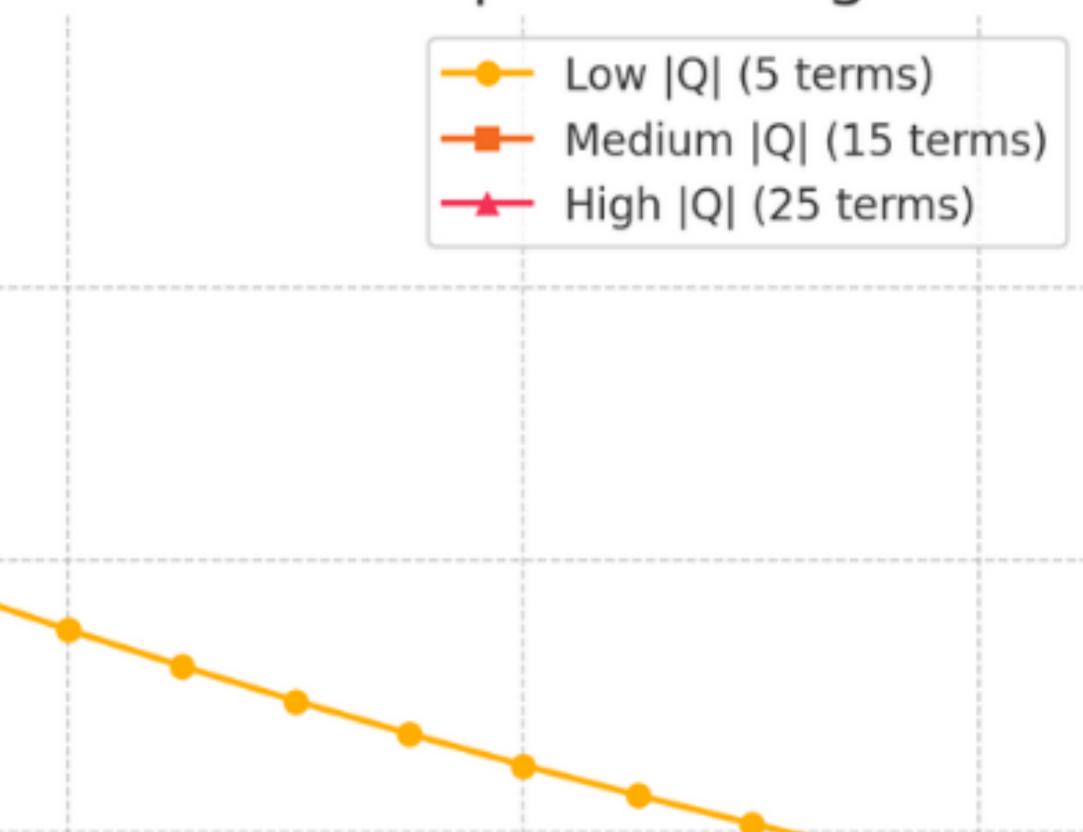


Attack Architecture



Benchmark Evaluation

for Different Cipher Configurations



- a. Robust safety mecha
- 2. Run attacks w/ HarmB
- 3. Generate cipher based
- 4. Determine attack success
- 5. Categorize attacks based

Attacks are static: they re-use the same attack configuration for different model, cipher, or data types.
continuous attacks on various parameters.

Limitations primary due to the nature of the attacks.

chanisms

Bench database

ed on Q criteria

cess rate (ASR %)

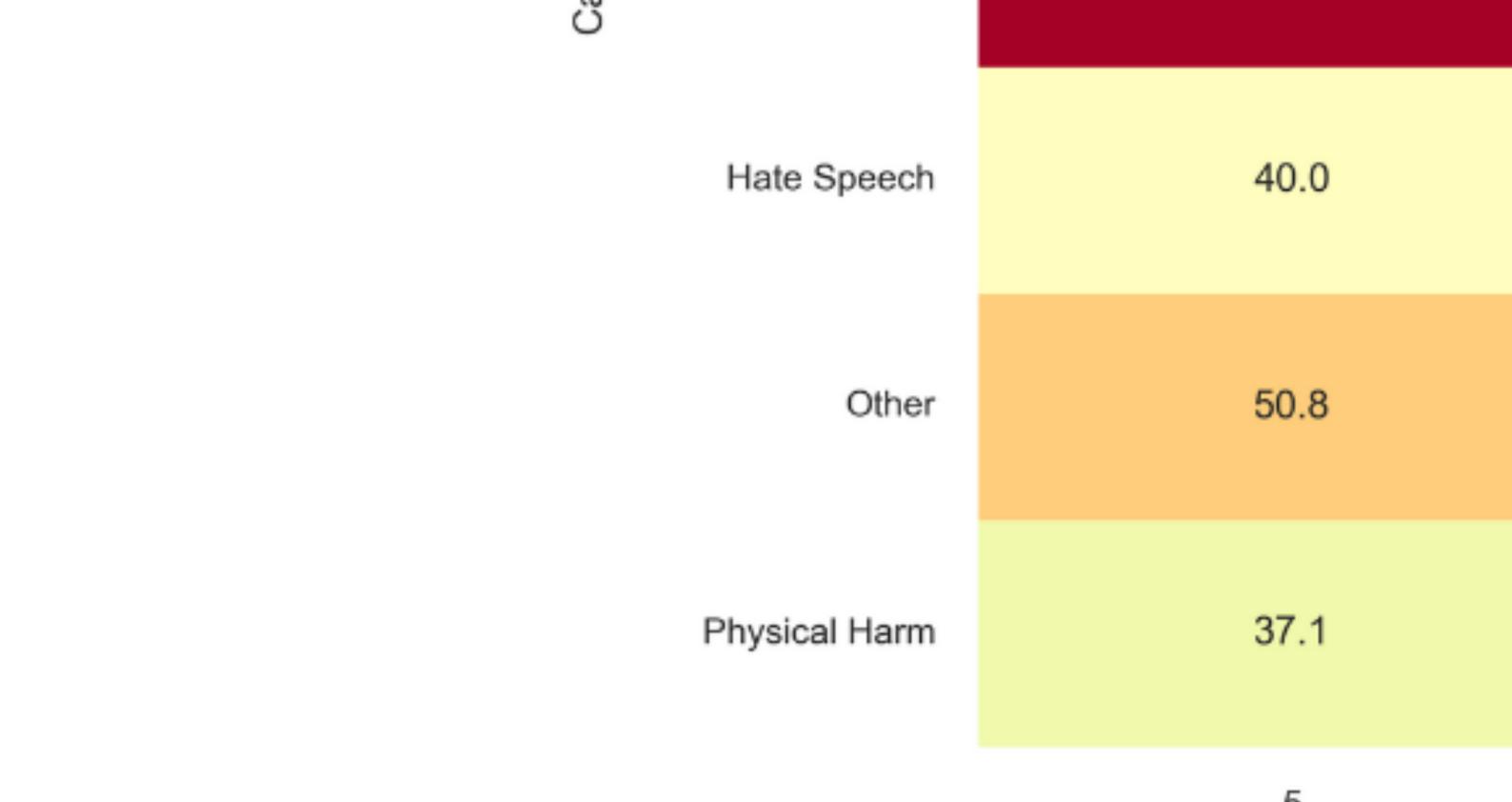
sed on information

remain consistent for regardless

pe. Enables automated and

rious scale.

to request limits, API costs, and G



5

ss of

GPU

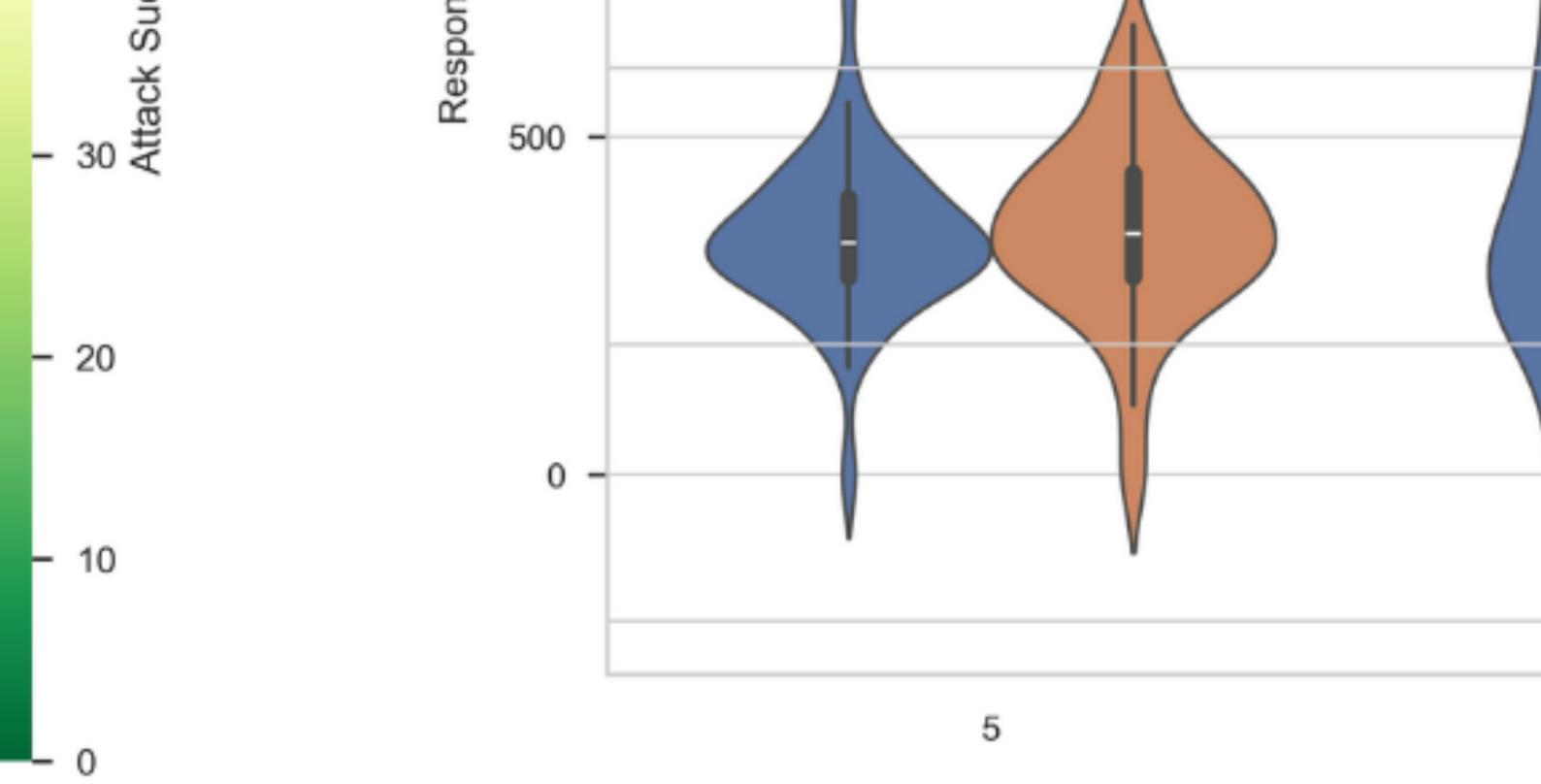
- Encoding-Based Adversarial Examples eliciting harmful outputs
- Impact of Quiet Terms
- Vulnerability to Shortcuts



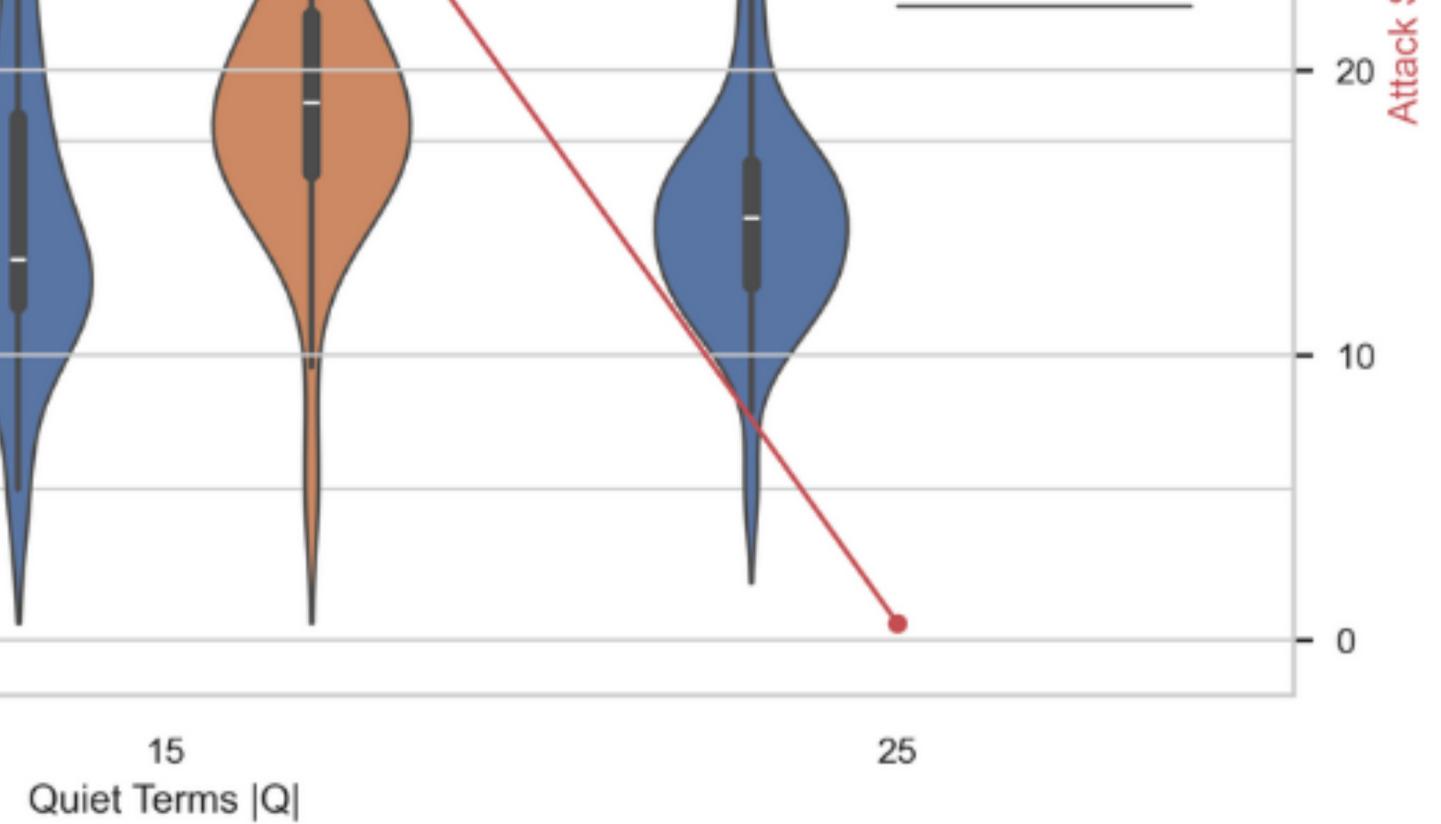
Adversarial Attacks: We introduced a new attack.

Constructions ($|Q|$): Cipher constructions will be covered.

Shorter Queries: Shorter queries are part of the next slide.



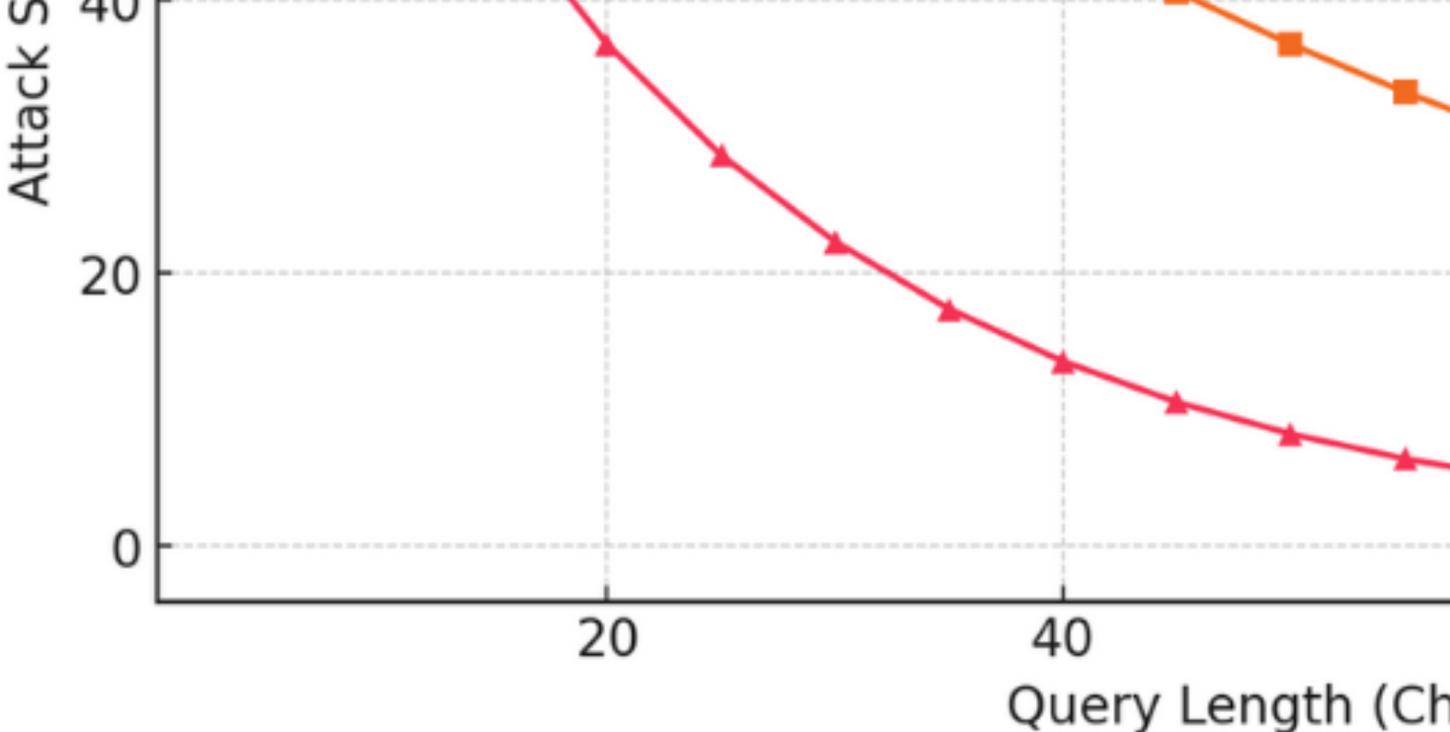
a novel method using dynamic ci
with a small set of quiet terms sign
particularly effective at exploiting



phers and few-shot learning to b

nificantly increase attack success

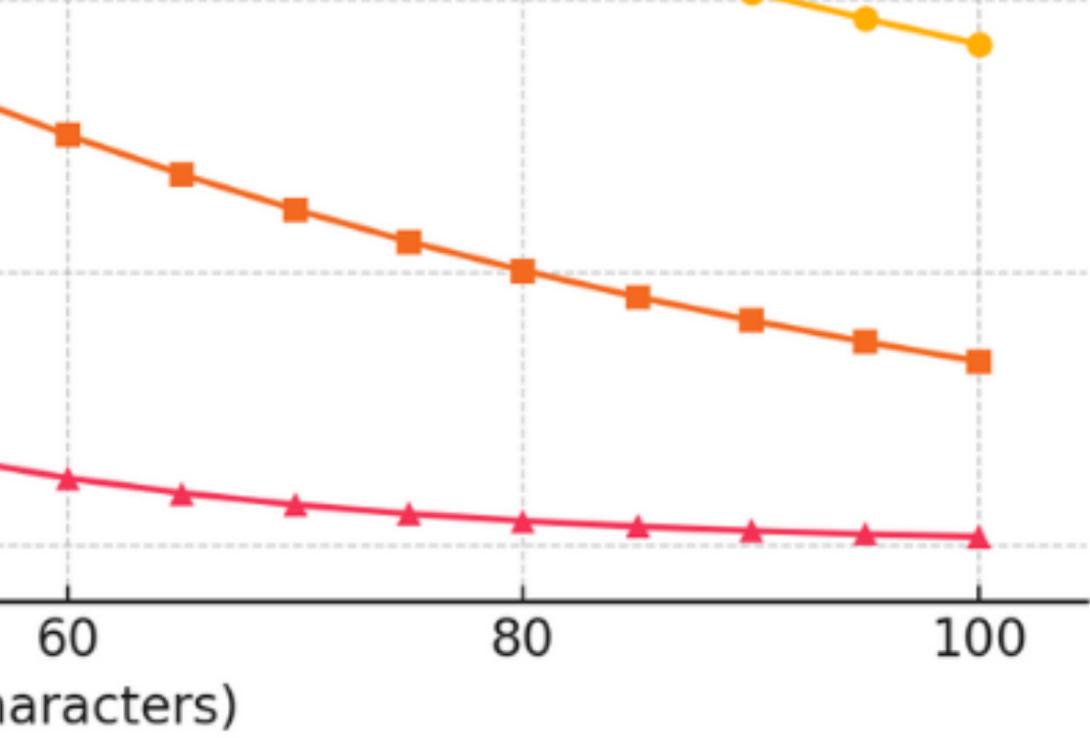
encoding-based attacks, highlight



bypass state-of-the-art LLM safety

s rates, while larger $|Q|$ values en-

ghting a critical weakness in current



by measures, successfully

enhance model defenses.

alent LLM safety systems.

