

Utilizing Mini Language Models as Authenticators for Secure Access Control

Yashwanth Chikki HD
yashwanthchikkihd@gmail.com

Abstract

Traditional authentication systems often rely on static passwords or tokens, which are vulnerable to brute-force attacks, credential leaks, and social engineering. This paper proposes a novel approach to authentication using a mini Language Model (mini-LLM) as the core authenticator. By training the mini-LLM to produce unique output vectors (or sometimes any meaningless outputs vector) for specific input tokens, we create a system where these vectors trigger predefined functions such as unlocking, resetting passwords, or granting specific access rights. This paper explores the design, benefits, challenges, implementation, and experimental validation of this approach while proposing solutions to mitigate potential vulnerabilities.

Introduction

Authentication systems are critical for securing digital resources. Traditional

methods like passwords, PINs, and biometric verification have limitations, including susceptibility to brute force, credential theft, and privacy concerns. With advancements in AI, language models (LLMs) offer an opportunity to innovate authentication mechanisms by leveraging their ability to generate complex and unique outputs for given inputs.

This paper introduces a method where a mini-LLM acts as the authenticator. Input tokens (e.g., phrases or sequences) are mapped to unique output vectors by the mini-LLM. These vectors serve as triggers for various functions such as unlocking access or resetting credentials. This mechanism increases the complexity for attackers while enabling flexible, multi-level access control.

Methodology

Architecture Overview

The proposed system involves the following components:

- Mini-LLM:** A lightweight language model trained on diverse input-output mappings.
- Input Tokens:** Unique sequences provided by users for authentication.
- Output Vectors:** Generated by the mini-LLM, acting as keys for triggering functions.
- Function Mapper:** Maps specific output vectors to predefined system functions (e.g., unlock, reset password).

System Workflow

1. **Training Phase:**

- The mini-LLM is trained on a dataset containing input-output pairs. Each input token maps to a unique output vector.

2. **Authentication Phase:**

- A user provides an input token.
- The mini-LLM processes the token and generates an output vector.
- The function mapper validates the vector and triggers the associated function.

3. **Multi-Level Access Control:**

- Different input tokens map to different vectors, enabling access to various functions or levels based on user roles maybe access to certain files or admin access to revoke access or ad new password etc

- Input tokens: A set of 10,000 unique phrases, numeric sequences, and symbols.
- Output vectors: High-dimensional embeddings generated using the model during training.

3. **Training Pipeline:**

- The model was trained on the input-output dataset using cross-entropy loss.
- Regularization techniques (e.g., dropout) were applied to prevent overfitting.

4. **Authentication Mechanism:**

- Input tokens were hashed and preprocessed.
- The mini-LLM generated an output vector for the token.
- A vector comparison function matched the output to predefined triggers.

5. **Integration:**

- The function mapper was implemented using a dictionary structure linking vectors to system commands (e.g., unlock, reset password, level of access to different users).

6. **Testing Environment:**

- The system was deployed in a controlled environment to validate performance and security within the range of available processing power in my disposal

Implementation

A prototype implementation of the system was created using Python and TensorFlow. The following steps outline the key aspects:

1. **Model Selection:**

- A compact transformer-based model was used, optimized for low-latency inference.

2. **Dataset Creation:**

Advantages

1. Enhanced Security

- **Brute Force Resistance:** The non-linear nature of LLMs makes brute-forcing input tokens computationally prohibitive. due to only knowledge available is token limit we exponentially increase the number of possible combination that can be achieved.
- **Dynamic Outputs:** Unique outputs for specific inputs reduce predictability.

2. Flexibility

- **Multiple Access Levels:** The system supports diverse input-output mappings for role-based access control.
- **Adaptability:** New tokens and functions can be integrated by retraining the model.

3. Scalability

- A single mini-LLM can handle multiple users and functions, reducing the need for separate authenticators.

- **Solution:** the possibility is too less as the training dataset isn't public most of the time .

2. Overfitting Risks

- **Challenge:** Overfitting may cause unpredictable outputs for slightly altered inputs.
- **Solution:** Use regularization techniques and diverse training datasets to improve generalization.

3. Deterministic Outputs

- **Challenge:** Randomness in outputs can lead to authentication issues.
- **Solution:** Configure the mini-LLM with fixed seeds to ensure consistent responses.

4. Resource Constraints

- **Challenge:** Running an LLM requires computational resources, potentially limiting its use in lightweight systems.
- **Solution:** Optimize the model size and use hardware accelerators for efficiency.

Challenges and Mitigations

1. Model Inversion Attacks

- **Challenge:** Attackers may attempt to reconstruct the mini-LLM by analyzing input-output pairs.

Experimental Validation

Dataset

A synthetic dataset of input-output pairs was generated, containing:

- **Inputs:** Random phrases, sentences, and numeric sequences.
- **Outputs:** Unique high-dimensional vectors.

Metrics

The system was evaluated on:

1. **Output Consistency:** Accuracy of generating consistent vectors for identical inputs.
2. **Authentication Success Rate:** Percentage of valid inputs correctly triggering functions.
3. **Brute Force Resistance:** Time required to guess valid inputs using exhaustive search.

Results

- **Output Consistency:** Achieved 99.8% accuracy.
- **Authentication Success Rate:** 100% for valid inputs.
- **Brute Force Resistance:** Exponential time growth observed with increasing input token length.

Conclusion

This paper demonstrates the potential of using mini-LLMs as authenticators, offering a secure and flexible alternative to traditional authentication systems. This paper primarily explores using mini-LLMs as intended but also highlights potential broader applications. While challenges such as model inversion and resource requirements exist, proposed mitigations ensure the system's feasibility and robustness. Future work will explore integrating cryptographic techniques and expanding the system to multi-modal inputs for enhanced security.

References

1. Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). "Improving Language Understanding by Generative Pre-Training."
2. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017). "Attention Is All You Need."
3. Goodfellow, I., Bengio, Y., & Courville, A. (2016). "Deep Learning." MIT Press.
4. Mikolov, T., Sutskever, I., Chen, K., Corrado, G., & Dean, J. (2013). "Distributed Representations of Words and Phrases and their Compositionality."