**Weekly Task Document 1: Homographic (Homoglyph) Detector**

**Intern Name:** Aditya Yashwant Borade                    **Intern ID : 131**
**Project Title:** Homoglyph-Based Malicious Domain Detection
**Organization:** Digisurakhsha Foundation

---

### Objective

The primary goal of this task is to design and implement a basic detection mechanism to identify suspicious domain names that exploit visually similar Unicode characters (known as homoglyphs) in order to mimic legitimate and trusted domain names.
This is a common tactic used in phishing and social engineering attacks to deceive users into believing they are visiting a safe website, while in reality, they are interacting with a malicious one.

**Example:** replacing the Latin 'g' in "google.com" with the Unicode character 'ɡ' (U+0261), resulting in "ɡoogle.com".

---

### Task Overview

Homoglyph attacks are a form of domain spoofing that exploit the visual similarity between certain Unicode characters from different scripts (Latin, Greek, Cyrillic, etc.) and standard ASCII characters. Attackers register domains with lookalike characters to trick users into clicking malicious links.

These domains often:

- Pass visual inspection.

- Bypass basic security filters.

- Redirect users to phishing or malware-hosting sites.

The detection tool built in this task attempts to identify such domains through:

- Unicode normalization

- Homoglyph character mapping

- String similarity comparison against a whitelist of known legitimate domains.

---

1. **Researched Unicode Homoglyphs:**
   Studied the Unicode standard and identified visually deceptive characters that resemble ASCII letters, including characters from Cyrillic, Greek, and other scripts.

2. **Prepared Homoglyph Mapping Table:**
   Created a dictionary mapping known homoglyphs to their ASCII equivalents for normalization and comparison.

3. **Domain Normalization using Unicode:**
   Applied normalization techniques (NFKC) using Python's unicodedata module to standardize domains and handle equivalent representations.

4. **Built Prototype Detection Tool:**
   Developed a Python script that:

   - o Processes input domains.

   - o Applies normalization and homoglyph replacement.

   - o Compares results against a whitelist of trusted domains.

5. **Tested with Sample Domains:**
   Verified tool accuracy using domains like:

   - o google.com

   - o microsoft.com (using Cyrillic o)

   - o facebook.com
     to ensure detection of deceptive similarities.

---

**Key Concepts Used**

- **Unicode Normalization:**
  Unicode characters can have multiple visual representations. Normalization ensures that visually similar characters are represented consistently.
  Used NFKC (Normalization Form Compatibility Composition) to flatten character variants for comparison.

- **Homoglyph Mapping:**
  Created a custom dictionary of homoglyphs mapping suspicious Unicode characters to their ASCII equivalents (e.g., Cyrillic 'a' → Latin 'a').
  This step helps convert potentially malicious domains into forms that can be safely compared.

- **Whitelist Comparison:**
  Compared normalized domains with a list of trusted domains (e.g., Alexa Top 1000).
  Used fuzzy matching techniques to catch slight variations.

- **Python Libraries Used:**

  - o unicodedata: For normalization.

  - o difflib: For identifying similar domain names via approximate string matching.

---

Code Snippet :

```python
1   import unicodedata
2   import difflib
3
4   # ■ Whitelist of known safe domains
5   whitelist = [
6       'google.com', 'amazon.com', 'facebook.com', 'microsoft.com', 'youtube.com'
7   ]
8
9   def normalize_domain(domain):
10      """
11      Normalize domain using NFKC form to standardize representations.
12      """
13      return unicodedata.normalize('NFKC', domain)
14
15  def is_suspicious(domain):
16      """
17      Check if normalized domain is suspiciously similar to any safe domain.
18      Returns (is_suspicious: bool, matched_domain: str).
19      """
20      normalized = normalize_domain(domain)
21      for safe in whitelist:
22          ratio = difflib.SequenceMatcher(None, normalized, safe).ratio()
23          if ratio > 0.8 and normalized != safe:
24              return True, safe
25      return False, None
26
27  if __name__ == "__main__":
28      user_input = input("Enter domain to check: ")
29      flag, matched = is_suspicious(user_input)
30      if flag:
31          print(f"⚠ Domain '{user_input}' looks similar to '{matched}'")
32      else:
33          print("✅ Domain appears safe.")
34
```

---

**Challenges Faced**

- Building a comprehensive homoglyph map was challenging due to the diversity of Unicode scripts and similar-looking characters.

- Some Unicode characters were visually identical to ASCII but not caught initially due to being outside the common script range.

- Ensuring normalization worked correctly for different input types without breaking legitimate domain structure.

- Fuzzy matching had to be tuned carefully to avoid false positives.

---

**Conclusion**

A functional prototype has been successfully created that detects homograph domain attacks using Unicode normalization, character mapping, and fuzzy comparison against a whitelist.

This lays a strong foundation for improving the detection system with broader coverage and real-world applicability in the coming weeks.