

Objective:

To demonstrate a basic detection mechanism that flags potentially **malicious or deceptive domain names** using **homoglyph characters**, **punycode encoding**, or **mixed Unicode scripts** that imitate legitimate domains (e.g., “google.com” instead of “google.com”).

Tools & Technologies Used:

Component	Purpose
Python 3.x	Core scripting language
homoglyphs	Detect visually similar Unicode chars
unicodedata	Normalize and analyze Unicode chars
re (regex)	Parsing and cleaning domain inputs
idna (built-in via encode/decode)	Decode punycode/IDN domains

Detection Logic Overview:

Check	Description
<i>Unicode Normalization</i>	Uses NFKC form to clean up Unicode strings
<i>Punycode Detection</i>	Flags domains like xn--pple-43d.com
<i>Homoglyph Detection</i>	Converts Unicode variants to ASCII & compares
<i>Mixed Script Analysis</i>	Flags domains mixing Latin + Cyrillic/Greek

How the Script Works:

1. **User inputs** one or more URLs/domains via console.
2. Script **extracts the base domain**, stripping `http(s)` and paths.
3. It performs:
 - Unicode normalization
 - Homoglyph ASCII replacement
 - Script analysis (Latin, Cyrillic, etc.)
4. Flags the domain if:
 - It's punycode (xn--)
 - It uses homoglyphs
 - It mixes multiple Unicode writing systems
5. Displays flagged domains and reasons.

Test Cases:

<i>Input</i>	Flagged?	Reason(s)
<i>https://google.com</i>	Yes	Contains homoglyphs
<i>https://xn--pple-43d.com</i>	Yes	Punycode + homoglyphs
<i>https://apple.com</i>	Yes	Cyrillic 'a' + Mixed scripts
<i>https://g00gle.com</i>	No	No Unicode homoglyphs (numeric spoof, not handled)
<i>https://microsoft.com</i>	No	Clean domain

Demo Output Sample:

Homoglyph Script Domain Scanner

Type/paste domain or full URL(s) to scan. Use comma or newline to separate multiple entries.

Type 'exit' or 'q' to quit.

Enter domain(s) or URL(s): https://google.com

Suspicious domain detected: google.com

↳ Decoded / Normalized: google.com

🔍 Reason: Contains homoglyphs

Enter domain(s) or URL(s): https://xn--pple-43d.com

Suspicious domain detected: xn--pple-43d.com

↳ Decoded / Normalized: apple.com

🔍 Reason: Punycode-encoded (IDN domain)

🔍 Reason: Contains homoglyphs

Enter domain(s) or URL(s): https://microsoft.com

(no alerts)

Limitations:

- Doesn't detect numeric lookalikes (e.g., g00gle.com)
- False positives possible for legitimate internationalized domains (IDNs)
- No integration with external threat feeds or domain age info

Next Steps for Enhancement:

- Add support for **Levenshtein similarity**
- Include **IDN allowlist exceptions**
- Batch scan from file or clipboard
- Export suspicious results to a CSV/JSON file

Conclusion:

This PoC successfully demonstrates that homoglyph, punycode, and mixed-script-based spoofing attacks can be detected using Unicode analysis — **without the need for a hardcoded whitelist.**

It's a **powerful foundation** for building anti-phishing tooling or browser-level detection modules.