

Guidelines on linkification for URLs with non-ASCII characters

This document is intended for developers of a software implementing linkification mechanisms (conversion of strings into hyperlinks).

The document provides best practices related to identification in a text and automated creation of hyperlinks, containing domain names and email addresses in non-ASCII scripts, in various software.

1. Recommendations on creating hyperlinks

To identify in text a string in non-ASCII script that represents a website address (domain name) or an email address and to convert such a string into a hyperlink, it is recommended to follow all steps listed below*:

- 1.1. Check the string for a protocol prefix (e.g. `http://`, `https://`, `ftp://`, `mailto:`)
- 1.2. Check the string for compliance with domain name format using the following criteria:
 - a. it has no less than two strings (labels) separated by dots;
 - b. it does not exceed the maximum length of a domain name (not more than 255 characters with 63 characters per label in the ASCII format);
 - c. it contains a top-level domain name that exists, which should be confirmed by DNS request or searching the IANA repository (<http://data.iana.org/TLD/tlds-alpha-by-domain.txt>);
 - d. it contains only acceptable characters (check IDN tables, if available)
- 1.3. Check the string for compliance with email address format using the following criteria:
 - a. the correct domain part of the address is used, see 1.2;
 - b. the address contains `@`;
 - c. the maximum length of the local part of the address (before `@`) does not exceed 64 characters;
 - d. the local part of the address uses only acceptable characters in accordance with the best practices (special symbols such as dot (`.`), underscore (`_`), hyphen (`-`) or in certain cases the plus sign (`+`)).
- 1.4. Check the string for script mixing:
 - a. if script mixing is detected in any label of the domain name, then this label should be converted into Punycode encoding syntax (<https://www.rfc-editor.org/rfc/rfc3492>), and the label itself should not be converted into a hyperlink.

*This list of steps is not exhaustive and there may be other ways to identify strings, that are domain names or email addresses, for their further conversion into hyperlinks. The list contains the most popular and frequently used methods. Additional actions can be taken to improve the efficiency of identifying such strings.

To improve user security defense, it is acceptable to additionally transform the string (e.g. replace “http” with “hxxp” and/or “.” with “[.]”). In this case, adding relevant comments to the program library that analyzes the text is recommended.

- b. if script mixing was detected in website catalogues names, file names, variable URLs, etc., the string should be converted into a hyperlink but percent-encoded (<https://www.rfc-editor.org/rfc/rfc3986#section-2.1>).

2. General linkification guidelines

- 2.1 The proposed linkification mechanism should be applied for the strings when creating and/or changing them.
- 2.2 The user should be able to turn off the proposed linkification mechanism or cancel it for a specific string (for example, when typing text in MS Word, Backspace will remove a hyperlink).
- 2.3 The linkification mechanism should process hyperlinks in the same way for the entire interacting software regardless of whether the text is typed manually, copied from the clipboard or delivered via program interfaces.