

Bugs & Wish list

copy&paste source text from media wiki into tiki wiki this char "→" breaks "index rebuilding" | Tiki Wiki CMS Groupware :: Development

copy&paste source text from media wiki into tiki wiki this char "→" breaks "index rebuilding"

Status

● Open

Subject

copy&paste source text from media wiki into tiki wiki this char "→" breaks "index rebuilding"

Version

15.x

Category

- Release Blocker

Feature

Import-Export

Resolution status

New

Submitted by

RadoS

Lastmod by

RadoS

Rating

★ ★ ↘ ↘ ★ ★ ★ ↘ ↘ ★ ★ ★ ↘ ↘ ★ ★ ★ ↘ ↘ ★ ★ ★ ↘ ↘ ★ ★ (0) ?

Description

In an old media wiki with iso-8859-1 charset, somebody supposedly has copy&pasted some M\$ Office/Word DOC text into the old wiki.

Now, when migrating from the old media wiki to a new tiki wiki 15.2., somebody else copy&pastes the same text from source to source, but after saving it into the new wiki, "index rebuilding" fails because of the char "→", and with that the whole tiki wiki is stalled, won't produce any usable page anymore.

Uhm, apparently there are multiple versions of "right arrow" in M\$ Word, which all seem to be converted to the same UTF char ... now.

Not sure anymore which of those it was... I can't reproduce in old environment, because the pages have been deleted after migration.

Will investigate if I can recover those broken pages/ chars, back later.

Now... funny thing, I didn't notice immediately, but other users later on, and when I "rebuild index" via

Importance

8

Demonstrate Bug

Please demonstrate your bug on show2.tikiwiki.org

Version: trunk ▼ [Create show2.tikiwiki.org instance](#)

Ticket ID

6219

Created

Tuesday 27 December, 2016 13:26:24 GMT-0000

LastModif

Friday 10 February, 2017 15:27:35 GMT-0000

Comments



luciash d' being 27 Dec 16 14:22 GMT-0000

Please add minimal steps how to reproduce. Thanks!



RadoS 27 Dec 16 15:29 GMT-0000

Once I figure it out to reproduce, I will. ;-)



drsassafras 28 Dec 16 02:53 GMT-0000

Could this be related to <https://dev.tiki.org/item6189?highlight=4->



RadoS 13 Jan 17 13:54 GMT-0000

I don't know utf8 stuff well enough (i.e. in how many ways a right arrow can be encoded or how many right arrows exist), but sounds possible as there obviously are >2 variants. ;-)

How can I check my DB setup with regard to utf (4byte capable or not)?



drsassafras 14 Jan 17 03:09 GMT-0000

I opened the doc you posted, but the characters were the same. Could you copy them directly into <http://unicode.scarfboy.com/>, then copy the "URL-encoded UTF8" line and paste them here. That will show exactly what the character is before it goes through any encoding, somewhere along the way.

The code I get with the arrows looks like:URL-encoded UTF8 %E2%86%92

Tiki is not currently able to save 4-bytestring characters. (hence the bug report) It requires updating the database requirements that tiki currently asks for, and also some time applying the changes in tiki. If its done quickly, it could impact performance for every tiki site... The biggest issue however that someone needs to step up and commit to implementing the changes... as is always the case. But lets figure out if his is related first.



RadoS 17 Jan 17 14:43 GMT-0000

Oh, I misunderstood it then that the DB is misconfigured and that alone were relevant, not tiki.

I'll come back when I found a reproducible case.



RadoS 09 Feb 17 16:38 GMT-0000

I can't find it anymore, the user forgot the source.
I assume it was some leftover from migrating from iso to utf (db remained iso, WebServer switched, copied "raw" char from iso made no sense in utf8, pasted into tiki made it fail because invalid utf8).
Should I hit another such case, I'll re-open this report.



RadoS 10 Feb 17 15:29 GMT-0000

I updated the description with my new findings.
a) I cannot reproduce this "just like that" in show, the same string & actions don't lead to the same results.
b) I have no clue what to inspect/ look for at DB, PHP or any other level, pointer & wild guesses welcome.

Attachments

	filename	created	hits	comment	version	filetype
	kkk.pdf	27 Dec 16 15:28 GMT-0000	0	different "right arrow" entities in M\$ word		

The original document is available at

<http://dev.tiki.org/item6219-copy-paste-source-text-from-media-wiki-into-tiki-wiki-this-char-breaks-index-rebuilding>

