

**parseIDS.html:
JavaScript utility for IDS.txt**

suzuki toshiya
Hiroshima University

What is “parseIDS.html”?

- A short JavaScript utility to search Hanzis from ids.txt (a collection of IDS for CJK Unified Ideographs)

Why “grep” is insufficient?

- ids.txt is a collection of the most composed expressions
 - something like NFC of Unicode
 - characters including “林” cannot be found by searching “木”
- recursive searching is required
 - The maintainer (Kawabata-san) uses Emacs + LISP
 - They are slightly exotic software on some platforms

NFC-like IDS vs NFD-like IDS

- When one searches something, an IDS of un-coded Hanzi could be wanted
 - NFC-like IDS cannot include un-coded character
 - except of some CDP compatible glyphs
- Searching with NFD-like IDS would be generic and portable

parseIDS.html uses NFC-like IDS

- For some Hanzis, ids.txt provides multiple expressions (e.g. xxx[GT] yyy[J] zzz[K])
- Making NFD-like IDS for all Hanzis will generate huge collection of possible expressions
 - numVariants(component1) x numVariants(component2) x
- parseIDS.html is expected to be small utility

Too Many Variation Example

- ids.txt defines 6 variants for 𪗇 (U+4E87)

𪗇 𪗇 𪗇 | [GK] 𪗇 𪗇 𪗇 | [T] 𪗇 𪗇 𪗇 𪗇
𪗇 𪗇 𪗇 𪗇 𪗇 𪗇 𪗇 | 𪗇 𪗇 𪗇 |

- 竹 (U+7AF9) will have $6 \times 6 = 36$ variants

- 𪗇 𪗇 (U+25D12) will have...

$36^3 = 46,656$ variants

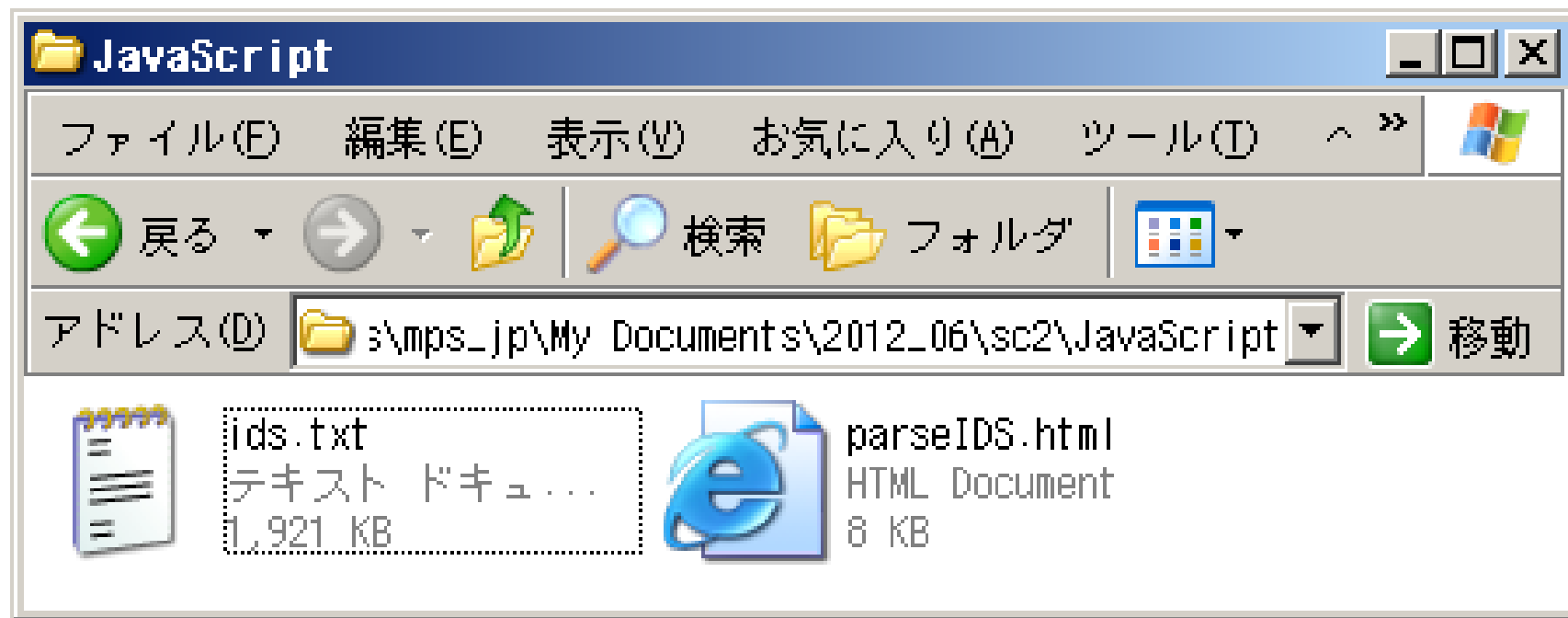
- 𪗇 𪗇 𪗇 𪗇 (U+25DF9) will have...

$36^4 = 1,679,161$ variants

(more than the number of UCS Hanzi!)

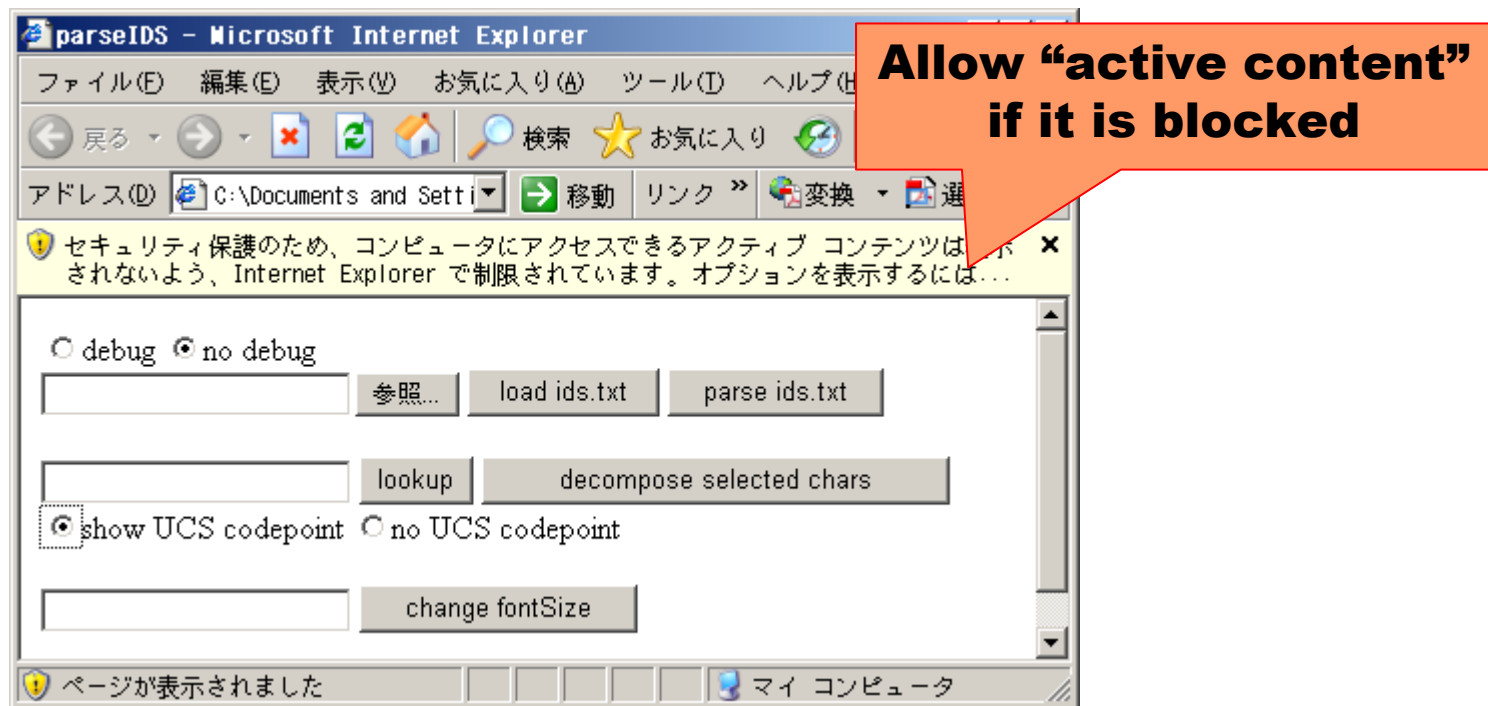
How to use? (installation)

- Place “ids.txt” and “parseIDS.html” to same folder



How to use? (startup)

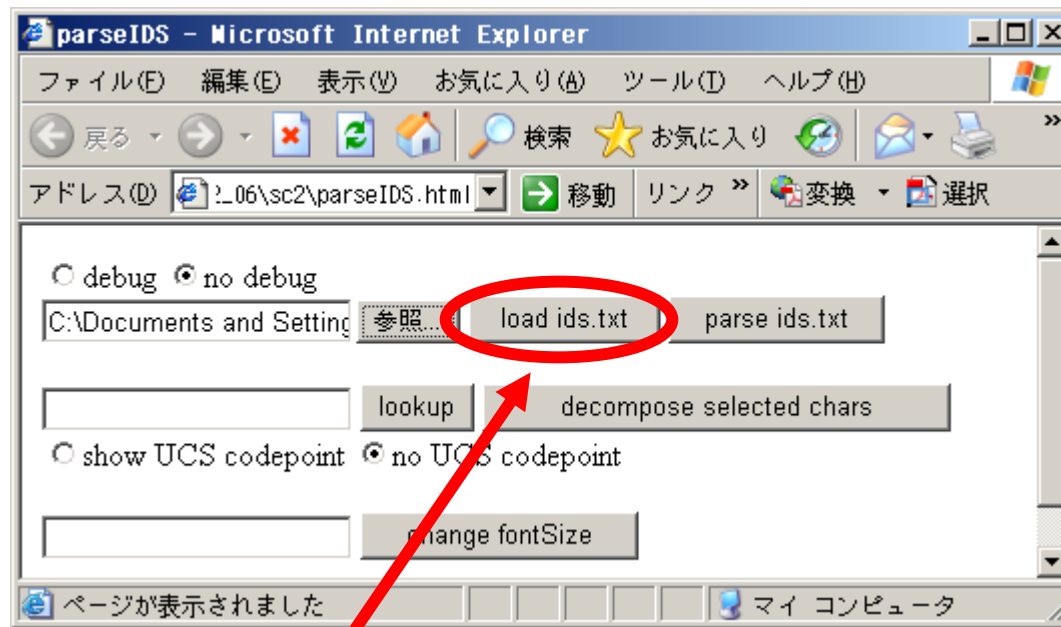
- Open “parseIDS.html” with your web browser



Because ActiveXObject(“Microsoft.XMLHTTP”) is invoked when executed by MSIE, the security warning will be issued.

How to use? (initialization(1))

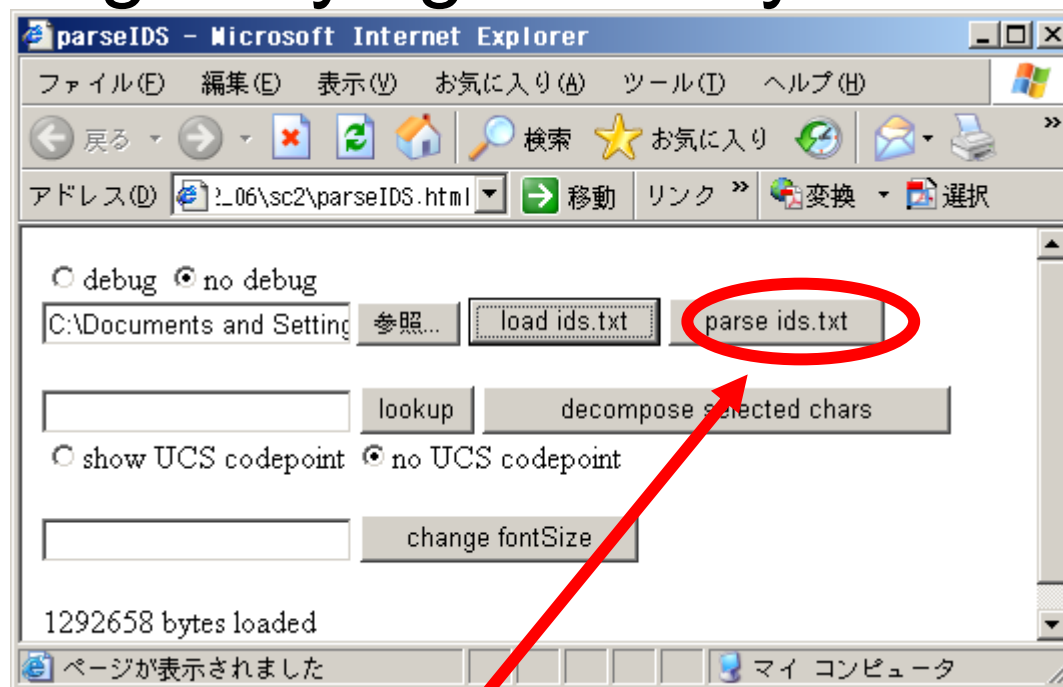
- Select “ids.txt” from the folder where parseIDS.html is placed.
 - The web browsers with HTML5 FileReader API (e.g. Firefox) can handle “ids.txt” in different folders.



- Push “load” button.

How to use? (initialization(2))

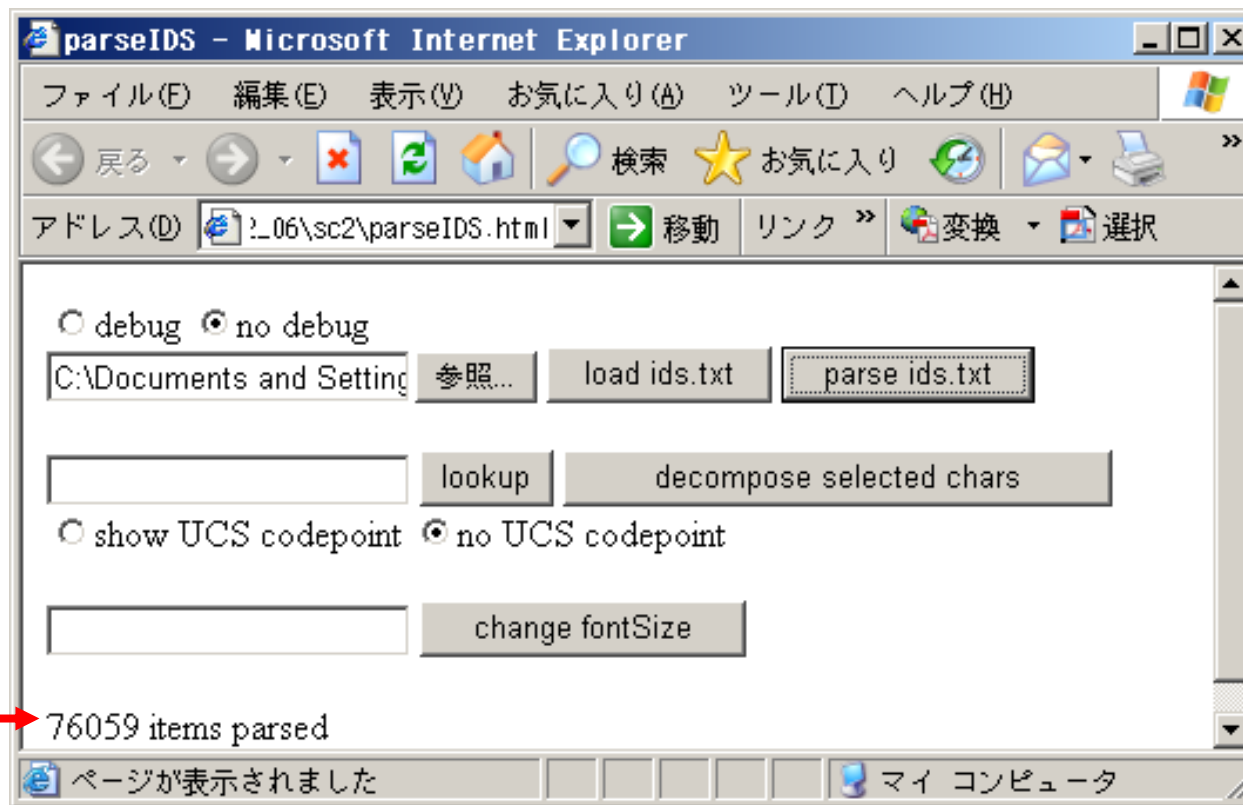
- If loading finishes successfully, you will get a message saying “xxxx bytes loaded”.



- Push “parse” button.

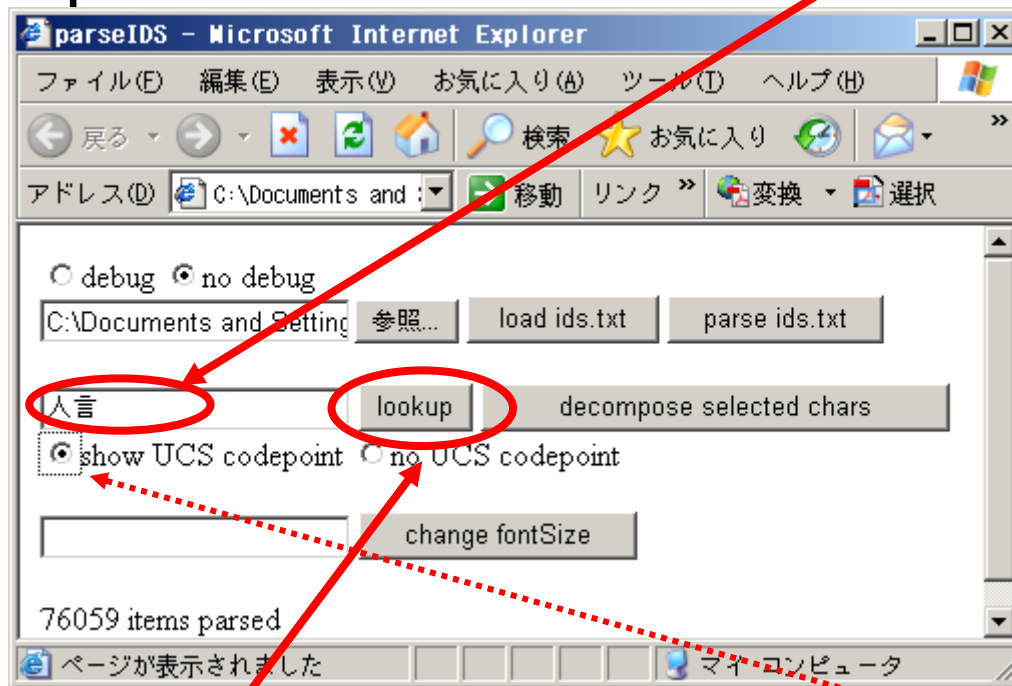
How to use? (initialization(3))

- If parsing finishes successfully, you will get a message saying “xxxx items parsed”.



How to use? (search)

- Enter some Hanzis to the form at the left of “lookup” button.



- Push “lookup” button.
 - If you want UCS codepoint too, check “show UCS”

Search Result

parseIDS - Microsoft Internet Explorer

ファイル(F) 編集(E) 表示(V) お気に入り(A) ツール(T) ヘルプ(H)

戻る 検索 お気に入り

アドレス(D) C:\Documents and Settings\mps_jp\My Documents\2012_06\sc2\pa 移動 リンク 変換 選択

debug no debug

C:\Documents and Setting 参照... load ids.txt parse ids.txt

人言 lookup decompose selected chars

show UCS codepoint no UCS codepoint

change fontSize

變(U+5911), 變(U+5971), 變(U+71EE), 燮(U+7215), 諷(U+85F9), 詆(U+8A1E), 詆(U+8A23), 詆(U+8A44), 詆(U+8A64), 詆(U+8A72), 詆(U+8A7C), 詆(U+8A87), 詆(U+8A8A), 詆(U+8A92), 詆(U+8AA3), 詆(U+8AA4), 詆(U+8AB6), 詆(U+8ABA), 詆(U+8AC7), 詆(U+8ADB), 詆(U+8B01), 詆(U+8B04), 詆(U+8B0D), 詆(U+8B11), 詆(U+8B28), 詆(U+8B29), 詆(U+8B47), 詆(U+8B4F), 詆(U+8B51), 詆(U+8B5B), 詆(U+8B63), 詆(U+8B6A), 詆(U+8B7A), 詆(U+8B82), 詆(U+8B83), 詆(U+8B96), 詆(U+8E9E), 詆(U+9744), 詆(U+39AA), 詆(U+3C14), 詆(U+3F4A), 詆(U+455B), 詆(U+46B6), 詆(U+46C8), 詆(U+46DF), 詆(U+46EE), 詆(U+46F3), 詆(U+46F4),

ページが表示されました

The result is hard to understand?

- You will get the list of Hanzis that include cover all given Hanzis as components.
- The result for “田丁” will include “町”, “學”, “畸”, etc.
 - Different from the result of “grep 田丁 ids.txt”

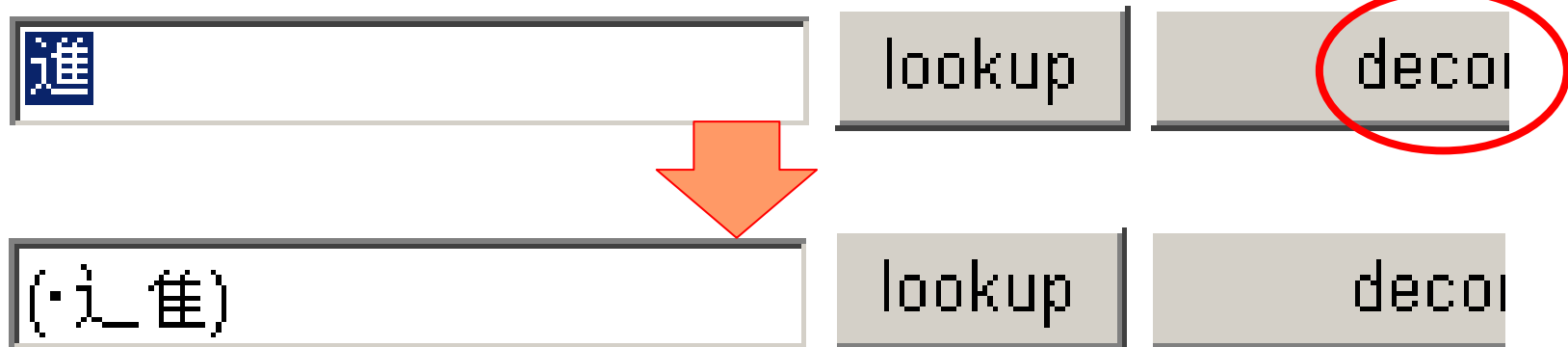
How to use “decomposition”

- Some Hanzi input methods are not easy to input the radicals.

➢ e.g. MS-IME for Japan: “shinnyou” → “之繞” (not “辶”)

- “Decompose” button replaces the selected Hanzis in “lookup” form

➢ After the decomposition, you can remove unrequired components.



Any comments?

feature requests, bugs, etc

→ mpsuzuki@hiroshima-u.ac.jp

appreciations

→ kanji-database-contact@lists.sourceforge.net