

105d67985f ▾

...

gpac / src / utils / utf.c



jeanlf fixed #2106

History

4 contributors



770 lines (703 sloc) | 20.4 KB

...

```
1  /*
2  *          GPAC - Multimedia Framework C SDK
3  *
4  *          Authors: Jean Le Feuvre
5  *          Copyright (c) Telecom ParisTech 2007-2012
6  *          All rights reserved
7  *
8  * This file is part of GPAC / common tools sub-project
9  *
10 * GPAC is free software; you can redistribute it and/or modify
11 * it under the terms of the GNU Lesser General Public License as published by
12 * the Free Software Foundation; either version 2, or (at your option)
13 * any later version.
14 *
15 * GPAC is distributed in the hope that it will be useful,
16 * but WITHOUT ANY WARRANTY; without even the implied warranty of
17 * MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
18 * GNU Lesser General Public License for more details.
19 *
20 * You should have received a copy of the GNU Lesser General Public
21 * License along with this library; see the file COPYING. If not, write to
22 * the Free Software Foundation, 675 Mass Ave, Cambridge, MA 02139, USA.
23 *
24 */
25
26 #ifndef GPAC_DISABLE_CORE_TOOLS
27
28 #include <gpac/utf.h>
29
```

```

30
31 #if 1
32
33
34 /*
35  * Copyright 2001-2004 Unicode, Inc.
36  *
37  * Disclaimer
38  *
39  * This source code is provided as is by Unicode, Inc. No claims are
40  * made as to fitness for any particular purpose. No warranties of any
41  * kind are expressed or implied. The recipient agrees to determine
42  * applicability of information provided. If this file has been
43  * purchased on magnetic or optical media from Unicode, Inc., the
44  * sole remedy for any claim will be exchange of defective media
45  * within 90 days of receipt.
46  *
47  * Limitations on Rights to Redistribute This Code
48  *
49  * Unicode, Inc. hereby grants the right to freely use the information
50  * supplied in this file in the creation of products supporting the
51  * Unicode Standard, and to make copies of this file in any form
52  * for internal or external distribution as long as this notice
53  * remains attached.
54  */
55
56 /* -----
57
58     Conversions between UTF32, UTF-16, and UTF-8. Source code file.
59     Author: Mark E. Davis, 1994.
60     Rev History: Rick McGowan, fixes & updates May 2001.
61     Sept 2001: fixed const & error conditions per
62         mods suggested by S. Parent & A. Lillich.
63     June 2002: Tim Dodd added detection and handling of incomplete
64         source sequences, enhanced error detection, added casts
65         to eliminate compiler warnings.
66     July 2003: slight mods to back out aggressive FFFE detection.
67     Jan 2004: updated switches in from-UTF8 conversions.
68     Oct 2004: updated to use UNI_MAX_LEGAL_UTF32 in UTF-32 conversions.
69
70     See the header file "ConvertUTF.h" for complete documentation.
71
72 ----- */
73
74 typedef u32 UTF32;      /* at least 32 bits */
75 typedef u16 UTF16;      /* at least 16 bits */
76 typedef u8 UTF8;        /* typically 8 bits */
77 typedef u8 Boolean; /* 0 or 1 */
78

```

```

79  /* Some fundamental constants */
80  #define UNI_REPLACEMENT_CHAR (UTF32)0x0000FFFD
81  #define UNI_MAX_BMP (UTF32)0x0000FFFF
82  #define UNI_MAX_UTF16 (UTF32)0x0010FFFF
83  #define UNI_MAX_UTF32 (UTF32)0x7FFFFFFF
84  #define UNI_MAX_LEGAL_UTF32 (UTF32)0x0010FFFF
85
86  typedef enum {
87      conversionOK,          /* conversion successful */
88      sourceExhausted,      /* partial character in source, but hit end */
89      targetExhausted,      /* insuff. room in target for conversion */
90      sourceIllegal         /* source sequence is illegal/malformed */
91  } ConversionResult;
92
93  typedef enum {
94      strictConversion = 0,
95      lenientConversion
96  } ConversionFlags;
97
98  static const int halfShift = 10; /* used for shifting by 10 bits */
99
100 static const UTF32 halfBase = 0x0010000UL;
101 static const UTF32 halfMask = 0x3FFUL;
102
103 #define UNI_SUR_HIGH_START (UTF32)0xD800
104 #define UNI_SUR_HIGH_END (UTF32)0xDBFF
105 #define UNI_SUR_LOW_START (UTF32)0xDC00
106 #define UNI_SUR_LOW_END (UTF32)0xDFFF
107 #define false 0
108 #define true 1
109
110 /*
111  * Index into the table below with the first byte of a UTF-8 sequence to
112  * get the number of trailing bytes that are supposed to follow it.
113  * Note that *legal* UTF-8 values can't have 4 or 5-bytes. The table is
114  * left as-is for anyone who may want to do such conversion, which was
115  * allowed in earlier algorithms.
116  */
117 static const char trailingBytesForUTF8[256] = {
118     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
119     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
120     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
121     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
122     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
123     0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0, 0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
124     1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,1, 1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,
125     2,2,2,2,2,2,2,2,2,2,2,2,2,2,2,2, 3,3,3,3,3,3,3,3,4,4,4,4,5,5,5,5
126 };
127

```

```

128  /*
129  * Magic values subtracted from a buffer value during UTF8 conversion.
130  * This table contains as many values as there might be trailing bytes
131  * in a UTF-8 sequence.
132  */
133  static const UTF32 offsetsFromUTF8[6] = { 0x00000000UL, 0x00003080UL, 0x000E2080UL,
134                                             0x03C82080UL, 0xFA082080UL, 0x82082080UL
135                                             };
136
137  /*
138  * Once the bits are split out into bytes of UTF-8, this is a mask OR-ed
139  * into the first byte, depending on how many bytes follow. There are
140  * as many entries in this table as there are UTF-8 sequence types.
141  * (I.e., one byte sequence, two byte... etc.). Remember that sequences
142  * for *legal* UTF-8 will be 4 or fewer bytes total.
143  */
144  static const UTF8 firstByteMark[7] = { 0x00, 0x00, 0xC0, 0xE0, 0xF0, 0xF8, 0xFC };
145
146  /* ----- */
147
148  /* The interface converts a whole buffer to avoid function-call overhead.
149  * Constants have been gathered. Loops & conditionals have been removed as
150  * much as possible for efficiency, in favor of drop-through switches.
151  * (See "Note A" at the bottom of the file for equivalent code.)
152  * If your compiler supports it, the "isLegalUTF8" call can be turned
153  * into an inline function.
154  */
155
156  /* ----- */
157
158  ConversionResult ConvertUTF16toUTF8 (
159      const UTF16** sourceStart, const UTF16* sourceEnd,
160      UTF8** targetStart, UTF8* targetEnd, ConversionFlags flags) {
161      ConversionResult result = conversionOK;
162      const UTF16* source = *sourceStart;
163      UTF8* target = *targetStart;
164      while (source < sourceEnd) {
165          UTF32 ch;
166          unsigned short bytesToWrite = 0;
167          const UTF32 byteMask = 0xBF;
168          const UTF32 byteMark = 0x80;
169          const UTF16* oldSource = source; /* In case we have to back up because of target overflow */
170          ch = *source++;
171          /* If we have a surrogate pair, convert to UTF32 first. */
172          if (ch >= UNI_SUR_HIGH_START && ch <= UNI_SUR_HIGH_END) {
173              /* If the 16 bits following the high surrogate are in the source buffer... */
174              if (source < sourceEnd) {
175                  UTF32 ch2 = *source;
176                  /* If it's a low surrogate, convert to UTF32. */

```

```

177         if (ch2 >= UNI_SUR_LOW_START && ch2 <= UNI_SUR_LOW_END) {
178             ch = ((ch - UNI_SUR_HIGH_START) << halfShift)
179                 + (ch2 - UNI_SUR_LOW_START) + halfBase;
180             ++source;
181         } else if (flags == strictConversion) { /* it's an unpaired high s
182             --source; /* return to the illegal value itself */
183             result = sourceIllegal;
184             break;
185         }
186     } else { /* We don't have the 16 bits following the high surrogate. */
187         --source; /* return to the high surrogate */
188         result = sourceExhausted;
189         break;
190     }
191 } else if (flags == strictConversion) {
192     /* UTF-16 surrogate values are illegal in UTF-32 */
193     if (ch >= UNI_SUR_LOW_START && ch <= UNI_SUR_LOW_END) {
194         --source; /* return to the illegal value itself */
195         result = sourceIllegal;
196         break;
197     }
198 }
199 /* Figure out how many bytes the result will require */
200 if (ch < (UTF32)0x80) {
201     bytesToWrite = 1;
202 } else if (ch < (UTF32)0x800) {
203     bytesToWrite = 2;
204 } else if (ch < (UTF32)0x10000) {
205     bytesToWrite = 3;
206 } else if (ch < (UTF32)0x110000) {
207     bytesToWrite = 4;
208 } else {
209     bytesToWrite = 3;
210     ch = UNI_REPLACEMENT_CHAR;
211 }
212
213 target += bytesToWrite;
214 if (target > targetEnd) {
215     source = oldSource; /* Back up source pointer! */
216     target -= bytesToWrite;
217     result = targetExhausted;
218     break;
219 }
220 switch (bytesToWrite) { /* note: everything falls through. */
221 case 4:
222     *--target = (UTF8)((ch | byteMark) & byteMask);
223     ch >>= 6;
224 case 3:
225     *--target = (UTF8)((ch | byteMark) & byteMask);

```

```

226         ch >>= 6;
227     case 2:
228         *--target = (UTF8)((ch | byteMark) & byteMask);
229         ch >>= 6;
230     case 1:
231         *--target = (UTF8)(ch | firstByteMark[bytesToWrite]);
232     }
233     target += bytesToWrite;
234 }
235 *sourceStart = source;
236 *targetStart = target;
237 return result;
238 }
239
240 /*
241  * Utility routine to tell whether a sequence of bytes is legal UTF-8.
242  * This must be called with the length pre-determined by the first byte.
243  * If not calling this from ConvertUTF8to*, then the length can be set by:
244  * length = trailingBytesForUTF8[*source]+1;
245  * and the sequence is illegal right away if there aren't that many bytes
246  * available.
247  * If presented with a length > 4, this returns false. The Unicode
248  * definition of UTF-8 goes up to 4-byte sequences.
249  */
250
251 Boolean isLegalUTF8(const UTF8 *source, int length) {
252     UTF8 a;
253     const UTF8 *srcptr = source+length;
254     switch (length) {
255     default:
256         return false;
257     /* Everything else falls through when "true"... */
258     case 4:
259         if ((a = (*--srcptr)) < 0x80 || a > 0xBF) return false;
260     case 3:
261         if ((a = (*--srcptr)) < 0x80 || a > 0xBF) return false;
262     case 2:
263         if ((a = (*--srcptr)) > 0xBF) return false;
264
265         switch (*source) {
266         /* no fall-through in this inner switch */
267         case 0xE0:
268             if (a < 0xA0) return false;
269             break;
270         case 0xED:
271             if (a > 0x9F) return false;
272             break;
273         case 0xF0:
274             if (a < 0x90) return false;

```

```

275         break;
276     case 0xF4:
277         if (a > 0x8F) return false;
278         break;
279     default:
280         if (a < 0x80) return false;
281     }
282
283     case 1:
284         if (*source >= 0x80 && *source < 0xC2) return false;
285     }
286     if (*source > 0xF4) return false;
287     return true;
288 }
289
290 /* ----- */
291
292 ConversionResult ConvertUTF8toUTF16 (
293     const UTF8** sourceStart, const UTF8* sourceEnd,
294     UTF16** targetStart, UTF16* targetEnd, ConversionFlags flags) {
295     ConversionResult result = conversionOK;
296     const UTF8* source = *sourceStart;
297     UTF16* target = *targetStart;
298     while (source < sourceEnd) {
299         UTF32 ch = 0;
300         unsigned short extraBytesToRead = trailingBytesForUTF8[*source];
301         if (source + extraBytesToRead >= sourceEnd) {
302             result = sourceExhausted;
303             break;
304         }
305         /* Do this check whether lenient or strict */
306         if (!isLegalUTF8(source, extraBytesToRead+1)) {
307             result = sourceIllegal;
308             break;
309         }
310         /*
311          * The cases all fall through. See "Note A" below.
312          */
313         switch (extraBytesToRead) {
314             case 5:
315                 ch += *source++;
316                 ch <<= 6; /* remember, illegal UTF-8 */
317             case 4:
318                 ch += *source++;
319                 ch <<= 6; /* remember, illegal UTF-8 */
320             case 3:
321                 ch += *source++;
322                 ch <<= 6;
323             case 2:

```

```

324         ch += *source++;
325         ch <= 6;
326     case 1:
327         ch += *source++;
328         ch <= 6;
329     case 0:
330         ch += *source++;
331     }
332     ch -= offsetsFromUTF8[extraBytesToRead];
333
334     if (target >= targetEnd) {
335         source -= (extraBytesToRead+1); /* Back up source pointer! */
336         result = targetExhausted;
337         break;
338     }
339     if (ch <= UNI_MAX_BMP) { /* Target is a character <= 0xFFFF */
340         /* UTF-16 surrogate values are illegal in UTF-32 */
341         if (ch >= UNI_SUR_HIGH_START && ch <= UNI_SUR_LOW_END) {
342             if (flags == strictConversion) {
343                 source -= (extraBytesToRead+1); /* return to the illegal v
344                 result = sourceIllegal;
345                 break;
346             } else {
347                 *target++ = UNI_REPLACEMENT_CHAR;
348             }
349         } else {
350             *target++ = (UTF16)ch; /* normal case */
351         }
352     } else if (ch > UNI_MAX_UTF16) {
353         if (flags == strictConversion) {
354             result = sourceIllegal;
355             source -= (extraBytesToRead+1); /* return to the start */
356             break; /* Bail out; shouldn't continue */
357         } else {
358             *target++ = UNI_REPLACEMENT_CHAR;
359         }
360     } else {
361         /* target is a character in range 0xFFFF - 0x10FFFF. */
362         if (target + 1 >= targetEnd) {
363             source -= (extraBytesToRead+1); /* Back up source pointer! */
364             result = targetExhausted;
365             break;
366         }
367         ch -= halfBase;
368         *target++ = (UTF16)((ch >> halfShift) + UNI_SUR_HIGH_START);
369         *target++ = (UTF16)((ch & halfMask) + UNI_SUR_LOW_START);
370     }
371 }
372 *sourceStart = source;

```



```

373     *targetStart = target;
374     return result;
375 }
376
377
378 GF_EXPORT
379 Bool gf_utf8_is_legal(const u8 *data, u32 length)
380 {
381     //we simply run ConvertUTF8toUTF16 without target
382     const UTF8** sourceStart = (const UTF8**) &data;
383     const UTF8* sourceEnd = (const UTF8*) ( data + length );
384     ConversionResult result = conversionOK;
385     const UTF8* source = *sourceStart;
386
387     while (source < sourceEnd) {
388         UTF32 ch = 0;
389         unsigned short extraBytesToRead = trailingBytesForUTF8[*source];
390         if (source + extraBytesToRead >= sourceEnd) {
391             result = sourceExhausted;
392             break;
393         }
394         /* Do this check whether lenient or strict */
395         if (! isLegalUTF8(source, extraBytesToRead+1)) {
396             result = sourceIllegal;
397             break;
398         }
399         /*
400          * The cases all fall through. See "Note A" below.
401          */
402         switch (extraBytesToRead) {
403             case 5:
404                 ch += *source++;
405                 ch <<= 6; /* remember, illegal UTF-8 */
406             case 4:
407                 ch += *source++;
408                 ch <<= 6; /* remember, illegal UTF-8 */
409             case 3:
410                 ch += *source++;
411                 ch <<= 6;
412             case 2:
413                 ch += *source++;
414                 ch <<= 6;
415             case 1:
416                 ch += *source++;
417                 ch <<= 6;
418             case 0:
419                 ch += *source++;
420         }
421         ch -= offsetsFromUTF8[extraBytesToRead];

```

```

422
423         if (ch <= UNI_MAX_BMP) { /* Target is a character <= 0xFFFF */
424             /* UTF-16 surrogate values are illegal in UTF-32 */
425             if (ch >= UNI_SUR_HIGH_START && ch <= UNI_SUR_LOW_END) {
426                 result = sourceIllegal;
427                 break;
428             }
429         } else if (ch > UNI_MAX_UTF16) {
430             result = sourceIllegal;
431             break; /* Bail out; shouldn't continue */
432         }
433     }
434     return (result==conversionOK) ? GF_TRUE : GF_FALSE;
435 }
436
437 GF_EXPORT
438 u32 gf_utf8_wcslen (const unsigned short *s)
439 {
440     const unsigned short* ptr;
441     if (!s) return 0;
442     for (ptr = s; *ptr != (unsigned short)'\0'; ptr++) {
443     }
444     return (u32) ( ptr - s );
445 }
446
447 GF_EXPORT
448 u32 gf_utf8_wcstombs(char* dest, size_t len, const unsigned short** srcp)
449 {
450     if (!srcp || !*srcp)
451         return 0;
452     else {
453         const UTF16** sourceStart = srcp;
454         const UTF16* sourceEnd = *srcp + gf_utf8_wcslen(*srcp);
455         UTF8* targetStart = (UTF8*) dest;
456         UTF8* targetEnd = (UTF8*) dest + len;
457         ConversionFlags flags = strictConversion;
458
459         ConversionResult res = ConvertUTF16toUTF8(sourceStart, sourceEnd, &targetStart, ta
460         if (res != conversionOK) return GF_UTF8_FAIL;
461         *targetStart = 0;
462         *srcp=NULL;
463         return (u32) strlen(dest);
464     }
465 }
466
467 GF_EXPORT
468 u32 gf_utf8_mbstowcs(unsigned short* dest, size_t len, const char** srcp)
469 {
470     if (!srcp || !*srcp)

```

```

471         return 0;
472     else {
473         const UTF8** sourceStart = (const UTF8**) srcp;
474         const UTF8* sourceEnd = (const UTF8*) ( *srcp + strlen( *srcp) );
475         UTF16* targetStart = (UTF16* ) dest;
476         UTF16* targetEnd = (UTF16* ) (dest + len);
477         ConversionFlags flags = strictConversion;
478         ConversionResult res = ConvertUTF8toUTF16(sourceStart, sourceEnd, &targetStart, ta
479         if (res != conversionOK) return GF_UTF8_FAIL;
480         *targetStart = 0;
481         *srcp=NULL;
482         return gf_utf8_wcslen(dest);
483     }
484 }
485
486
487 #else
488
489 GF_EXPORT
490 u32 gf_utf8_wcslen (const unsigned short *s)
491 {
492     const unsigned short* ptr;
493     for (ptr = s; *ptr != (unsigned short)'\0'; ptr++) {
494     }
495     return (u32) (ptr - s);
496 }
497
498 GF_EXPORT
499 u32 gf_utf8_wcstombs(char* dest, size_t len, const unsigned short** srcp)
500 {
501     /*
502     * Original code from the GNU UTF-8 Library
503     */
504     size_t count;
505     const unsigned short * src = *srcp;
506
507     if (dest != NULL) {
508         char* destptr = dest;
509         for (;;) src++ {
510             unsigned char c;
511             unsigned short wc = *src;
512             if (wc < 0x80) {
513                 if (wc == (wchar_t)'\0') {
514                     if (len == 0) {
515                         *srcp = src;
516                         break;
517                     }
518                     *destptr = '\0';
519                     *srcp = NULL;

```

```

520             break;
521         }
522         count = 0;
523         c = (unsigned char) wc;
524     } else if (wc < 0x800) {
525         count = 1;
526         c = (unsigned char) ((wc >> 6) | 0xC0);
527     } else {
528         count = 2;
529         c = (unsigned char) ((wc >> 12) | 0xE0);
530     }
531     if (len <= count) {
532         *srcp = src;
533         break;
534     }
535     len -= count+1;
536     *destptr++ = c;
537     if (count > 0)
538         do {
539             *destptr++ = (unsigned char)((wc >> (6 * --count)) & 0x3F
540             } while (count > 0);
541     }
542     return (u32) (destptr - dest);
543 } else {
544     /* Ignore dest and len. */
545     size_t totalcount = 0;
546     for (;;) src++ {
547         unsigned short wc = *src;
548         size_t count;
549         if (wc < 0x80) {
550             if (wc == (wchar_t)'\0') {
551                 *srcp = NULL;
552                 break;
553             }
554             count = 1;
555         } else if (wc < 0x800) {
556             count = 2;
557         } else {
558             count = 3;
559         }
560         totalcount += count;
561     }
562     return (u32) totalcount;
563 }
564 }
565
566
567 typedef struct
568 {

```

```

569     u32 count : 16;    /* number of bytes remaining to be processed */
570     u32 value : 16;    /* if count > 0: partial wide character */
571     /*
572         If WCHAR_T_BITS == 16, need 2 bits for count,
573         12 bits for value (10 for mbstowcs direction, 12 for wcstombs direction).
574     */
575 } gf_utf8_mbstate_t;
576
577 static gf_utf8_mbstate_t internal;
578
579 GF_EXPORT
580 u32 gf_utf8_mbstowcs(unsigned short* dest, size_t len, const char** srcp)
581 {
582     gf_utf8_mbstate_t* ps = &internal;
583     const char *src = *srcp;
584
585     unsigned short* destptr = dest;
586     for (; len > 0; destptr++, len--) {
587         const char* backup_src = src;
588         unsigned char c;
589         unsigned short wc;
590         size_t count;
591         if (ps->count == 0) {
592             c = (unsigned char) *src;
593             if (c < 0x80) {
594                 *destptr = (wchar_t) c;
595                 if (c == 0) {
596                     src = NULL;
597                     break;
598                 }
599                 src++;
600                 continue;
601             } else if (c < 0xC0) {
602                 /* Spurious 10XXXXXX byte is invalid. */
603                 goto bad_input;
604             }
605             if (c < 0xE0) {
606                 wc = (wchar_t)(c & 0x1F) << 6;
607                 count = 1;
608                 if (c < 0xC2) goto bad_input;
609             } else if (c < 0xF0) {
610                 wc = (wchar_t)(c & 0x0F) << 12;
611                 count = 2;
612             }
613             else goto bad_input;
614             src++;
615         } else {
616             wc = ps->value << 6;
617             count = ps->count;

```

```

618     }
619     for (;;) {
620         c = (unsigned char) *src++ ^ 0x80;
621         if (!(c < 0x40)) goto bad_input_backup;
622         wc |= (unsigned short) c << (6 * --count);
623         if (count == 0)
624             break;
625         /* The following test is only necessary once for every character,
626            but it would be too complicated to perform it once only, on
627            the first pass through this loop. */
628         if ((unsigned short) wc < ((unsigned short) 1 << (5 * count + 6)))
629             goto bad_input_backup;
630     }
631     *destptr = wc;
632     ps->count = 0;
633     continue;
634
635 bad_input_backup:
636     src = backup_src;
637     goto bad_input;
638 }
639 *srcp = src;
640 return (u32) (destptr - dest);
641
642 bad_input:
643     *srcp = src;
644     return GF_UTF8_FAIL;
645 }
646
647
648 #endif
649
650
651 GF_EXPORT
652 GF_Err gf_utf_get_utf8_string_from_bom(const u8 *data, u32 size, char **out_ptr, char **result)
653 {
654     u32 unicode_type = 0;
655     if (!out_ptr || !result || !data) return GF_BAD_PARAM;
656     *out_ptr = NULL;
657     *result = (char *) data;
658
659     if (size >= 5) {
660         /*0: no unicode, 1: UTF-16BE, 2: UTF-16LE*/
661         if ((data[0] == 0xFF) && (data[1] == 0xFE)) {
662             if (!data[2] && !data[3]) {
663                 return GF_OK;
664             } else {
665                 unicode_type = 2;
666             }

```

```

667         } else if ((data[0]==0xFE) && (data[1]==0xFF)) {
668             if (!data[2] && !data[3]) {
669                 return GF_OK;
670             } else {
671                 unicode_type = 1;
672             }
673         } else if ((data[0]==0xEF) && (data[1]==0xBB) && (data[2]==0xBF)) {
674             *result = (char *) (data+4);
675             return GF_OK;
676         }
677     }
678
679     if (!unicode_type) {
680         *result = (char *) data;
681         return GF_OK;
682     }
683
684     if (size%2) size--;
685     u16 *str_wc = gf_malloc(size+2);
686     if (!str_wc) return GF_OUT_OF_MEM;
687     u16 *srcwc;
688     char *dst = gf_malloc(size+2);
689     if (!dst) {
690         gf_free(str_wc);
691         return GF_OUT_OF_MEM;
692     }
693     *out_ptr = dst;
694     u32 i;
695     for (i=0; i<size; i+=2) {
696         u16 wchar=0;
697         u8 c1 = data[i];
698         u8 c2 = data[i+1];
699
700         /*Little-endian order*/
701         if (unicode_type==2) {
702             if (c2) {
703                 wchar = c2;
704                 wchar <<=8;
705                 wchar |= c1;
706             }
707             else wchar = c1;
708         } else {
709             wchar = c1;
710             if (c2) {
711                 wchar <<= 8;
712                 wchar |= c2;
713             }
714         }
715         str_wc[i/2] = wchar;

```

```

716     }
717     str_wc[i/2] = 0;
718     srcwc = str_wc;
719     u32 res = gf_utf8_wcstombs(dst, size, (const unsigned short **) &srcwc);
720     gf_free(str_wc);
721     if (res==GF_UTF8_FAIL) {
722         gf_free(dst);
723         *out_ptr = NULL;
724         return GF_IO_ERR;
725     }
726     *result = dst;
727     return GF_OK;
728 }
729
730
731 #if defined(WIN32)
732
733 GF_EXPORT
734 wchar_t* gf_utf8_to_wcs(const char* str)
735 {
736     size_t source_len;
737     wchar_t* result;
738     if (str == 0) return 0;
739     source_len = strlen(str);
740     result = gf_calloc(source_len + 1, sizeof(wchar_t));
741     if (!result)
742         return 0;
743     if (gf_utf8_mbstowcs(result, source_len, &str) == GF_UTF8_FAIL) {
744         gf_free(result);
745         return 0;
746     }
747     return result;
748 }
749
750 GF_EXPORT
751 char* gf_wcs_to_utf8(const wchar_t* str)
752 {
753     size_t source_len;
754     char* result;
755     if (str == 0) return 0;
756     source_len = wcslen(str);
757     result = gf_calloc(source_len + 1, UTF8_MAX_BYTES_PER_CHAR);
758     if (!result)
759         return 0;
760     if (gf_utf8_wcstombs(result, source_len * UTF8_MAX_BYTES_PER_CHAR, &str) == GF_UTF8_FAIL)
761         gf_free(result);
762     return 0;
763 }
764 return result;

```



```
765     }  
766     #endif  
767  
768     #endif /* GPAC_DISABLE_CORE_TOOLS */  
769  
770
```

