

Accelerating UTF-8 Decoding Using SIMD Instructions

<u>Hiroshi Inoue</u> IBM Research – Tokyo Mar 17, 2008

© 2008 IBM Corporation



Background and Goal

- Background
 - UTF-8 encoding is commonly used to exchange text data (e.g. HTML, XML)
 - But, Java VM uses UTF-16 as its internal representation



- Conversion between UTF-8 and UTF-16 matters

→Goal

- Accelerate conversion from UTF-8 to UTF-16



UTF-8 and UTF-16

- What is UTF-8 and UTF-16
 - UTF-8: variable length encoding (from 1 byte to 3 bytes per character)
 - UTF-16: constant length encoding (2 bytes, other than surrogate pair)
- Mapping UTF-8 onto UTF-16

range	UTF-8	UTF-16
00-7F (ASCII char)	0xxxxxx	000000000xxxxxxx
80-7FF (Greece char etc)	110yyyyy 10xxxxx	00000yyyyyxxxxx
800-FFFF (Chinese, Japanese char etc)	1110zzzz 10yyyyyy 10 <mark>xxxxx</mark>	zzzzyyyyyyxxxxx



Decoding UTF-8 to UTF-16

Naive implementation of decoding process based on *conditional branches*



Main source of overhead

 Branch mispredictions caused by data-dependent conditional branches that are hard to predict (when having multiple types of characters)



simple encoding format for explanation

We use a simplified format "**simple encoding**" for ease of explanation

simple encoding

- variable length encoding for integers up to 30-bits long
- I value is encoding to 1 to 4 bytes
- Most significant 2 bits of the first byte show the length

length	representation in simple encoding	
1 byte		00xxxxx
2 bytes	01уууууу	xxxxxxx
3 bytes	10zzzzz yyyyyyy	xxxxxxx
4 bytes	11 aaaaaa zzzzzzz yyyyyyyy	xxxxxxx



Our Technique without conditional branches

Steps for conversion

Step 1: Identify data length for next few data (by scalar)

- Step 2: Load required constants based on the data length identified in Step 1 (by scalar)
- Step 3: Move data using permute instruction (by SIMD)
- Step 4: Mask off unused bits (by SIMD)

SIMD Permute (Shuffle) instruction

reorder input bytes based on a pattern specified by register at runtime (not compile time!)

many SIMD ISAs have similar instructions (e.g. VMX, SSE)













Apply our technique for UTF-8 to UTF-16 conversion

- Not that different from the case for simple encoding...
- Some difference:
 - use <u>two permute</u> and <u>two select</u> instructions instead of just one permute instruction
 - use four constant vector values per iteration
 - convert 8 characters (2 bytes each) in each iteration



UTF-8 to UTF-16 Conversions





Additional optimization for real-world text

- Same type (i.e. same length in UTF-8 representation) of characters tends to appear repeatedly in real-world data
 - If all characters converted in an iteration have same type, we use specialized code for this data type
 - Specialized code for ASCII characters are especially efficient because it simply inserts 0 before each byte (0xxxxxxx → 0000000 0xxxxxxx)
 - Specialized code fall backs to the default code if input values include a value of other data type

Applied for both vectorized version and serial version



Evaluation

- Tested on PowerPC 970MP 2.5 GHz running Linux 2.6.18
- Implemented using VMX intrinsics
- Compiled with gcc-4.0.1



Decode performance for simple encoding format



Scalar processing caused frequent branch mispredictions for random input
Performance of our technique does not depend on input data



Performance of UTF-8 decoding with artificial text





Performance of UTF-8 decoding with realistic text





Future work

To integrated into Java environment

- Conversion is implemented in Java, where programmer cannot directly write platform-dependent SIMD code
- JIT compiler must generate SIMD instructions for each platform

UTF-16 to UTF-8 conversion

- As important as UTF-8 to UTF-16 conversion
- Our technique is applicable. However, writing UTF-8 sequence to memory require costly unaligned memory store



Summary

- We developed decoding technique for variable-length encoding format without using hard-to-predict conditional branches
- We observed significant performance boost in UTF-8 to UTF-16 conversion with real-world text data

backup



Comparison with Cameron [PPoPP 2008]

