



Confusion In Conversion

How Does Occasional Character Corruption Happen ?

Ken Glidden, Teruhiko Kurosaka, & Kazuo Hirai

25th Internationalization and Unicode Conference,
Washington, DC, March/April 2004

strategy h process h technology h results

www.basistech.com

Introductions

- Basis Technology – www.basistech.com
- Ken – manager of engineering services
- Kuro – principal software engineer
- Kazuo – principal software engineer
- Kazuo and Kuro: native Japanese speakers

Outline

- Characterization
- Business impact
- Detection
- Prevention

Symptom of the Problem

```
Command Prompt
C:\>type CP932.txt
彌 日本 TEL

C:\>jdk1.3\bin\java EchoFile CP932.txt Shift_JIS
彌 日本 TEL

C:\>j2sdk1.4.2_03\bin\java EchoFile CP932.txt Shift_JIS
?g 日本 ??

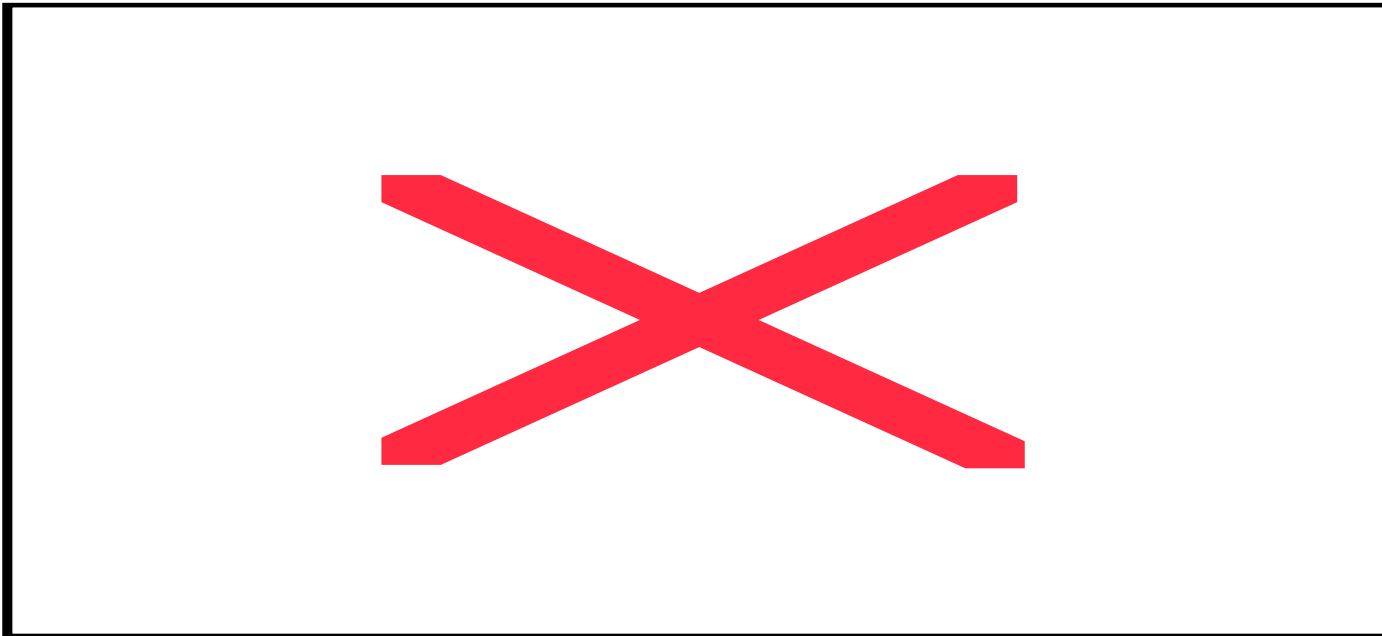
C:\>
```

Problem Domain

- Some definitions
 - Confusion: 2 systems disagree
 - Character: implies readable text
 - Conversion: legacy - > Unicode -> legacy
 - Corruption: code points are wrong
 - Occasional: edge cases
- Domain Summary
 - Multi-tier systems
 - Transcoders disagree for identical encodings



Location of the Problem



Impact of the Problem

- Legacy encodings are still used
- System may need to input/output in a legacy encoding
- Examples:
 - Email-based help desk
 - Document storage/retrieval system
 - Network monitoring programs
 - Mixture of Java and C++ subsystems
- These issues are deal breakers

Cause of the Problem

- Generic issues
 - Conflicting encoding definitions
 - Conflicting maps
 - Vendor extensions
 - Undefined Code Points
 - Duplicates
- Japanese examples
 - 0x5C Yen or backslash
 - Some JIS characters mapped differently
 - IBM/NEC extensions
 - Half-width katakana and ISO-2022
 - CP932

Conflicting Encoding Definitions

- ASCII or JIS-Roman?
 - 0x5C = \ or ¥ (half-width yen)
 - 0x7e = ~ or _

Vendor Platform	Source Encoding	ASCII / JIS-Roman
RedHat Linux	EUC-JP	ASCII
RedHat Linux	SJIS	JIS-Roman
Solaris	EUC-JP	ASCII
Windows	CP932	ASCII
DEC OSF	EUC-JP	JIS-Roman

Conflicting Maps

- Differences in JIS X 0208-1990 to Unicode maps for 7 characters. (nn-nn) = JIS code point
 - \
 ~ (01-33)
 - _ (01-34)
 - (01-61)
 - _ (01-81)
 - _ (01-82)
 - _ (02-44)
- e.g. 01-33
 - U+301C, ~, RedHat Linux
 - U+FF5E, ~, Microsoft
- Chinese problems, too.

Vendor Extensions

- IBM/NEC
- Microsoft's CP932 includes these
- Others do not
- Potential for OCC



Undefined Code Points

- Some transcoders handle this differently
- Example: Halfwidth Katakana (ㇿ vs. ㇾ)
 - Not defined in ISO-2022-JP
 - But, some transcoders convert halfwidth to fullwidth

Duplicates

- Microsoft's CP932 has duplicate and triplicate characters in the IBM/NEC extension range
- The CP932 code map does not allow these to roundtrip.
- Visually this is not a problem
- Bit-wise it is.



Detecting OCC

- White-box testing (QA)
- Black-box testing (development)



White-box Testing

- Casual tests are insufficient
- Carefully planned tests are needed:
 - Loopback tests
 - Use all characters
 - Ensure all transcoders are exercised
- Test all subsystems

White-box Testing Caveats

- User-Defined Characters.
- Duplicated code points in MS 932 extension blocks typically do not round trip; need special treatment in judging test result.



Black-box Testing

- Test each transcoding module
- Characterize each transcoder
 - Test every character
 - Roundtrip every character



Preventing OCC

- Use the same transcoder everywhere
 - ICU
 - Basis Tech's RCLU
 - Gnu's libiconv
- Use transcoder's exception handling features.
- Use roundtrip-enabled transcoders

Conclusion

- OCC is a real issue for Unicode-enabled systems
- Recommendations for application developers
 - Include OCC tests during QA
 - Use one cross-platform library everywhere
 - Characterize the transcoder and handle its anomalies
- Recommendations for vendors
 - Standardize maps (e.g. iconv)

Thanks for coming!

- Thanks for the opportunity to share.