



Distributed Computing

Lesson 7: Text Encoding

Thomas Weise · 汤卫思

tweise@hfu.edu.cn · <http://www.it-weise.de>

Hefei University, South Campus 2
Faculty of Computer Science and Technology
Institute of Applied Optimization
230601 Shushan District, Hefei, Anhui, China
Econ. & Tech. Devel. Zone, Jinxiu Dadao 99

合肥学院 南艳湖校区/南2区
计算机科学与技术系
应用优化研究所
中国 安徽省 合肥市 蜀山区 230601
经济技术开发区 锦绣大道99号

1 Character Encoding



website

- How can we deal with text?
- How can we know that a sequence of bits stands for “汤卫思” and not for “Ölüberschussländer”?

- Different languages have different characters

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]
- Original idea: bytes have different meaning, depending on language

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]
- Original idea: bytes have different meaning, depending on language (for German, we can e.g., replace some less important characters with “ä” and “ß” ...)

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]
- Original idea: bytes have different meaning, depending on language
- GB2312 ^[2] encoding especially for Chinese characters (2B for each non-ASCII char)

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]
- Original idea: bytes have different meaning, depending on language
- GB2312 ^[2] encoding especially for Chinese characters (2B for each non-ASCII char)
- These approaches are insufficient for other languages

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: ASCII / ISO/IEC 8859-1 ^[1]
- Original idea: bytes have different meaning, depending on language
- GB2312 ^[2] encoding especially for Chinese characters (2B for each non-ASCII char)
- These approaches are insufficient for other languages
- Universal Coded Character Set (UCS) ^[3] and Unicode ^[4-6]

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: `ASCII` / `ISO/IEC 8859-1` ^[1]
- Original idea: bytes have different meaning, depending on language
- `GB2312` ^[2] encoding especially for Chinese characters (2B for each non-`ASCII` char)
- These approaches are insufficient for other languages
- **Universal Coded Character Set (UCS)** ^[3] and **Unicode** ^[4-6]
- Encoded as `UTF-7` , `UTF-8` ^[7] (compatible to `ASCII`), `UTF-16` , and `UTF-32`

- Different languages have different characters
- Originally, storage of text data was mainly designed for US English
- Here, 1B per character is sufficient: `ASCII` / `ISO/IEC 8859-1` ^[1]
- Original idea: bytes have different meaning, depending on language
- `GB2312` ^[2] encoding especially for Chinese characters (2B for each non-`ASCII` char)
- These approaches are insufficient for other languages
- `Universal Coded Character Set (UCS)` ^[3] and `Unicode` ^[4-6]
- Encoded as `UTF-7` , `UTF-8` ^[7] (compatible to `ASCII`), `UTF-16` , and `UTF-32`
- When sending text data, we need to make sure to use the right encoding!

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input**

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader`s read one or multiple *unicode characters*

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream
- **Output**

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream
- **Output:**
 - `Writer` s write one or multiple unicode characters

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream
- **Output:**
 - `Writer` s write one or multiple unicode characters
 - `OutputStreamWriter` s are writers which store all characters written to them in a specified encoding into byte-based output streams

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream
- **Output:**
 - `Writer` s write one or multiple unicode characters
 - `OutputStreamWriter` s are writers which store all characters written to them in a specified encoding into byte-based output streams
- TCP sockets: plug the `InputStreamReader` and `OutputStreamWriter` s directly into the streams that the socket offers to us

- In Java, we can use the more general Character Stream API ^[8] to deal with data conversation
- **Input:**
 - `Reader` s read one or multiple *unicode characters*
 - `InputStreamReader` s are readers which take their data from a byte-based input stream
- **Output:**
 - `Writer` s write one or multiple unicode characters
 - `OutputStreamWriter` s are writers which store all characters written to them in a specified encoding into byte-based output streams
- TCP sockets: plug the `InputStreamReader` and `OutputStreamWriter` s directly into the streams that the socket offers to us
- UDP sockets: create the packets in memory

- Usually determined at compile-time

- Usually determined at compile-time
- Different types for different characters and encodings: `char` , `TCHAR` , `wchar_t` , ...

Listing: MinHTTPClient.java HTTP Client: Java

```
import java.io.BufferedReader;      import java.io.InputStreamReader;
import java.io.OutputStreamWriter; import java.net.Socket;

public class MinHTTPClient {//this is a minimum web client; see lesson 07 coming later

    public static final void main(final String[] args) {
        String dest, request, response;      Socket sock;
        OutputStreamWriter w;                BufferedReader r;

        dest    = "www.baidu.com";           // a random example for a Chinese host
        request = "GET_/index.html_HTTP/1.1\nHost:_ " + dest + "\n\n";

        try {
            sock = new Socket(dest, 80);      // web servers are usually listening at port 80
            w     = new OutputStreamWriter(sock.getOutputStream());
            w.write(request);                  // write the HTTP request [9-11]
            w.flush();                         // make sure that all data has been sent
            sock.shutdownOutput();            // closing down the channel for sending data to the server

            r = new BufferedReader(new InputStreamReader(sock.getInputStream(), "UTF-8")); // Baidu uses UTF-8 encoding
                                                    // before they used GB2312 [9] encoding

            while ((response = r.readLine()) != null) { // read strings line-by-line until connection closed by server
                System.out.println(response);          // print to output
            }

            sock.close();
        } catch (Throwable t) {
            t.printStackTrace();
        }
    }
}
```

Listing: MinHTTPClientJava17.java Min HTTP Client + Try-With-Resource

```
import java.io.BufferedReader;           import java.io.InputStreamReader;
import java.io.OutputStreamWriter;       import java.net.Socket;

public class MinHTTPClientJava17 {//this is a minimum web client; see lesson 07 coming later
    public static final void main(final String[] args) {
        String dest, request, response;

        dest    = "www.baidu.com";       // a random example for a Chinese host
        request = "GET_/index.html_0HTTP/1.1\nHost:_" + dest + "\n\n";

        try(Socket sock = new Socket(dest, 80)) { // web servers are usually listening at port 80
            try(OutputStreamWriter w = new OutputStreamWriter(sock.getOutputStream())) {
                w.write(request);           // write the HTTP request [9-11]
                w.flush();                 // make sure that all data has been sent
                sock.shutdownOutput();     // closing down the channel for sending data to the server

                try (InputStreamReader is = new InputStreamReader(sock.getInputStream());
                    BufferedReader r = new BufferedReader(is)) { // Baidu uses UTF-8 encoding
                    while ((response = r.readLine()) != null) { // read strings line-by-line until connection closed by server
                        System.out.println(response);           // print to output
                    }
                }
            }
        }
    }
}
```

- Text is a very complex variable-length data structure.
- Historically, there exist many different mappings from characters to bits and bytes.
- Unicode assigns an integer number to a character.
- UTF-8 defines how such a number can be translated to a variable-length list of bits.
- UTF-8 is now the prevalent text encoding in the internet, i.e., you should store all your text-based documents (txt, html, xml, ...) in UTF-8 encoding.

谢谢

Thank you

Thomas Weise [汤卫思]
tweise@hfu.edu.cn
<http://www.it-weise.de>

Hefei University, South Campus 2
Institute of Applied Optimization
Shushan District, Hefei, Anhui,
China



Caspar David Friedrich, "Der Wanderer über dem Nebelmeer", 1818
http://en.wikipedia.org/wiki/Wanderer_above_the_Sea_of_Fog



1. *ISO/IEC 8859-1 – Final Text of DIS 8859-1, 8-bit Single-Byte Coded Graphic Character Sets – Part 1: Latin Alphabet No.1*, volume ISO/IEC 8859-1:1997 (E). Geneva, Switzerland: International Organization for Standardization (ISO), February 12, 1998. URL <http://std.dkuug.dk/jtc1/sc2/wg3/docs/n411.pdf>.
2. Ken Lunde. *CJKV Information Processing*. Sebastopol, CA, USA: O'Reilly & Associates, Inc., 1999. ISBN 0-596-51447-6 and 1-56592-224-7. URL http://examples.oreilly.de/english_examples/cjkvinfo/AppE/gb2312.pdf.
3. *Information Technology – Universal Coded Character Set (UCS) (ISO/IEC 10646:2011)*. Geneva, Switzerland: International Organization for Standardization (ISO), 2011.
4. USA: The Unicode Consortium Mountain View, CA and Julie D. Allen. *The Unicode Standard, Version 5.0*. Reading, MA, USA: Addison-Wesley Professional, fifth edition, 2007. ISBN 0-321-48091-0 and 978-0-321-48091-0. URL <http://books.google.de/books?id=Yn1UAAAAMAAJ>.
5. The unicode consortium, 2011. URL <http://www.unicode.org/>.
6. Jukka K. Korpela. *Unicode Explained*. Internationalization Documents, Programs, and Web Sites. Sebastopol, CA, USA: O'Reilly Media, Inc., June 28, 2006. ISBN 059610121X and 9780596101213. URL <http://books.google.de/books?id=PcWU2yxc8WkC>.
7. François Yergeau. *STD 63: UTF-8, A Transformation Format of ISO 10646*, volume 3629 of *Request for Comments (RFC)*. Network Working Group, November 2003. URL <https://tools.ietf.org/html/rfc3629>.
8. Herbert Schildt. *Java 2: A Beginner's Guide*. Essential Skills for First-Time Programmers. Maidenhead, England, UK: McGraw-Hill Ltd., 2002. ISBN 0072225130 and 9780072225136. URL <http://books.google.de/books?id=YWDJGYaLG4C>.
9. Timothy John Berners-Lee, R. Fielding, and H. Frystyk. *Hypertext Transfer Protocol – HTTP/1.0*, volume 1945 of *Request for Comments (RFC)*. Network Working Group, May 1996. URL <http://tools.ietf.org/html/rfc1945>.
10. R. Fielding, J. Gettys, Jeffrey Mogul, H. Frystyk, L. Masinter, P. Leach, and Timothy John Berners-Lee. *Hypertext Transfer Protocol – HTTP/1.1*, volume 2616 of *Request for Comments (RFC)*. Network Working Group, June 1999. URL <http://tools.ietf.org/html/rfc2616>.
11. David Gourley and Brian Totty. *HTTP: The Definitive Guide*. Definitive Guide. Sebastopol, CA, USA: O'Reilly Media, Inc., 2002. ISBN 1565925092 and 9781565925090. URL http://books.google.de/books?id=qEo019bcV_cC.