

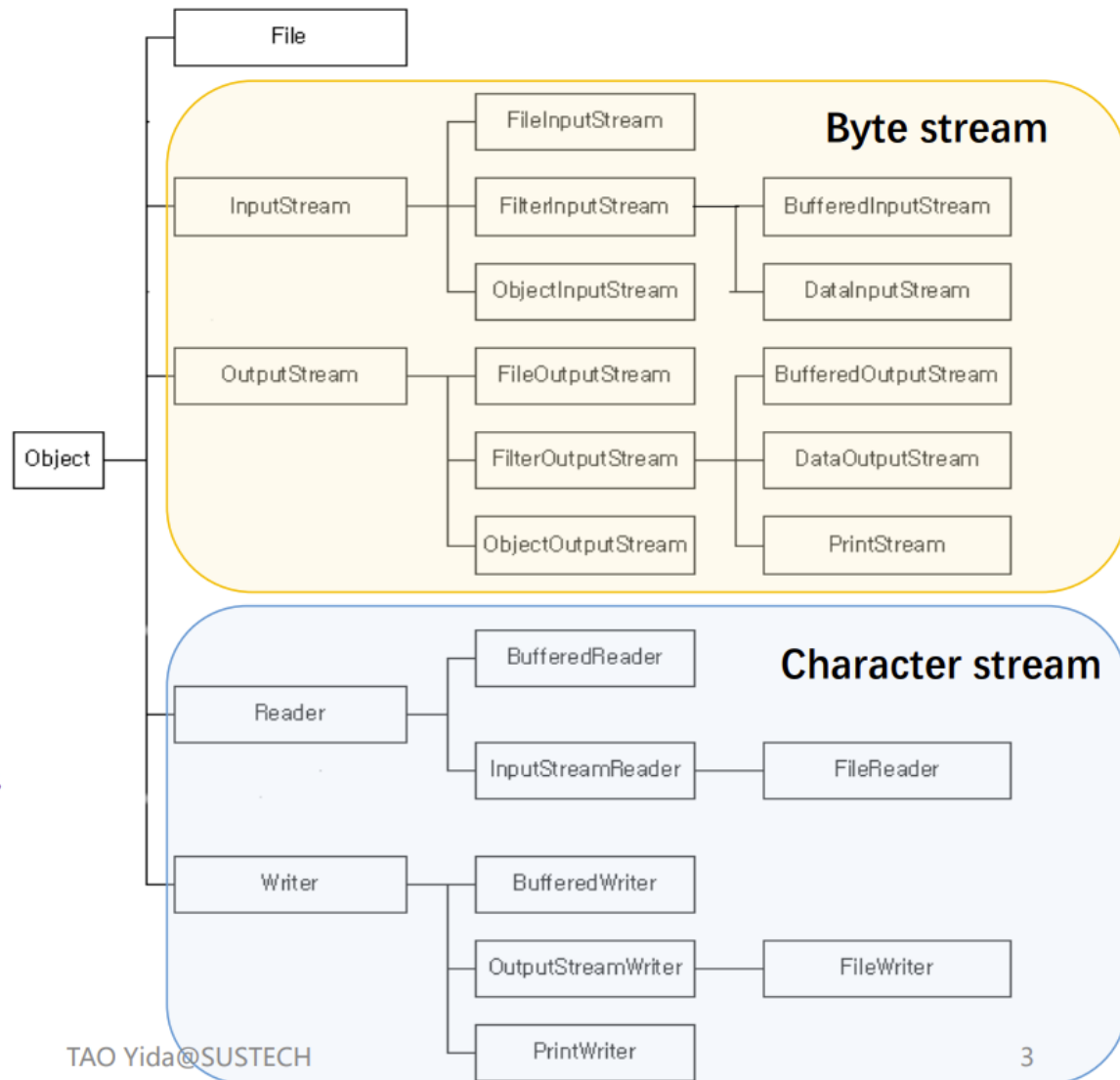
Lecture 5 & Lab 5

I/O Overview and encoding

1. overview

in `java.io` package

hierarchy:



2. Encoding

convert characters to other formats, often numbers, in order to store and transmit them more effectively.

1. ASCII

- represent text in computers
- using 7 bits to represent 128 chars
- extended ASCII uses 8 bits for 256 characters

2. GB2312, GBK, GB18030

- GB stands for 国标
- GB2312 uses 2 bytes(cover 99% usages)
- GBK extends GB2312
- GB18030 extends GBK

3. problems

different contries with different language systems implement their own character encoding

4. Unicode

- motivate by the need to encode characters in all languages consistently without conflicts
- Unicode is a standard (defines the mapping to code point)

UTF-8

- Uses a minimum of 1 byte, but if the character is bigger, then it can use 2, 3 or 4 bytes.
- is compatible with the ASCII table

UTF-16

- uses a minimum of 2 bytes. UTF-16 can not take 3 bytes, it can either take 2 or 4 bytes
- is not compatible with the ASCII table

UTF-32

- always uses 4 bytes
- is not compatible with the ASCII table

5. utf-8

- UTF-8 stands for “Unicode Transformation Format – 8-bit”
- Characters are encoded with varied lengths (1~4 bytes)

Character Range	Encoding
0...7F	0a ₆ a ₅ a ₄ a ₃ a ₂ a ₁ a ₀
80...7FF	110a ₁₀ a ₉ a ₈ a ₇ a ₆ 10a ₅ a ₄ a ₃ a ₂ a ₁ a ₀
800...FFFF	1110a ₁₅ a ₁₄ a ₁₃ a ₁₂ 10a ₁₁ a ₁₀ a ₉ a ₈ a ₇ a ₆ 10a ₅ a ₄ a ₃ a ₂ a ₁ a ₀
10000...10FFFF	11110a ₂₀ a ₁₉ a ₁₈ 10a ₁₇ a ₁₆ a ₁₅ a ₁₄ a ₁₃ a ₁₂ 10a ₁₁ a ₁₀ a ₉ a ₈ a ₇ a ₆ 10a ₅ a ₄ a ₃ a ₂ a ₁ a ₀

6. Java char implementation

- 16-bit unsigned int (U+0000~U+FFFF), corresponding to Unicode code points
- Conversion between int and char refers to the Unicode mapping
- Characters whose code points are **greater than U+FFFF** are called **supplementary** characters

- Supplementary characters are represented **as a pair of char values** (16 + 16 = 32 bits / 4 bytes)

```
int v1 = 0x0454; // Hex
System.out.printf("%c\n", v1); //e
System.out.printf("%c\n", (char)v1); //e

int v2 = 1108; // Decimal
System.out.printf("%c\n", v2); //e
System.out.printf("%c\n", (char)v2); //e

int v3 = 0x10454; // Hex
System.out.printf("%c\n", v3); //ð
System.out.printf("%c\n", (char)v3); //e
```

```
int v1 = 0x0454;
int v3 = 0x10454;
char[] c1 = Character.toChars(v1); // length 1
char[] c2 = Character.toChars(v3); // length 2
System.out.println(c1); //e
System.out.println(c2); //ð
```

Byte streams & Character Streams

1. difference

Byte Stream

- for processing **bytes** file, such as picture, music, or large file.
- Byte streams are **inconvenient for processing info stored** in Unicode
- Used to

Character Stream

- for processing **characters** file, such as `.txt`, `.csv`, `.html`...
- These classes have read and write operations that are based on **char** values rather than **byte values**

2. Similarity

- Byte Stream: `InputStream` & `OutputStream` are abstract classes
Character stream: `Reader` & `Writer` are abstract classes
- Byte Stream: Subclasses are all called "xxxStream"
Character stream: "xxxReader" & "xxxWriter"
- Subclasses for `InputStream` or `Reader` must implement `read()`

- Subclasses for `OutputStream` or `Writer` must implement `write()`

3. `FileInputStream`: Used for reading streams of raw **bytes**

- `fileInputStream.read()` will read a bytes a time even if text encoding in **utf-8**
- Same file with different encoding will generate different outcome in the reading.
- `System.out.println((char)input.read());` will not give a right answer for **utf-8** encoding file
- useful methods:

```
1 // Constructor
2 FileInputStream(String name);
3 // open the file named by the path name name in the file system.
4
5 // read from file by bytes
6 public int read() throws IOException;
7 // Return the next byte of data, or -1 if the end of the file is reached
8 public int read(byte[] b) throws IOException;
9 // Reads up to b.length bytes of data from this input stream into an array
  of bytes.
10 public long skip(long n) throws IOException;
11 // Skips over and discards n bytes of data from the input stream.
12
13 // close the file
14 public void close() throws IOException;
```

4. `FileReader`: Used for reading streams of **characters** (instead of streams of **bytes**)

1. Java system/platform default encoding

- Differs from **OS** and **language settings** (e.g., **GBK** on 中文操作系统)
- Could be changed (environment variable, IDE, code)

2. File encoding

- Independent from Java
- Could be changed
- When using Java to read a file, the Java system default encoding and the file encoding should be **consistent**

3. Set a right encoding in `FileReader`

```
1 Reader r = new FileReader("src/test.txt", Charset.forName("gb2312"))
```

4. `InputStreamReader`:

- superclass of `FileReader`
- use to transfer `InputStream` to `reader`

```

1 // Constuctor of InputStreamReader
2 InputStreamReader(InputStream in, String charsetName);
3
4 Reader reader = new InputStreamReader(new
  FileInputStream("src/test.txt"), "UTF-8")

```

5. Useful methods in `FileReader`: almost the same

```

1 // Constuctor
2 // using the platform's default charset
3 public FileReader(String fileName) throws FileNotFoundException;
4 // use the specified charset
5 public FileReader(String fileName, Charset charset) throws IOException;
6
7 // method extends from InputStreamReader
8 public int read() throws IOException;
9 public int read(char[] cbuf, int off, int len) throws IOException;
10
11 // close
12 public void close() throws IOException;

```

5. Similarly, methods in `FileOutputStream` and `FileWriter`

- `FileOutputStream`

```

1 // constuctor
2 // If append is true, then bytes will be written to the end of the
  file
3 public FileOutputStream(String name, boolean append) throws
  FileNotFoundException;
4
5 // write()
6 public void write(int b) throws IOException;
7 public void write(byte[] b, int off, int len) throws IOException;
8
9 // close()
10 public void close() throws IOException;

```

- `FileWriter`

```

1 // constructor
2 FileWriter(File file, Charset charset, boolean append);
3
4 // write
5 // a character takes at most 4 bytes
6 public void write(int c) throws IOException;
7 public void write(char[] cbuf, int off, int len) throws IOException;
8 public void write(String str, int off, int len) throws IOException;
9
10 // close()
11 public void close() throws IOException;

```

6. Refer to more details

[FileInputStream \(Java SE 17 & JDK 17\)_\(oracle.com\)](#)

[FileReader \(Java SE 17 & JDK 17\)_\(oracle.com\)](#)

[FileOutputStream \(Java SE 17 & JDK 17\)_\(oracle.com\)](#)

[FileWriter \(Java SE 17 & JDK 17\)_\(oracle.com\)](#)

FilterInputStream

FilterInputStream allows us to **wrap an existing input stream** by extending the InputStream class, and **filter** or **modify** it when reading data.

This approach allows us to perform some **transformations** or **manipulations** on the data as we read it, such as **decrypting or decompressing the data**.

1. FilterInputStream Example I

```

1 InputStream zfile = new GZIPInputStream( // + gzip functionality
2     new BufferedInputStream( // + buffered functionality
3         new FileInputStream(
4             "src/test.zip"
5         )
6     )
7 )

```

2. FilterInputStream Example II

```

1 PushbackInputStream pbin =
2 new PushbackInputStream( // + pushback functionality
3     new BufferedInputStream( // + buffered functionality
4         new FileInputStream("test.txt")));
5
6 int b = pbin.read();
7 pbin.unread(b); // added functionality to push back bytes

```

3. FilterInputStream Example III

```

1  DataInputStream din =
2  new DataInputStream( // + read-numbers functionality
3      new ZipInputStream( // + zip functionality
4          new FileInputStream("test.zip")));
5
6  // added abilities to read numbers, not only bytes
7  din.readDouble();
8  din.readInt();

```

4. direct known subclasses

- `BufferedInputStream`

Buffer the input in an array, enhance the speed of reading

- `DataInputStream`

A data input stream lets an application read **primitive Java data types** from an underlying input stream in a machine-independent way.

- `DigestInputStream`

record the digest message (like a summary)

- `InflaterInputStream`

decompress data in `.zip`, `.gzip`, `.zlib`

- `LineNumberInputStream`

`read()` line by line

- `PushbackInputStream`

Add functionality that push back something has just read.

- etc.....

Reading/Writing Text Input/Output

1. When working with I/O, we often work with **human-readable** text rather than binary data

2. Java provide two APIs to assist working with text I/O

- **Scanning**: useful for breaking down formatted input into **tokens** and translating individual tokens according to their data type (`Scanner`)
- **Formatting**: assembles data into nicely formatted, **humanreadable** form (`PrintWriter`)

3. Using Scanner for reading text files

```

1  Scanner in = new Scanner(new File("input.txt"));
2  // read the text file in the specified size
3  while(in.hasNextDouble()){
4      double value = in.nextDouble();
5      // process the value
6  }

```

4. Using PrintWriter for writing text files

```

1 | PrintWriter out = new PrintWriter("output.txt");
2 | out.println("hello, world");
3 | out.printf("total: %8.2f\n", total);

```

The `PrintWriter` class is an enhancement of the `PrintStream` class that you already know—`System.out` is a `PrintStream` object. You can use the familiar `print`, `println`, and `printf` methods with any `PrintWriter` object:

5. Constructing a `Scanner` and `PrintWriter`

```

1 | PrintWriter(String file_name);
2 | Scanner(File file); // remember that it's not a filename
3 | Scanner(InputStream source);

```

I/O from the Command Line

Standard Streams

- Standard streams read input from the keyboard and write output to the display.
- Java platform supports three Standard Streams
- `System` is a `final` class with **private constructor**

Fields		
Modifier and Type	Field and Description	
static <code>PrintStream</code>	<code>err</code> The "standard" error output stream.	<code>System.err</code>
static <code>InputStream</code>	<code>in</code> The "standard" input stream.	<code>System.in</code>
static <code>PrintStream</code>	<code>out</code> The "standard" output stream.	<code>System.out</code>

`System.in`

- Standard input, often read keyboard input
- To use **Standard Input** as a **character stream**, wrap `System.in` in `InputStreamReader`.
- `System.in` with `Scanner`: Parse the input into different primitive types and strings

```

1 | Scanner in= new Scanner(System.in);

```

`System.out`

- `System.out` is defined as a `PrintStream` object
- Although it is technically a byte stream, `PrintStream` **utilizes an internal character stream object to emulate many of the features of character streams** (same for `PrintWriter`)
- Could use `setOut()` to **redirect** the output to other resources


```
// construct a new PrintStream with a specified file
PrintStream out = new PrintStream(new File("src/sysout.txt"));
// re-assign the standard output from console to file
System.setOut(out);
// this will be written to file
System.out.println("where am I?");
```

notes: out is a final field, for a specialization, compiler will change it if we use `setOut()`

- `System.out.println`: All things will be printed with no filter.