

CS209

Computer system design and application

Stéphane Faroult
faroult@sustc.edu.cn

Zhao Yao zhaoy6@sustc.edu.cn

PERSISTENCE

Working in memory is nice, but everything goes when the computer is turned off (or crashes)

Need to ^{save, keep} persist
program results
somewhere

As I said last time, that's all the idea of persistence.
"Somewhere" may take multiple shapes (including a remote computer)

So we have to work as much as we can in memory, and only in memory, for speed ...

Mostly work in
memory for speed

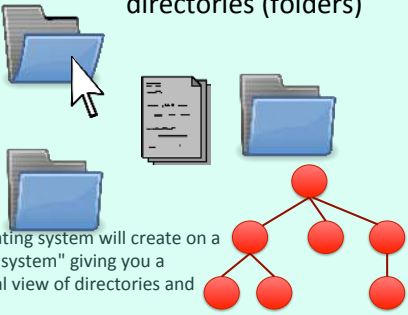
Memory

... while we still need a safety net.

Write to file for
safety

File System

directories (folders)



Your operating system will create on a disk a "file system" giving you a hierarchical view of directories and files.

Stream redirection

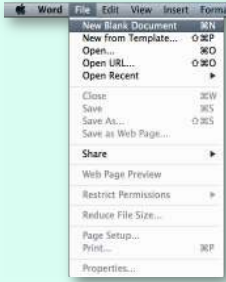
Instead of keyboard

```
$ java MyProgram < input_file
```

Instead of screen

```
$ java MyProgram > output_file
```


One very easy way to save data and to restore it is, instead of doing it in Java, to let the system do it through what is known as "Stream redirection". Your program may read from what it thinks is the keyboard when in reality input comes from a file, and what is written to the screen can also be redirected to a file. System.out and System.err can be separated.



Files may be handled directly by programs.


Many program, though, will directly handle file either when you specifically request it, or in your back (for instance "autosave" every ten minutes).

Two traditional types of files:



Text files

only printable characters



Binary files

If there are tons of file extensions and many types of files, they all fall in one of two categories: text or binary. The only problem is, whatever they are, they are all made of 0s and 1s. So really it's all a matter of interpretation.

Interpretation is the big, hard question. You cannot guess the meaning of 1s and 0s just by looking at them. You must have an idea already. And even with text, there are many different ways to encode one single character (and don't believe that the problem doesn't exist even with basic Latin letters – there is another encoding system than ASCII called EBCDIC and the bits meaning 'a' in ASCII mean '/' in EBCDIC). If you haven't the key allowing you to decrypt the bits, you are lost.

HOW to understand the 0's and 1's?

Two traditional types of files:



Text files ← only printable characters
WHEN DECODING AS
CHARACTERS



Binary files

So the true definition of "text file" is that it only contains characters that you can print (including spaces and carriage returns) when you decrypt the file as a bunch of characters.

Text files

There are many types of text files – not only files with extension .txt !

Can be opened by a "text editor" (eg Notepad) or displayed using **more** (Linux) or **type** (Windows)

Program code (.c, .h, .py, .php, .java, .bat, .sh, ...)

Plain text (.txt, .ini)

Text with readable tags (.html, .rtf, .xml)

Data as text (.csv)

Only contain printable characters

Binary files

There are also many types of binary files, including those used by documents that are supposed to be mostly text ...

Can only be opened by a special program

Compiled program (.o, .exe, .class)

Archives, compressed files (.tar, .tgz, .zip)

Crypted files

Text with non readable formatting (.docx, .pdf, .xlsx, .pptx)

Multimedia (.gif, .jpg, .png, .mp3, .wav, .mpg, .flv, ...)

Binary files

Very often (but not always) binary files are a basic "dump" of what you have in memory. When the file may be written on one system and read on a very different one, some standard encoding may be applied.

Memory structures written "as is"

More compact (no difference for text)

No conversion during I/Os

May be portability issues between computers (Windows/Mac/Linux)

One big problem is the "small endian"/"big endian" issue, which is a hardware issue. The 4 bits that make up half a byte may be swapped.

Binary files

Descriptive Header

DATA

Bunch of attributes

Most often a binary file contains a header part that describes the structure, followed by data proper.

Bunch of 0s and 1s

Stream redirection

Text files
(most often)

```
$ java MyProgram < input_file
```

```
$ java MyProgram > output_file
```

Although nothing forbids writing to/reading from binary files, stream redirection is mostly used with text files.

What is important to understand is that whether you are calling `println` with an integer, whether you are explicitly calling `toString()` or not, a conversion occurs.

Integer.toString(int_val);

Binary internal computer representation

The number has to be turned into a string of digits for output.

```

if number is negative
    display '-'
loop on decreasing powers of 10
    get the result r of the integer division of the number by the power of 10
    if we have already displayed a non zero digit
        display the digit corresponding to the code of '0' plus r
    else
        if the digit is not zero
            record that we have found a non zero digit
            display the digit corresponding to the code of '0' plus r
        decrease the number by r times the power of 10 processed
end of loop

```

Input requires the opposite.

Integer.parseInt() or the method **nextInt()** of a **Scanner** object perform the reverse operation

HOW to understand the 0's and 1's?

```

public class Hello {

    public static void main(String[] args) {
        System.out.println("Hello!");
    }
}

```

Let's check **Hello.class**

A program has, to "understand a file", a number of options.

HOW to understand the 0's and 1's?



Assume that it's what we expect

for instance, text ...

The simple one is "when the only tool you have is a hammer, everything looks like a nail".

```
$ cat Hello.class
????4

<init>()VCodeLineNumberTablemain([Ljava/lang/String;)V
SourceFile
Hello.java

    Hello!
        ellojava/lang/Objectjava/lang/
SystemoutLjava/io/PrintStream;java/io/
PrintStreamprintln(Ljava/lang/String;)V!    *??

    %   ???

$ "cat" writes everything as text to the screen. The result looks
like garbage.
```

HOW to understand the 0's and 1's?

2 Assume that the extension is correct

I renamed Hello.class to Hello.c and tried to compile it. I got 105 warnings and 11 errors but the compiler tried.

HOW to understand the 0's and 1's?

3 Check the file header for a "magic number"

Binary files usually contain a "signature" in their header, a small number of bytes that are very specific to one type of files. You don't need to trust the extension.

```
$ od -x Hello.class
00000000 feca beba 0000 3400 1d00 000a 0006 090f
00000020 1000 1100 0008 0a12 1300 1400 0007 0715
00000040 1600 0001 3c06 6e69 7469 013e 0300 2928
00000060 0156 0400 6f43 6564 0001 4c0f 6e69 4e65
00000100 6d75 6562 5472 6261 656c 0001 6d04 6961
00000120 016e 1600 5b28 6a4c 7661 2f61 616c 676e
00000140 532f 7274 6e69 3b67 5629 0001 530a 756f
00000160 6372 4665 6c69 0165 0a00 6548 6c6c 2e6f
00000200 616a 6176 000c 0007 0708 1700 000c 0018
00000220 0119 0600 6548 6c6c 216f 0007 0c1a 1b00
...
```

od is a Unix command that dumps a file.
All .class files start with the same bytes.

feca beba
cafe babe

Most file-related classes are in the java.io package (there is also a java.nio). Two types of streams, Character or Byte, which can all be buffered or unbuffered.

Handling files in Java

Key concept: **Stream**

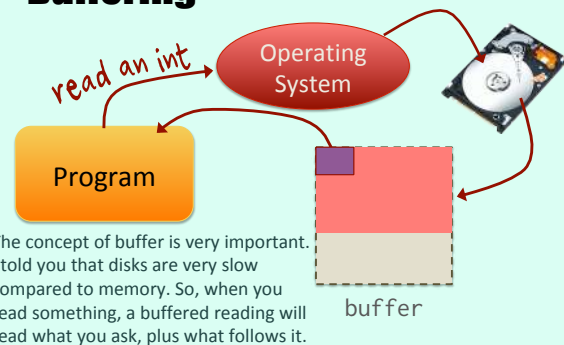
Character streams

Buffered or unbuffered

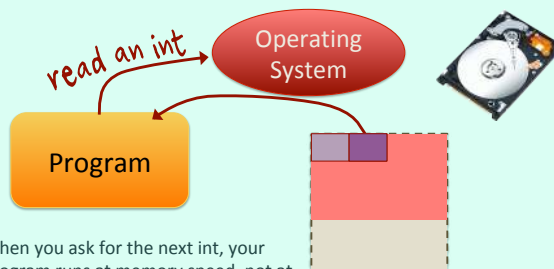
Byte streams

```
import java.io.*
```

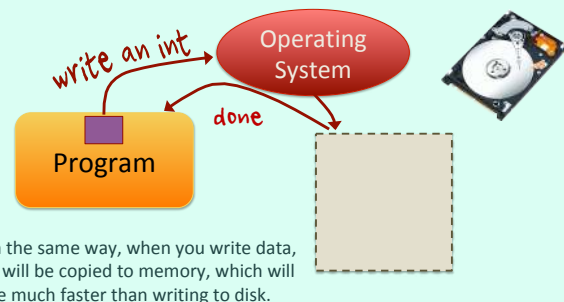
Buffering



Buffering

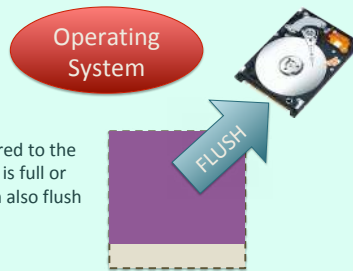


Buffering



Buffering

Data will be bulk-transferred to the disk only when the buffer is full or you close the file (you can also flush buffers explicitly)

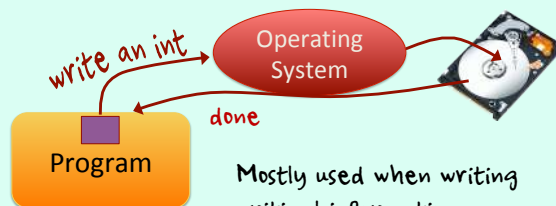


Buffering

What happens if the system restarts after a crash?

The answer of course is that what was in buffers is lost. Not a problem when reading, big problem when writing. It's not always easy to know what to replay.

No buffering



If you freak about crashes, you should use slow unbuffered operations.

Mostly used when writing critical information

Logs
Messaging

Performance - copying a 11M CSV file character by character

Test on my Mac (internal SSD)

| | |
|------------|-------------|
| Unbuffered | about 34.5s |
| Buffered | about 1s |

Test on my Mac (External USB HD)

| | |
|------------|------------|
| Unbuffered | about 36s |
| Buffered | about 1.4s |

How often does your computer crash? Would it be a complete disaster to run the program again after restart?


Performance - copying a 11M CSV file
character by character

For 99% of cases, you should use buffered input/output operations (we could even say 100% for input).

Use **buffered** operations
unless writing safely is a
critical concern.

Or for debugging.

Sometimes hard to say what's going
on today

Big disk systems have their own,
battery protected, buffers (also
called **CACHE**)  French for Hideout

Note that especially with high-end storage you rarely have
one level of buffering (in which case unbuffered wouldn't
be what it seems). It's a bit hard sometimes to know if the
data is on disk or not, and the Cloud doesn't make it any
simpler.

```
InputStream in = null;
OutputStream out = null;
```

UNBUFFERED

```
in = new InputStream(...);
out = new OutputStream(...);
```

The basis for all binary Input/Output operations are
InputStreams and OutputStreams, which are unbuffered.

One thing that should not be forgotten with file
operations is that it's probably the part of a program
where everything can fail.

Lots of things can **GO WRONG**

Wrong file/directory path
Not allowed 

Content not as expected
Hardware problem (rare)

UNBUFFERED

```

InputStream in = null;
OutputStream out = null;

try {
    in = new InputStream(...);
    out = new OutputStream(...);
} catch ... {
} finally {
    if (in != null) {
        try {
            in.close();
        } catch (IOException ioe) {
            // Just ignore
        }
    }
} ...

```

So you should really do everything in a try block, either a "try with resources" or a try with a "finally" block to make sure that files are cleanly closed and not left corrupt.

IMPORTANT! ← Flush everything and close properly

UNBUFFERED

```

FileInputStream in = null;
FileOutputStream out = null;

try {
    in = new FileInputStream("filename");
    out = new FileOutputStream("filename");
} catch ... {
} finally {
    ...
}

```

File location is always a practical problem. Think of reflection and properties file.

Looks in the current directory unless you provide a full path

BUFFERED

```

BufferedInputStream in = null;
BufferedOutputStream out = null;

try {
    in = new BufferedInputStream(new InputStream(...));
    out = new BufferedOutputStream(new ...);
} catch ... {
} finally {
    ...
}

```

To turn an unbuffered stream into a buffered one, you just wrap the call to the stream constructor into a call to a buffered stream constructor.

InputStream and OutputStream are the parent classes for all byte streams.

Byte Streams

(not the most used)

There are libraries for multimedia

You may not use byte streams very often. Remember JavaFx: when you create a new Image or Media object, a binary file is read into memory by the constructor. There is necessarily a byte stream behind the scene, but it's all done by the constructor.

BYTE STREAM
(unbuffered)

```

FileInputStream in = null;
FileOutputStream out = null;

try {
    in = new FileInputStream("filename");
    out = new FileOutputStream("filename");

    ...

} catch ... {
} finally {
    ...
}

```

The examples I have previously shown are for byte streams ...

BYTE STREAM
(buffered)

```

BufferedInputStream in = null;
BufferedOutputStream out = null;

try {
    in = new BufferedInputStream(new FileInputStream());
    out = new BufferedOutputStream(new ...);

    ...

} catch ... {
} finally {
    ...
}

```

... including the buffered version.

CAUTION

When you talk about "bytes" in input/output operations, you are really dealing with **int** variables, not **byte** variables.

HISTORICAL REASON

Don't be misled by the "byte" in "byte stream". Operations deal with more than one byte.

Object Serialization

```

class Obj2 {
    private String name;
    private int value;
    ...
}

class Obj1 {
    private Obj2[] o = null;
    private int count = 0;

    public Obj1() {
        o = new Obj2[10];
    }
    ...
}

```

One application of byte streams is "serialization", dumping a memory object to file.

The diagram illustrates the memory heap. It shows a box for 'Obj1' at memory address 6300. Inside this box is an array of 'Obj2' objects. The array contains three references: one to 'Amazon' (at 6992), one to 'Nile' (at 6853), and one to 'Yangtse' (at 6300). Arrows indicate the references from the array in Obj1 to the individual objects in memory.

Object Serialization

When doing so you only want to store data, not memory addresses that will change when you reload.

```

class Obj2 {
    private String name;
    private int value;
    ...
}

class Obj1 {
    private Obj2[] o = null;
    private int count = 0;

    public Obj1() {
        o = new Obj2[10];
    }
    ...
    You can also declare some attributes
    as "transient" to say they shouldn't
    be dumped.
  
```

| | | |
|--------|---------|------|
| Amazon | 6992 | Nile |
| 6853 | Yangtse | 6300 |

Object Serialization: requisites

```

class Obj2 implements java.io.Serializable {
    private String name;
    private int value;
    ...
}

class Obj1 implements java.io.Serializable {
    private Obj2[] o = null;
    private int count = 0;

    public Obj1() {
        o = new Obj2[10];
    }
    ...
  
```

1 INTERFACE

To be able to do so you must first say that every object implements the Serializable interface.

NO method!

It's just a declaration, there is no method to implement (in fact, it just tells javac to generate what is needed)

Object Serialization: requisites

```

class Obj2 implements java.io.Serializable {
    private String name;
    private int value;
    ...
}

class Obj1 implements java.io.Serializable {
    private Obj2[] o = null;
    private int count = 0;

    public Obj1() {
        o = new Obj2[10];
    }
    ...
  
```

2 CONSTRUCTOR

YES!

A default constructor is helpful when recreating the object (although it looks that implementing Serializable seems to take care of that)

Object Serialization: requisites

And you need an "ObjectOutputStream" that is a special flavor of byte stream. This one comes by default with a buffer.

3 Stream

Buffer included

```

FileOutputStream fileOut = new FileOutputStream("file.dat");
ObjectOutputStream out = new ObjectOutputStream(fileOut);
out.writeObject(o);
out.close();
fileOut.close();
  
```

A file written on one computer can be read on any computer!

Now, I'm not impressed by performance. There may be cases when serialization performs very well, but it requires testing.

Character Streams

(most often)

Handle 16-bit unicode characters (**char** datatype)

You'll probably use Character Streams more often than byte streams. If you have a C background, don't forget that Java chars are 2 bytes, not one as in C.

CHARACTER STREAM

(unbuffered)

```
FileReader in = null;
FileWriter out = null;
```

```
try {
    in = new FileReader("filename");
    out = new FileWriter("filename");
```

```
} catch ... {
} finally {
    if (in != null) {
        in.close();
    }
    if (out != null) {
        out.close();
    }
}
```

What used to be "Input" and "Output" with byte streams becomes "Reader" and "Writer" with character streams.

```
BufferedReader in = null;
BufferedWriter out = null;
```

CHARACTER STREAM

(buffered)

```
try {
    in = new BufferedReader(new FileReader(...));
    out = new BufferedWriter(new FileWriter(...));
```

```
} catch ... {
} finally {
    if (in != null) {
        in.close();
    }
    if (out != null) {
        out.close();
    }
}
```

Otherwise it's exactly the same thing.

So what can character streams do that byte streams can't do?

Good question. In fact character streams have three features absent from byte streams.



1

They can read (write) lines

Need to use a **BufferedReader****readLine() method**

When buffered they can read or write a full line at once – because text files are usually a sequence of lines.

```
BufferedReader in = null;
String line;

try {
    in = new BufferedReader(new FileReader("..."));
    while ((line = in.readLine()) != null) {
        ...
    }
} finally {
    if (in != null) {
        try {
            in.close();
        } catch (IOException e) {
        }
    }
}
```

But beware of lines when text files were written on a system and are read on a **different** system.

%n in Java format

Which brings the interesting problem of lines. You probably know the carriage return character, '\n'. Java prefers '%n' with `printf()`, which may be one or two characters depending on the system it runs on.

It all dates back to the glorious days of the typewriter, in which letters were always typed at the same place. It's the "carriage" around which the sheet was wrapped that moved from right to left.

Where is the Life we have lost in living? ↵
 Where is the wisdom we have lost in knowledge? ↵
 Where is the knowledge we have lost in information? ↵

TS Eliot ↵
 (1888-1965)

What was happening at the end of the line? You needed to scroll the paper (line feed) and push back the carriage to the right (carriage return).

Guess what, the first printers were computer-controlled typewriters (sort of).



Where everything becomes fun, it's that a Unix system separates lines with a single 'carriage return' character.

Where is the Life we have lost in living? Where is the wisdom we have lost in knowledge? Where is the knowledge we have lost in information? TS Eliot (1888-1965)



\n

0x0A

00001010

While a Windows system still stores the two separate commands that used to be sent to the printer - line feed (represented by \r), then carriage return.

Where is the Life we have lost in living? Where is the wisdom we have lost in knowledge? Where is the knowledge we have lost in information? TS Eliot (1888-1965)



\r\n

0x0D0A

0000110100001010

Linux file in a Windows editor

Where is the Life we have lost in living? Where is the wisdom we have lost in knowledge? Where is the knowledge we have lost in information? TS Eliot (1888-1965)

A basic Windows editor (as Linux became more common, many editors became cleverer), looking for \r\n as separator between lines, will see in a Linux text file only one big line with \n characters that it doesn't know how to represent and that it will replace with squares.

Windows file in a Linux editor

Where is the Life we have lost in living?^W
Where is the wisdom we have lost in knowledge?^W
Where is the knowledge we have lost in information?^W
^W

TS Eliot^W
(1888-1965)^W

While on Linux the editor will understand correctly the \n in the Windows file, but will not know what to do with \r and will show it as ^W. There is nothing simple in IT.

*There are conversion programs
Many programs understand both.*

MOST character files are organized in lines (variable or constant length)

BUT a big character file can sometimes contain a single line (HTML, JSON, XML ...)

You would of course be wrong to believe that a big text file always contains many lines.

Another interesting characteristic is that the parent class of `FileReader` allows you to specify the character encoding.

2 They can change encoding

You can specify the character set in the constructor

```
graph TD; InputStreamReader[InputStreamReader] -.-> FileReader[FileReader];
```

The third characteristic will come next time ...