

# BYTE DATA COMMUNICATION

**Instructor: Prasun Dewan (FB 150, [dewan@unc.edu](mailto:dewan@unc.edu))**



# CASE STUDIES

XINU IPC: Design and Implementation

Unix Pipes

Java Sockets

Java Non-blocking IO

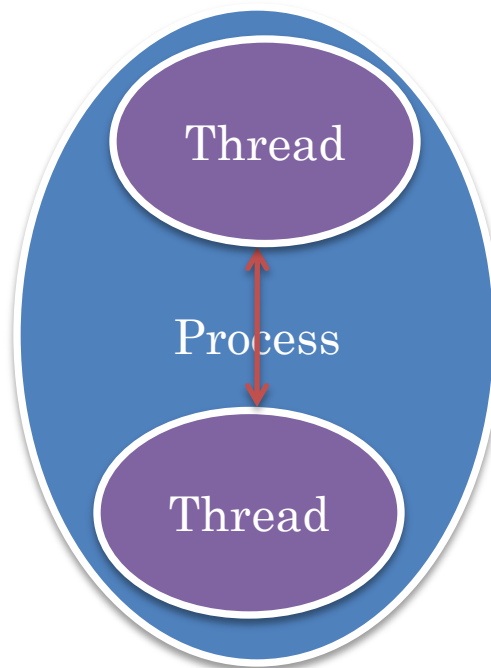


# XINU LOW-LEVEL MESSAGE PASSING

Focus is on simplicity of design and use



# LOCATION OF COMMUNICATING THREADS



Multiple address spaces not supported in XINU

Also true in Several PC OS's

No notion of a separate process;  
process = thread

Intra-address communication

# OPERATIONS

send (<thread\_id>, <int expression>)

Non blocking

int receive ()

Synchronous

int recvclr ()

Non blocking, polling, returns either message if it exists, otherwise a special value

Implementation?



# DATA STRUCTURE

stack

registers

program counter

priority

status

int receivedWord

bool hasMessage

Thread



# SEND OPERATION

stack

registers

program counter

priority

status

int receivedWord

bool hasMessage

Thread

```
send (tid, 5 )
```

```
send (tid, 6 )
```

```
send (tid, intExpression>)
```

```
If (! wordReceived) {  
    receivedWord = intExpression;  
    hasMessage = true;  
    if (waiting(tid)) {  
        ready(tid)  
    }  
}
```



# RECEIVE OPERATION

stack

registers

program counter

priority

READY

int receivedWord

bool hasMessage

Thread

receive( )

5

```
int receive ()
```

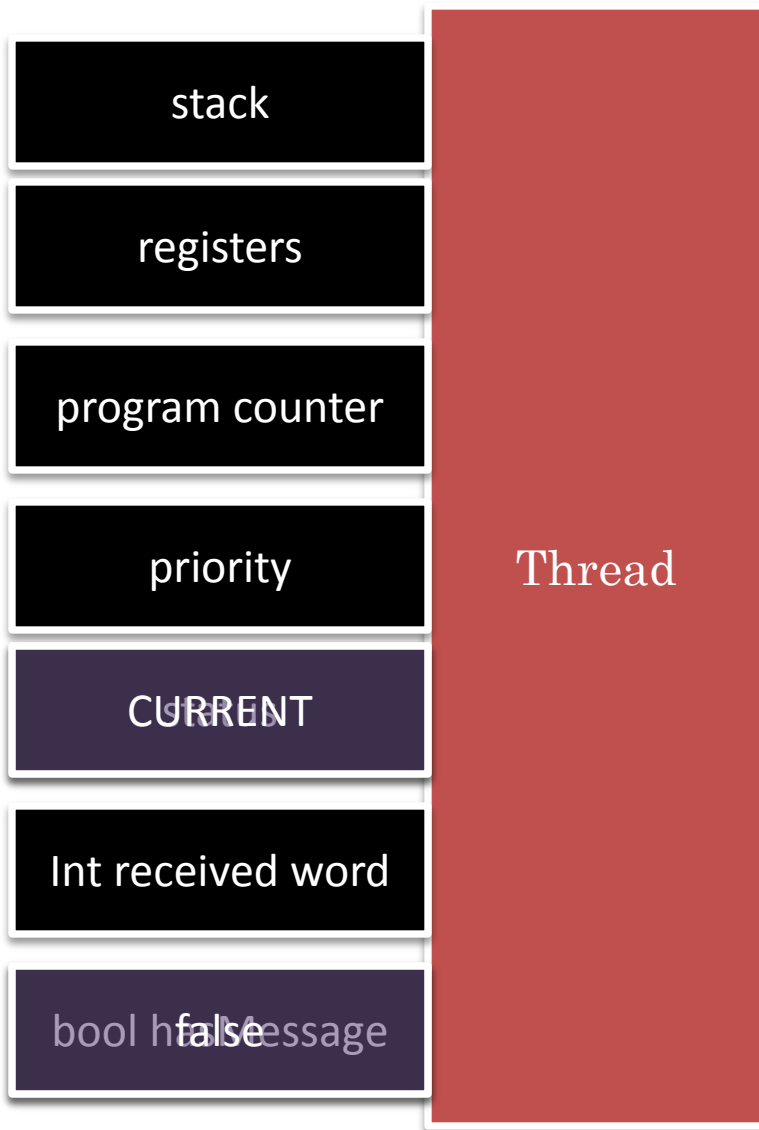
```
if (!wordReceived) {  
    status = RECEIVE;  
    blockAndResched(currentPid);  
}
```

```
hasMessage= false;  
return receivedWord;
```





# RECVCLR OPERATION



receive( )	NO_VAL
------------	--------

```
int receive ()
```

```
if (!word received) {  
    return NO_VAL  
}
```

```
hasMessage= false;  
return receivedWord;
```



# UNIX PIPES

```
bluetang(65)% man 2 pipe | more
PIPE(2)                               Linux Programmer's Manual
```

PIPE(2)

User command

## NAME

pipe, pipe2 - create pipe

## SYNOPSIS

```
#include <unistd.h>
```

```
int pipe(int pipefd[2]);
```

```
#define _GNU_SOURCE
```

```
#include <unistd.h>
```

```
int pipe2(int pipefd[2], int flags);
```

API call

## DESCRIPTION

pipe() creates a pipe, a unidirectional data channel that can be used for interprocess communication. The array pipefd is used to return two file descriptors referring to the ends of the pipe. pipefd[0] refers to the read end of the pipe. pipefd[1] refers to the write end of the pipe. Data written to the write end of the pipe is buffered by the kernel until it is read from the read end of the pipe. For further details, see pipe(7).

User command exposes part of  
the functionality of API call



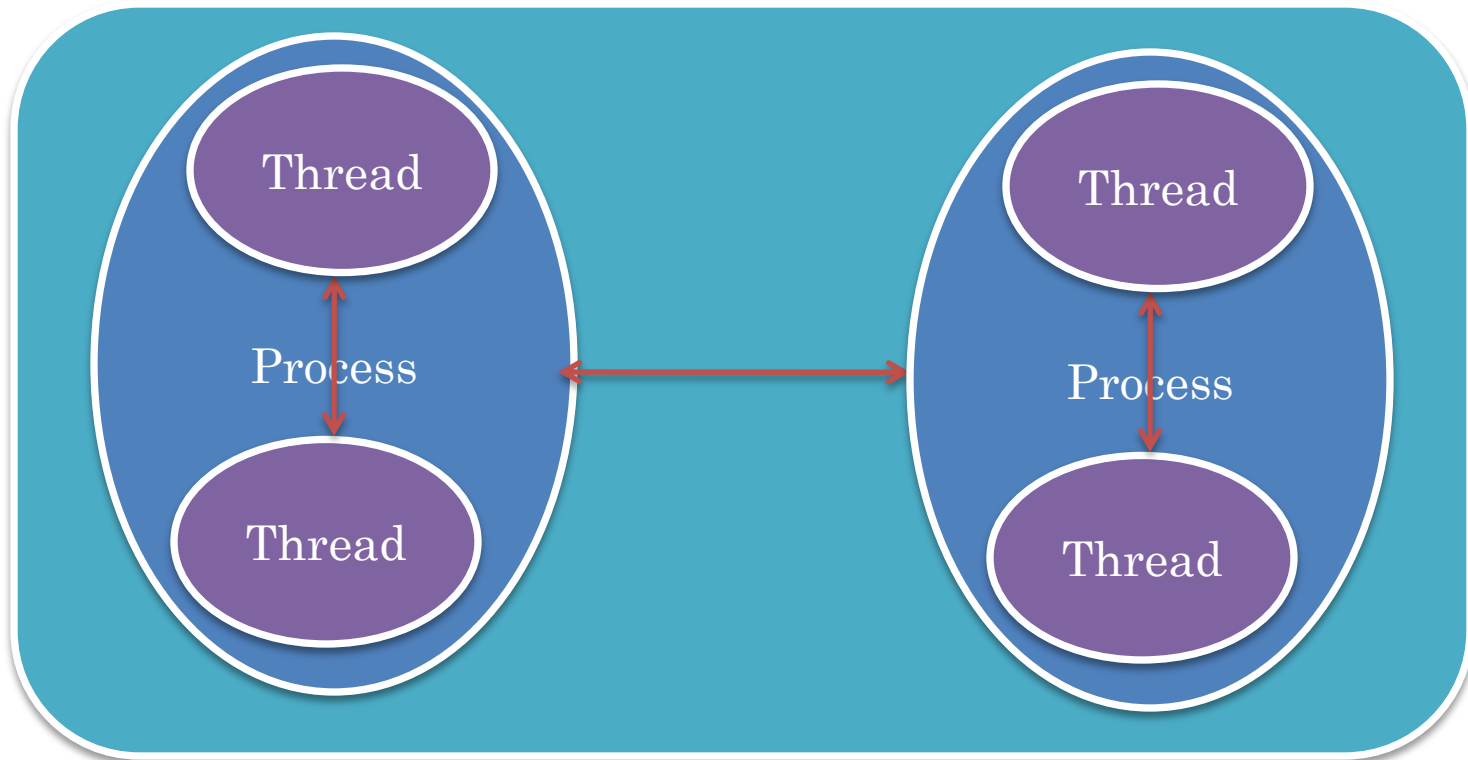
# UNIX PIPES

Focus is on Late Binding of Teletype I/O Source and Sink

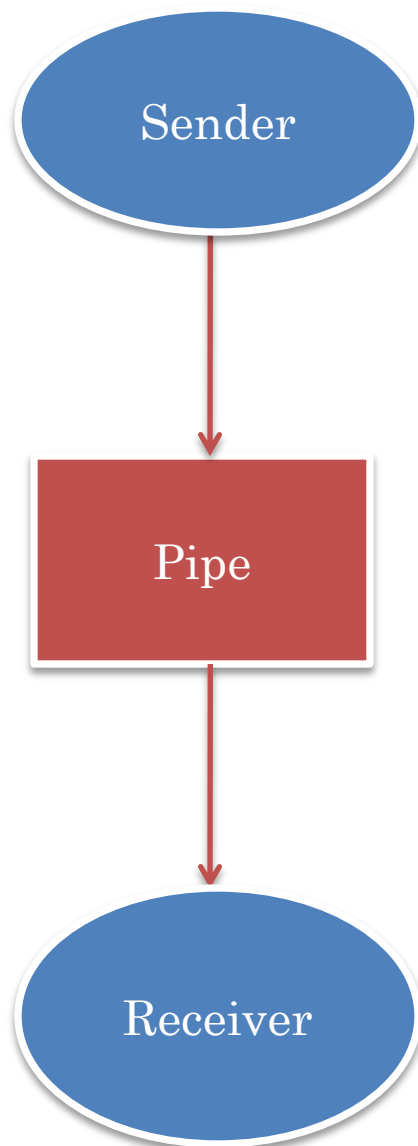
Stream-based communication



# INTRA-COMPUTER MESSAGE PASSING



# SIMPLEX PIPE?

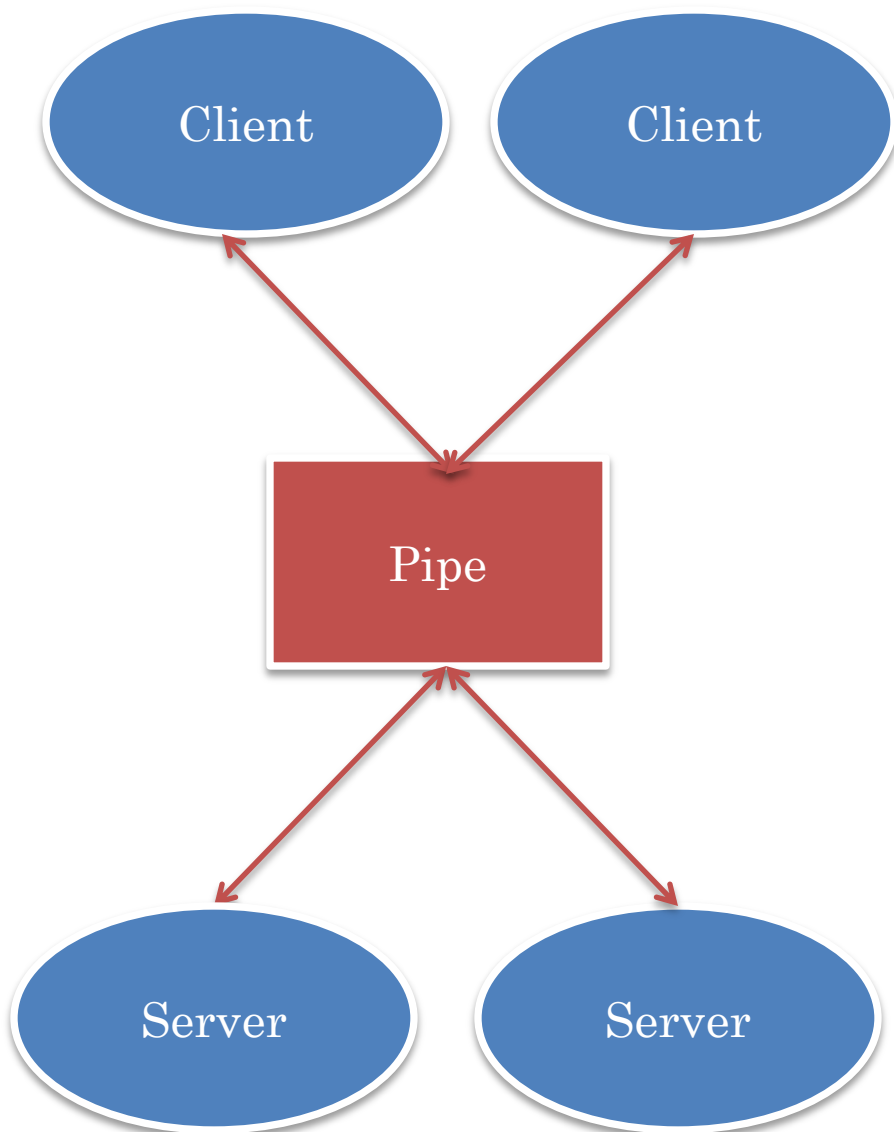


man 2 pipe | more

Command interpreter connects  
output of a process to input of  
another process



# NAMING?



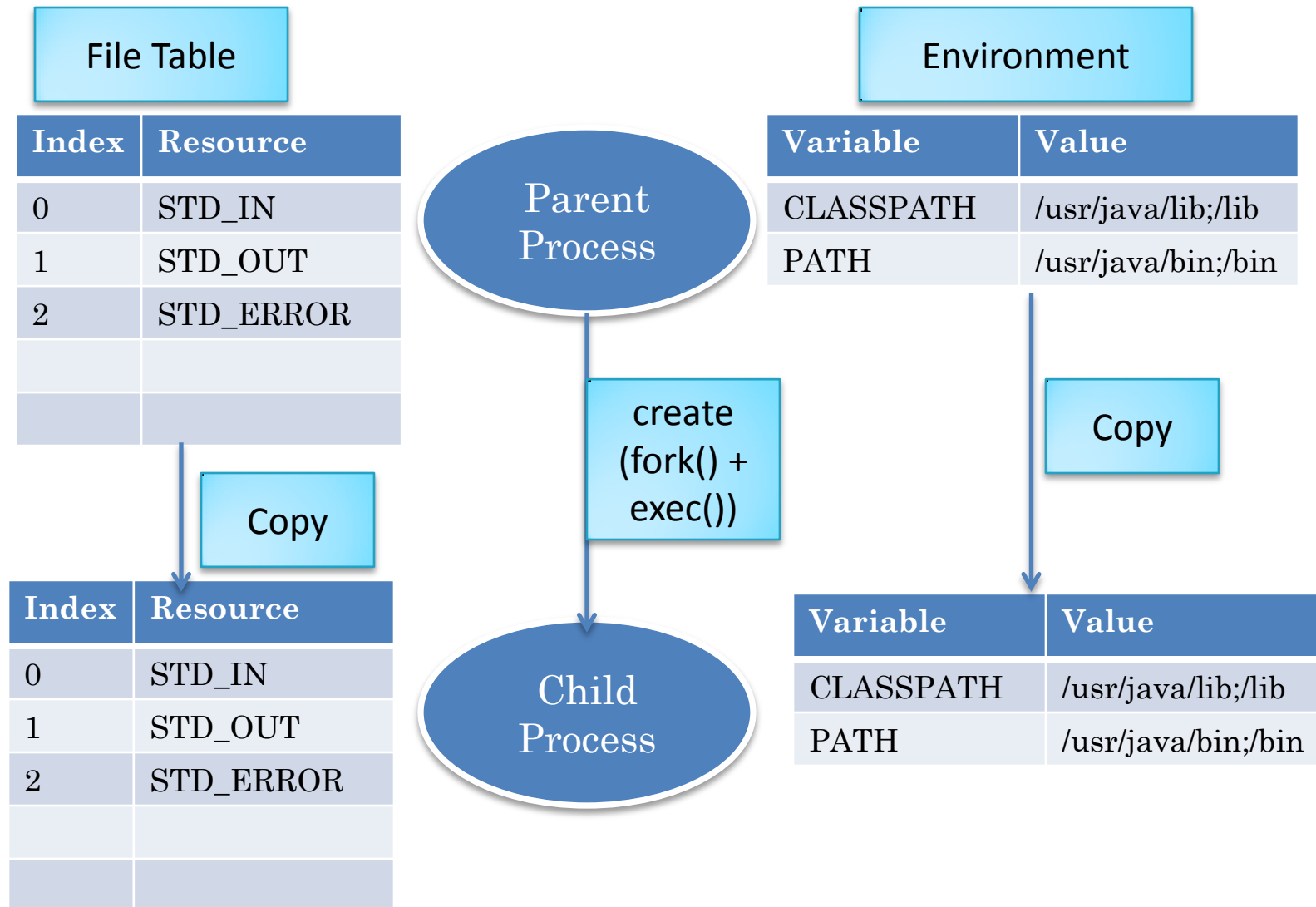
How do writers and readers name the pipe?

Memory is (usually) not shared by processes

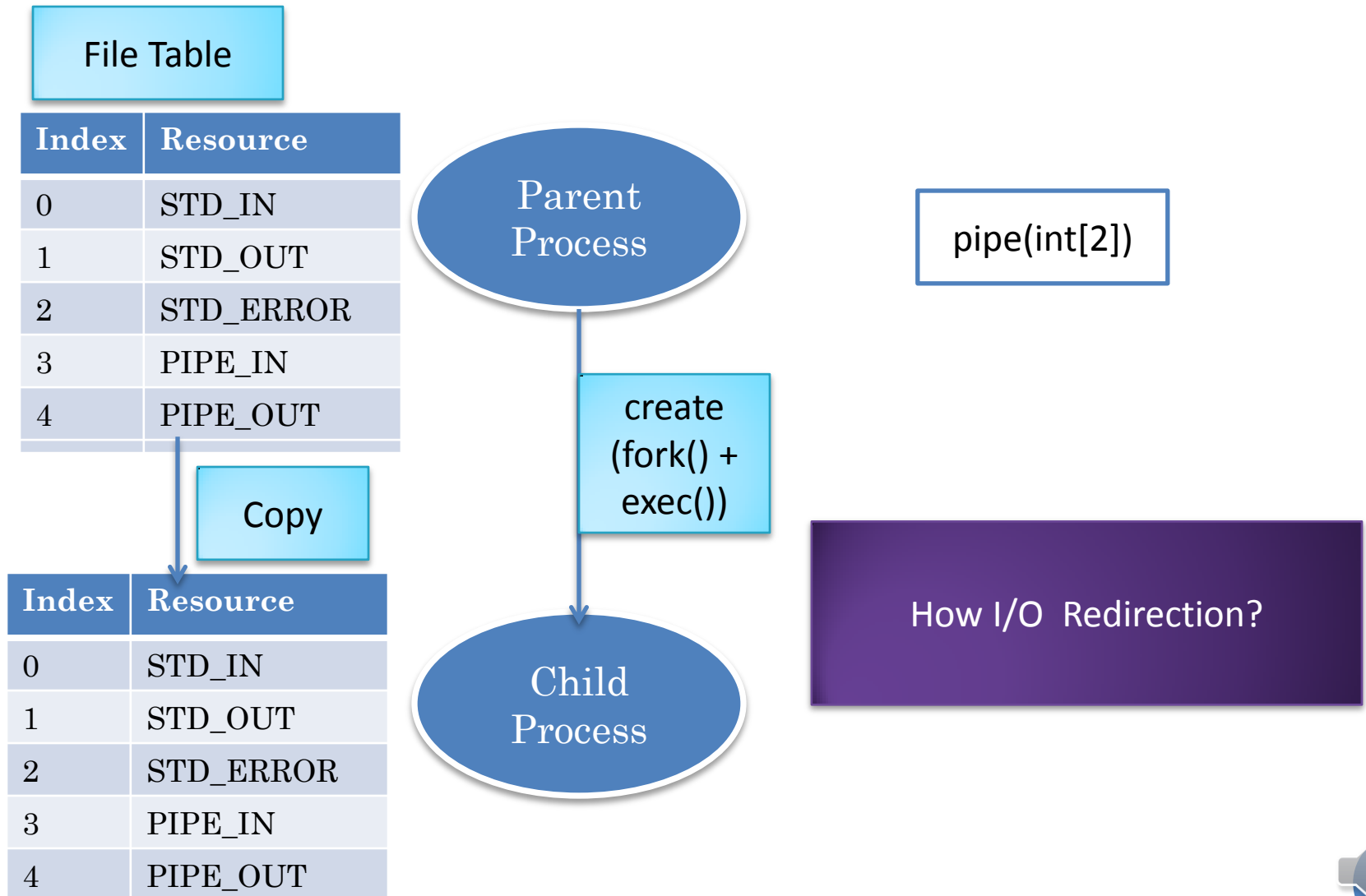
Child process inherits file descriptors (and environment variables) from parent process



# FILE AND ENVIRONMENT TABLE



# SHARING PIPE DESCRIPTORS





# COPYING FILE TABLE ENTRIES

File Table

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Copy

Index	Resource
0	STD_IN
1	PIPE_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Parent Process

create  
(Fork +  
Exec)

Child Process

pipe(int[2])

Can copy file table entry to another position (standard input/output) and can close it.



# I/O REDIRECTION DETAILS

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

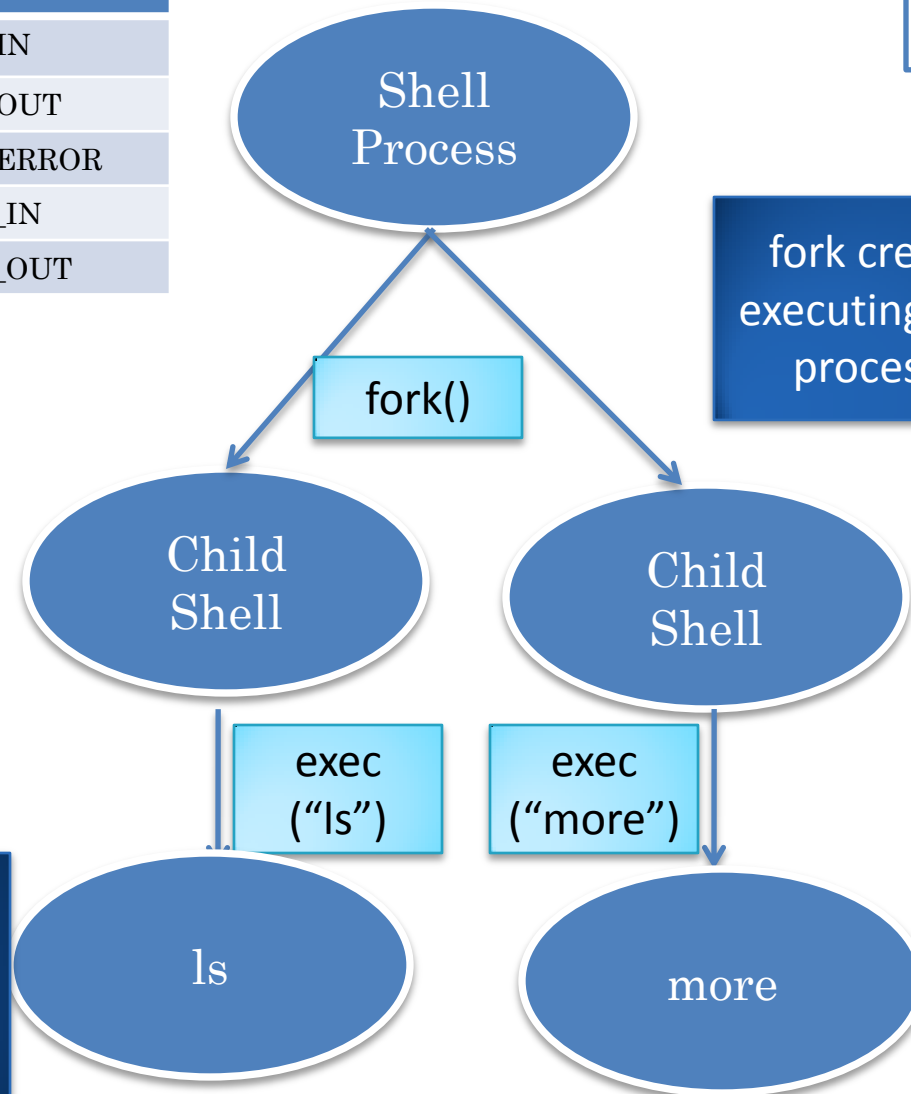
cat file | more

fork creates new process  
executing same program as  
process executing fork

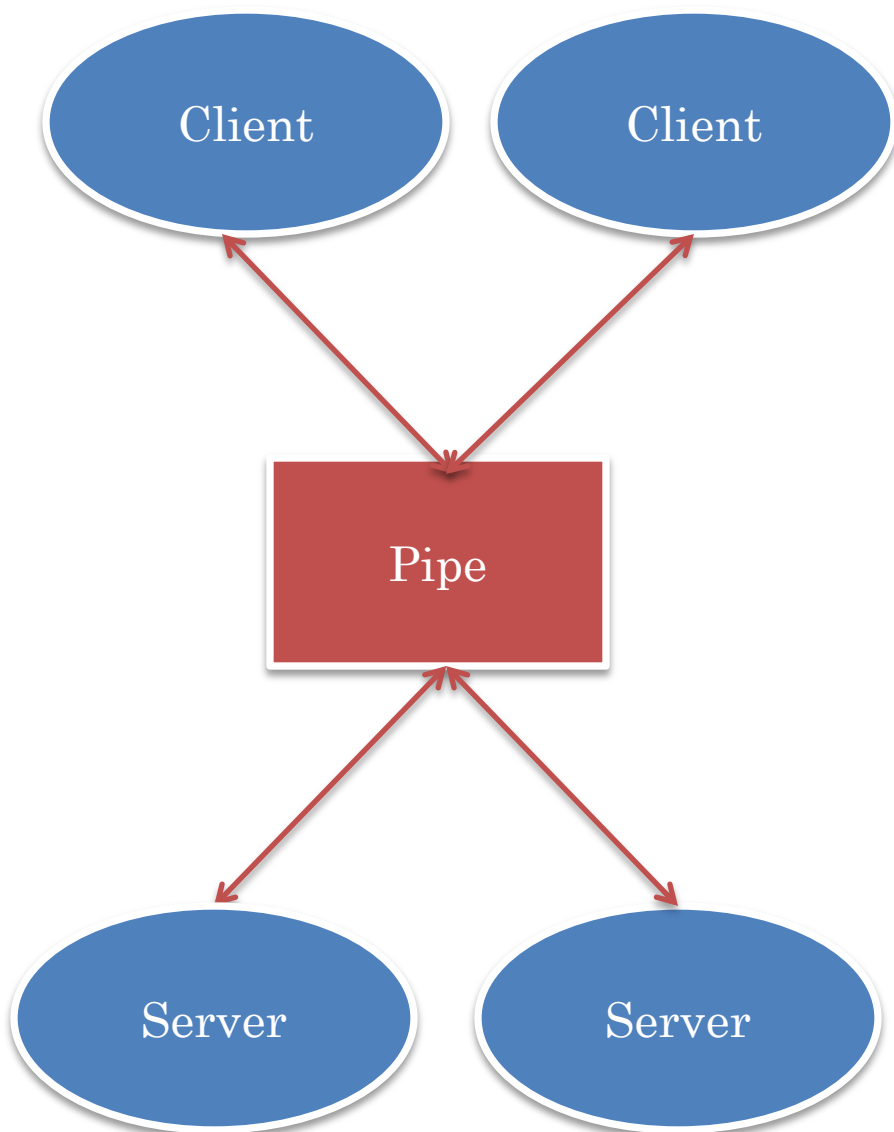
Index	Resource
0	STD_IN
1	PIPE_OUT
2	STD_ERROR

Index	Resource
0	PIPE_IN
1	STD_OUT
2	STD_ERROR

exec makes same  
process execute  
another program,  
keeping the file table



# PIPE (REVIEW)



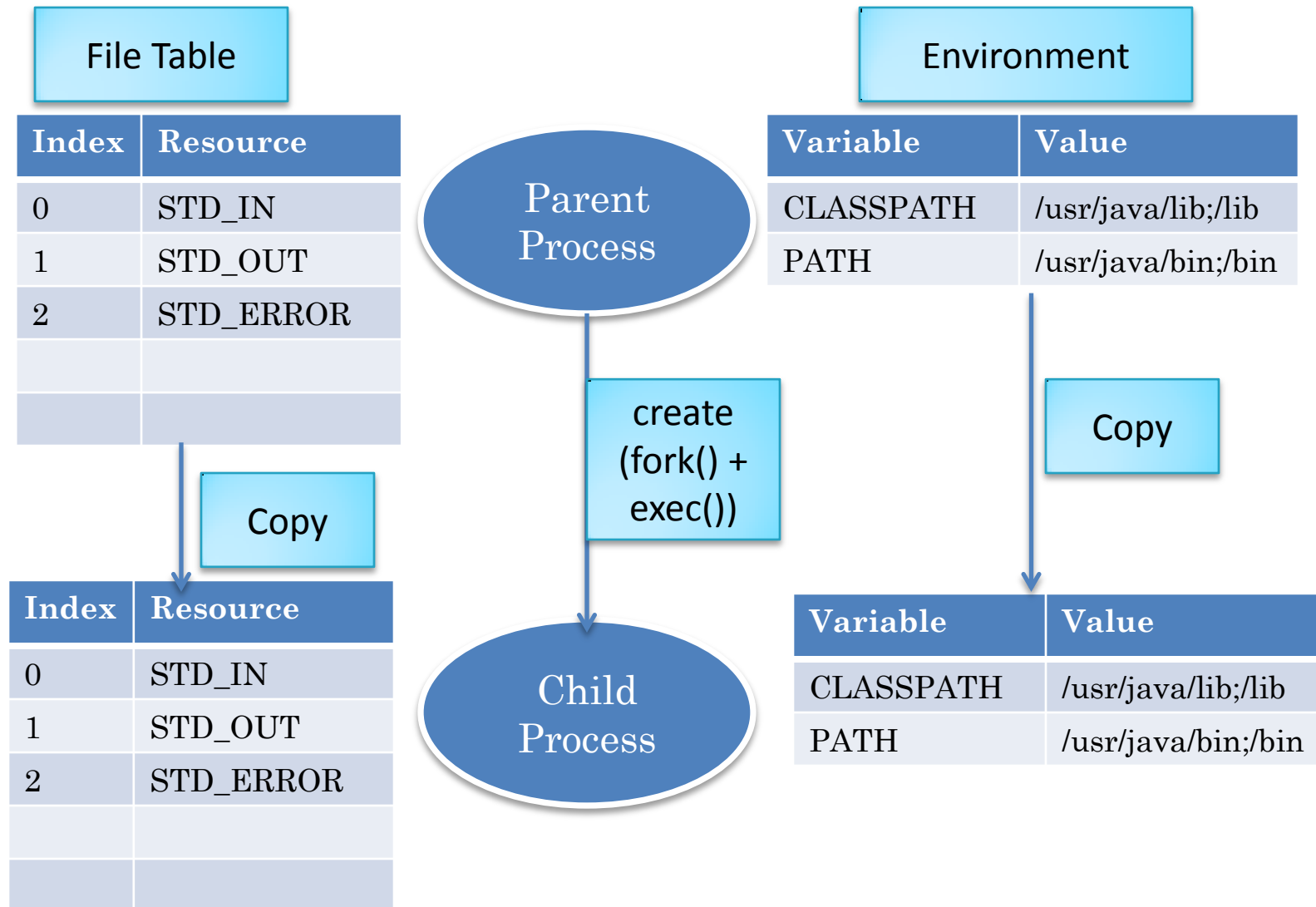
How do writers and readers name the pipe?

Memory is (usually) not shared by processes

Child process inherits file descriptors (and environment variables) from parent process



# FILE AND ENVIRONMENT TABLE (REVIEW)



# SHARING PIPE DESCRIPTORS (REVIEW)

File Table

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Copy

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Parent Process

create  
(fork() +  
exec())

Child Process

pipe(int[2])

How I/O Redirection?



# COPYING FILE TABLE ENTRIES (REVIEW)

File Table

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Copy

Index	Resource
0	STD_IN
1	PIPE_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

Parent Process

create  
(Fork +  
Exec)

Child Process

pipe(int[2])

Can copy file table entry to another position (standard input/output) and can close it.



# I/O REDIRECTION DETAILS (REVIEW)

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

cat file | more



fork()

fork creates new process  
executing same program as  
process executing fork



exec  
("ls")

exec  
("more")



exec makes same  
process execute  
another program,  
keeping the file table

Index	Resource
0	STD_IN
1	PIPE_OUT
2	STD_ERROR

Index	Resource
0	PIPE_IN
1	STD_OUT
2	STD_ERROR

# SEND/RECEIVE BLOCKING TIMES?: MESSAGE PIPE LINE

write(fd, byte)

send buffer

Process  
Thread

Operation started

system buffer

Message in Source  
System Buffer

Bounded buffer semantics

Sender waits for non full  
buffer

Receiver waits for non  
empty buffer

Allows lazy evaluation

Do not to computation that  
is not needed





# BLOCKING TIMES?: MESSAGE PIPE LINE

infinite\_output\_producer

```
static void main (String args[] {  
    while (true) {  
        printf("infinite output");  
    }  
}
```

infinite\_output\_producer | head - 2

infinite\_output\_producer blocks  
after filling buffer

head grabs first two lines from  
buffer, closes pipe, and  
terminates

Parent shell process waiting for  
head unblocks and kills  
infinite\_output\_producer



# PIPES: IMPLEMENTATION

Pipes

Pre 4.2 Unix BSD

Pipes

Sockets

4.2 Unix BSD

# PIPES: PROS AND CONS

I/O Redirection to Different Processes

Processes on same computer with  
common ancestor



# SOCKETS

Introduced by Berkeley Unix (4.2 BSD)

All OS's seem to have them in their basic form

Languages such as Java provide a layer above them



# SOCKETS

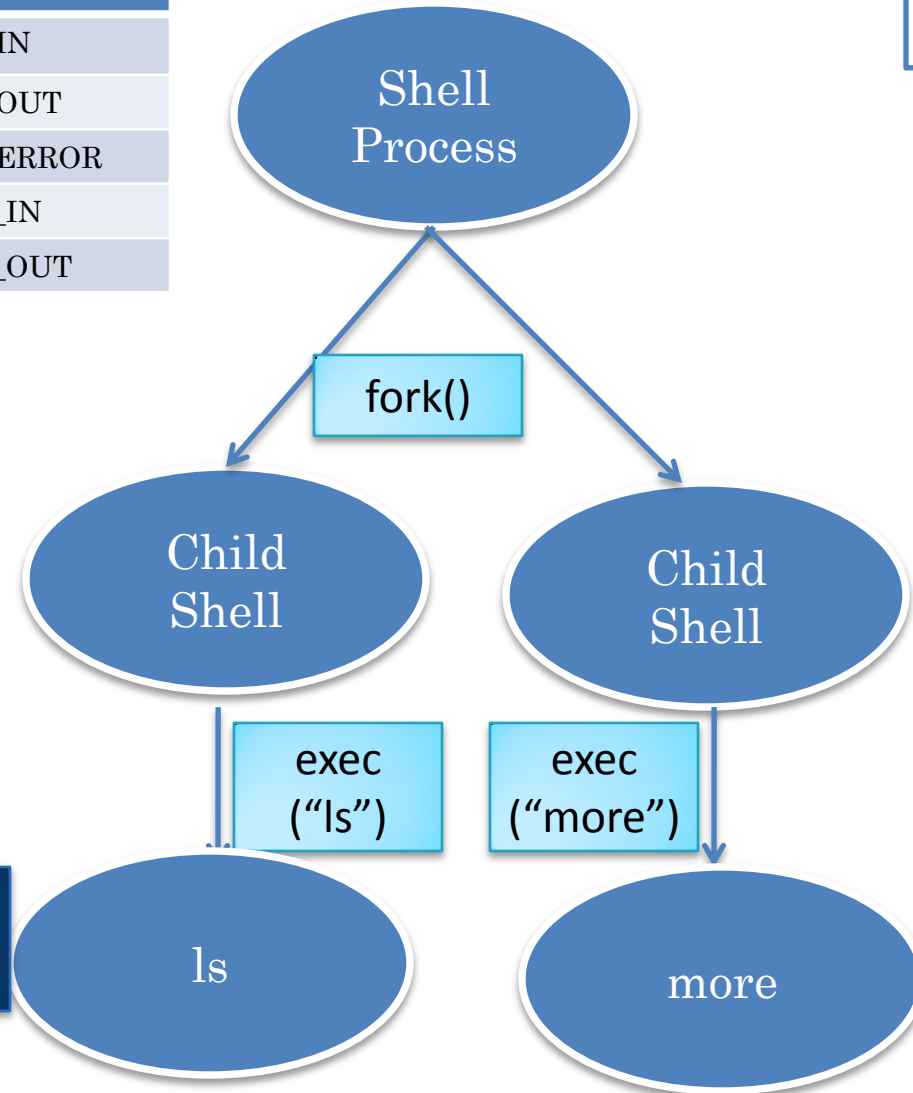
Focus is on generality and integration with  
File and Teletype I/O



# NAMING AND SHARING IN PIPES

Index	Resource
0	STD_IN
1	STD_OUT
2	STD_ERROR
3	PIPE_IN
4	PIPE_OUT

```
cat file | more
```



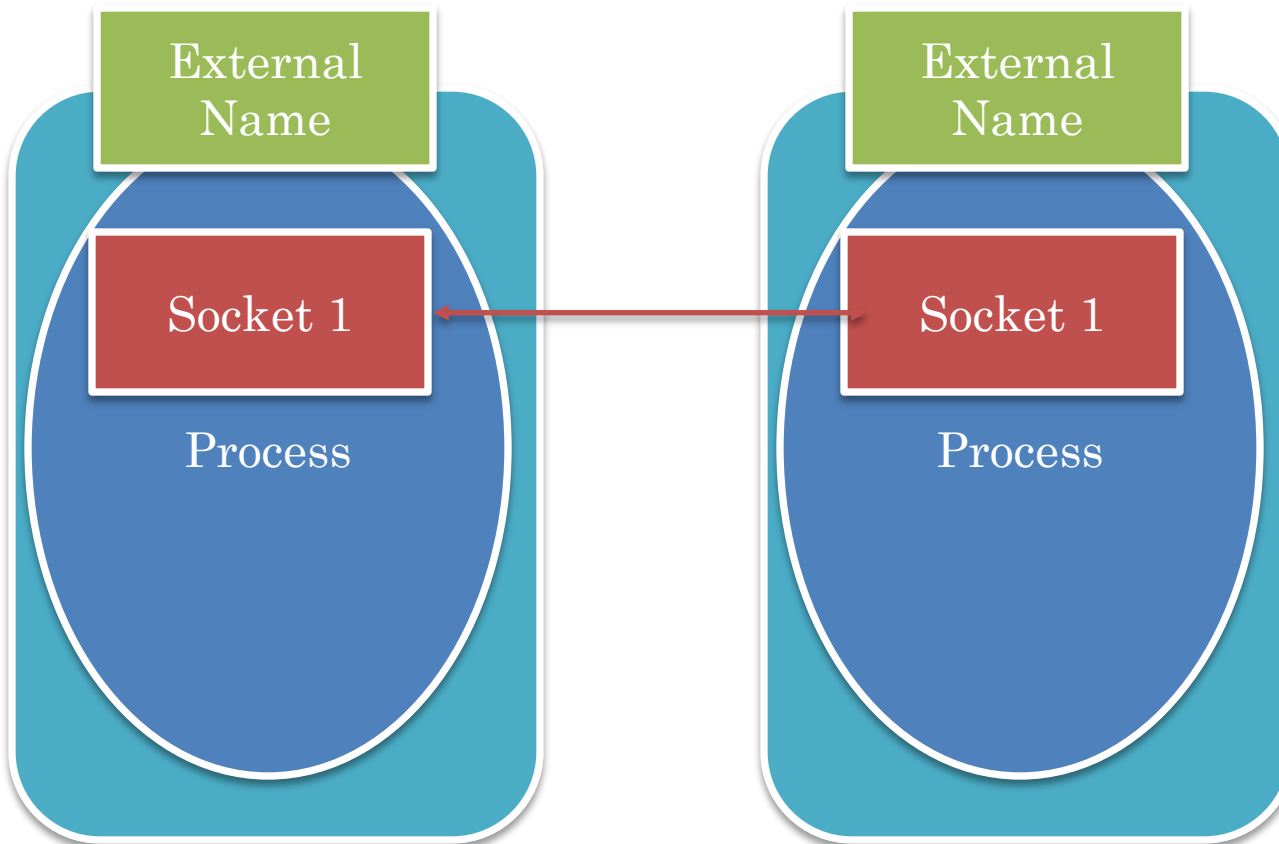
Index	Resource
0	STD_IN
1	PIPE_OUT
2	STD_ERROR

Index	Resource
0	PIPE_IN
1	STD_OUT
2	STD_ERROR

## Different hosts?



# NEED TO CONNECT DESCRIPTORS

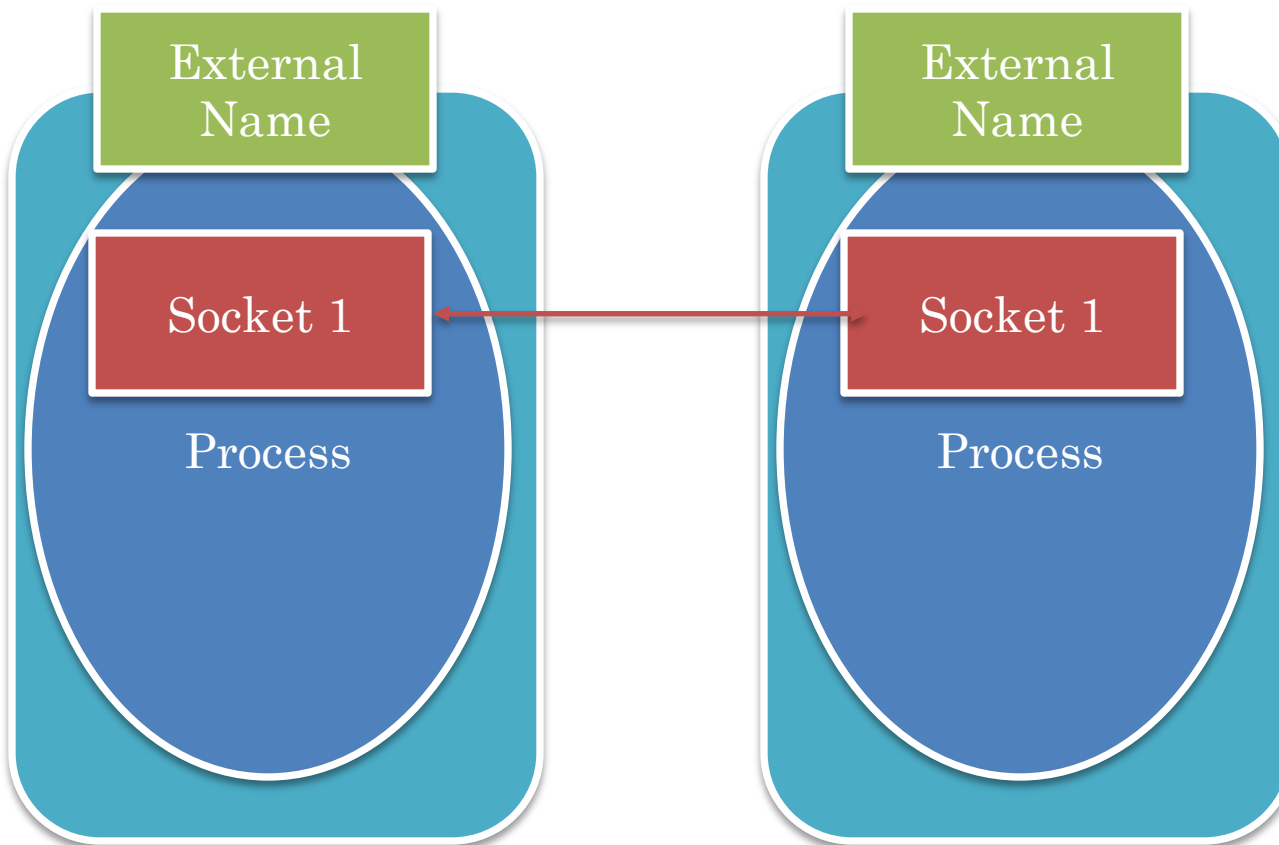


Some how message sent to socket in one process must be received at socket in a another process with no common ancestor

Need external names



# EXTERNAL NAME



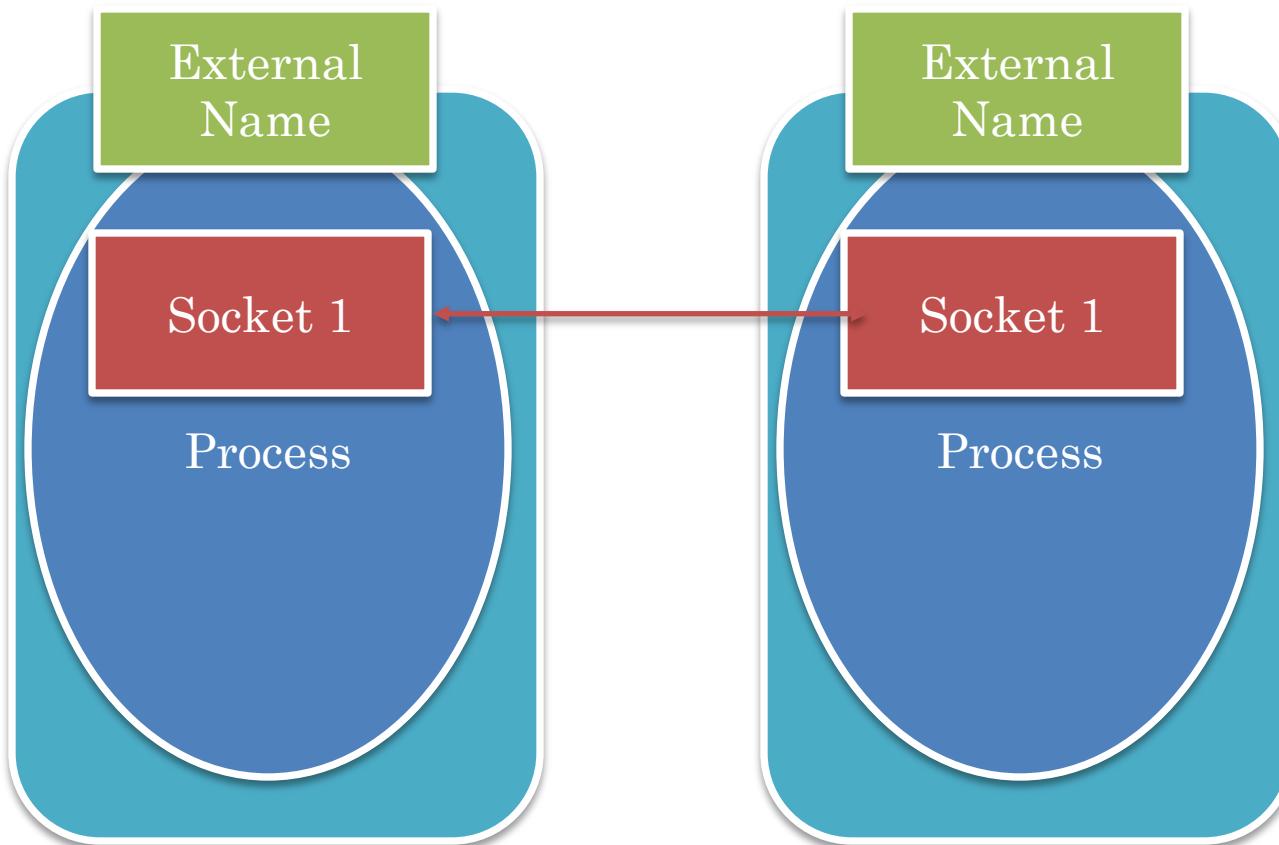
AF\_INET address family: host, port number (Java, Unix)

AF\_UNIX address family: file name (Unix)





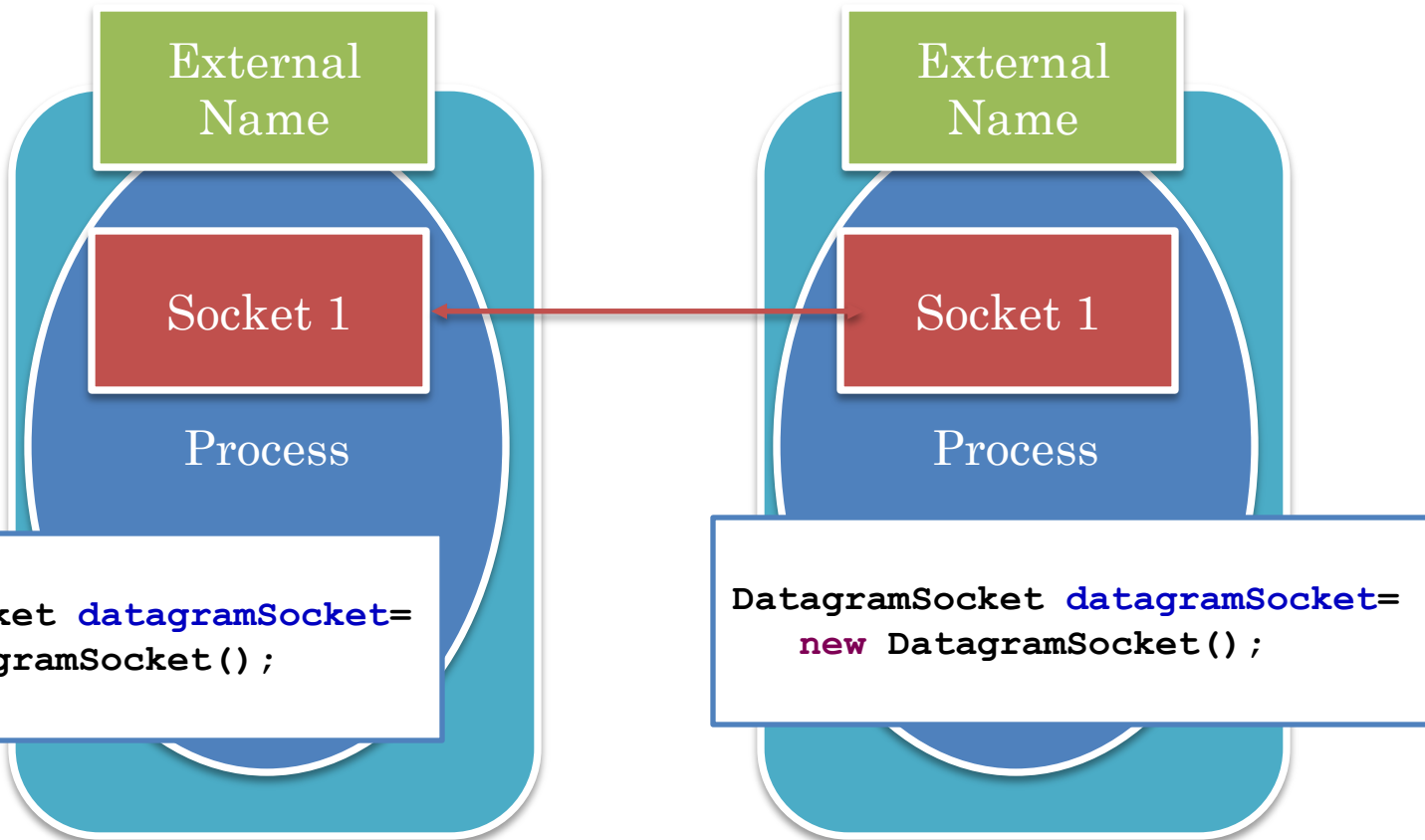
# EXTERNAL NAME



Somehow message sent to socket 1 must be received at socket 2



# DATAGRAM SOCKET: SEND GIVES DESTINATION

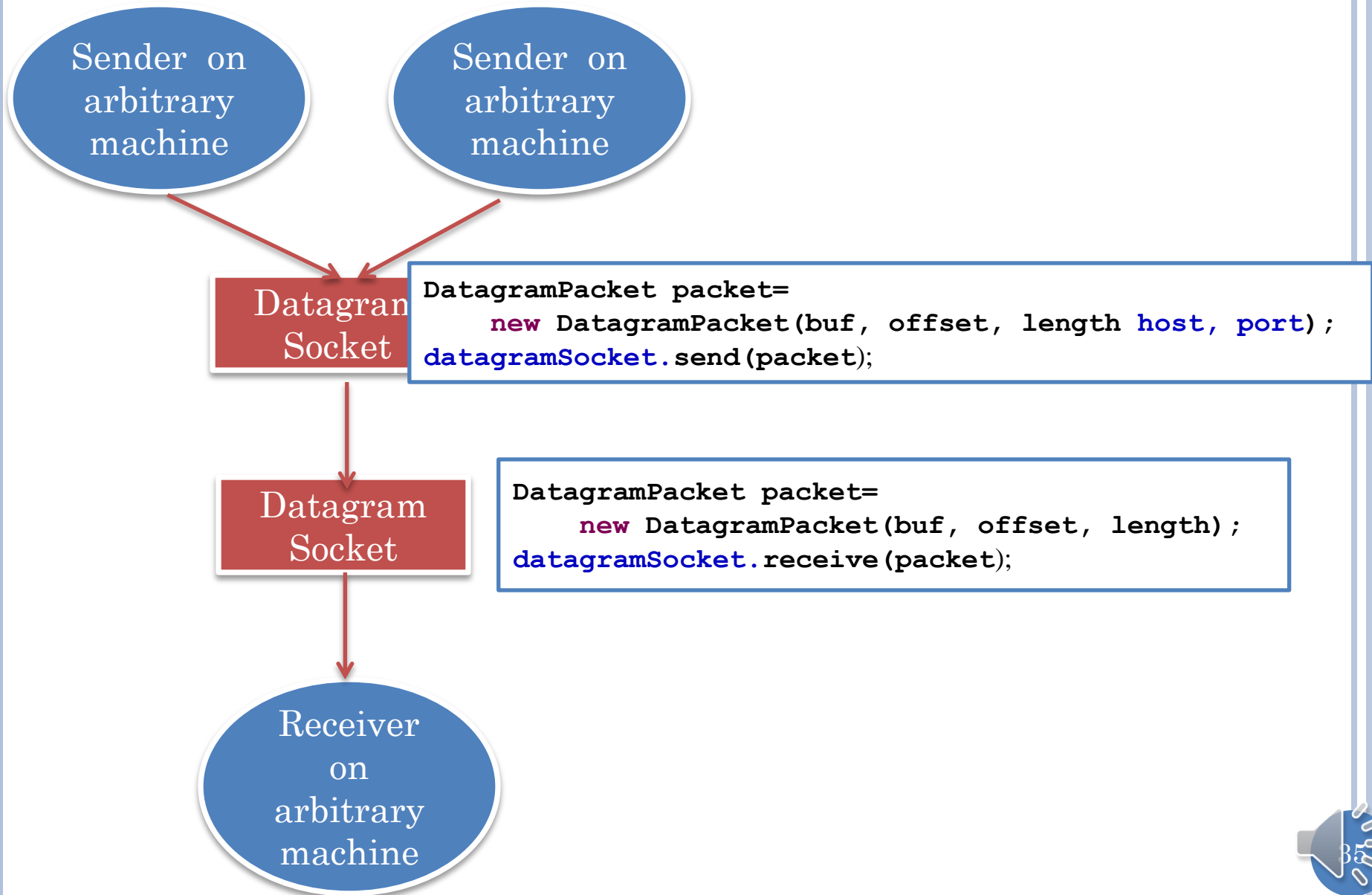


```
DatagramPacket packet= new DatagramPacket(buf, offset,  
length);  
datagramSocket.receive(packet);
```

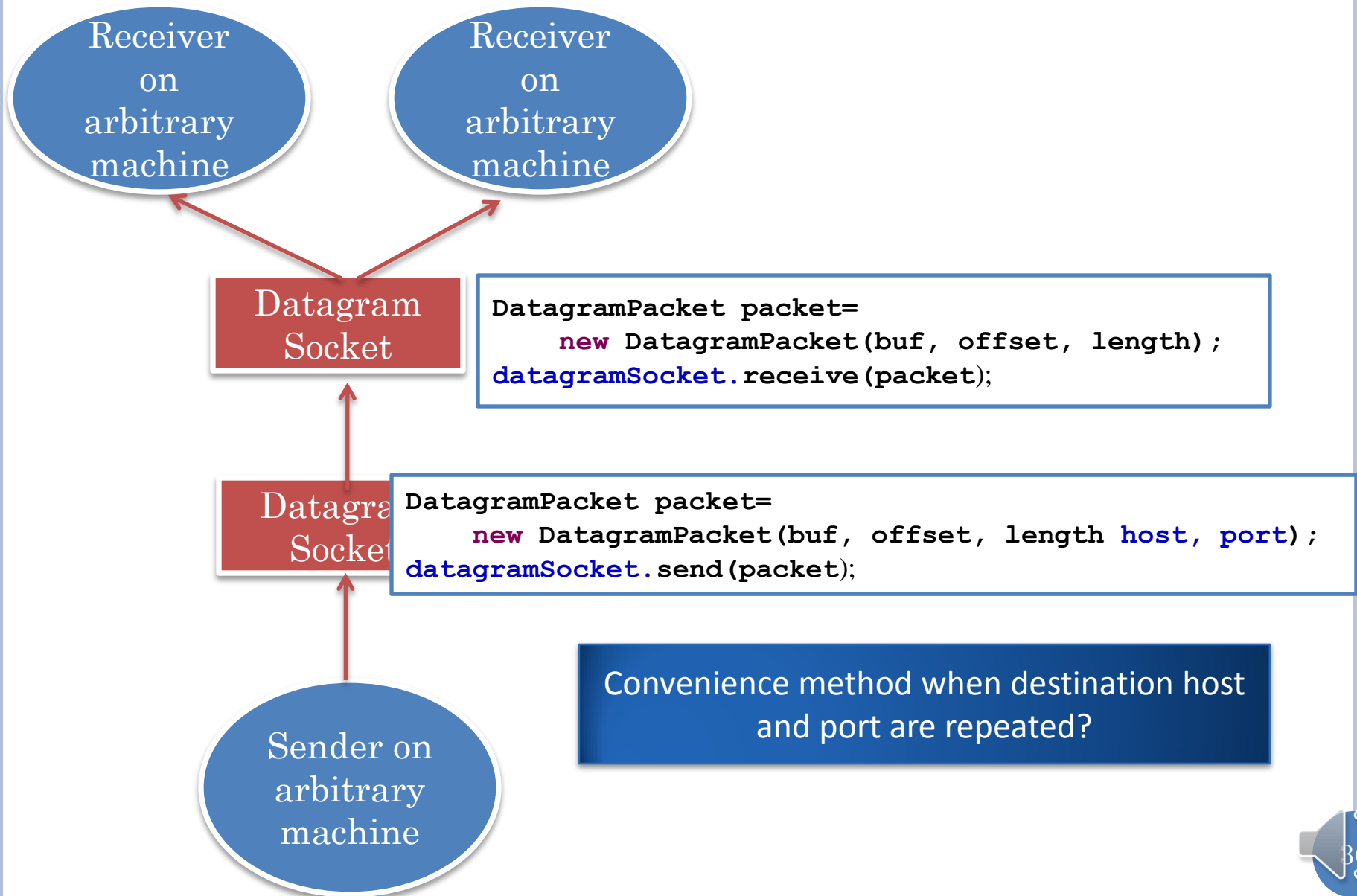
```
DatagramPacket packet= new DatagramPacket(buf, offset,  
length host, port);  
datagramSocket.send(packet);
```



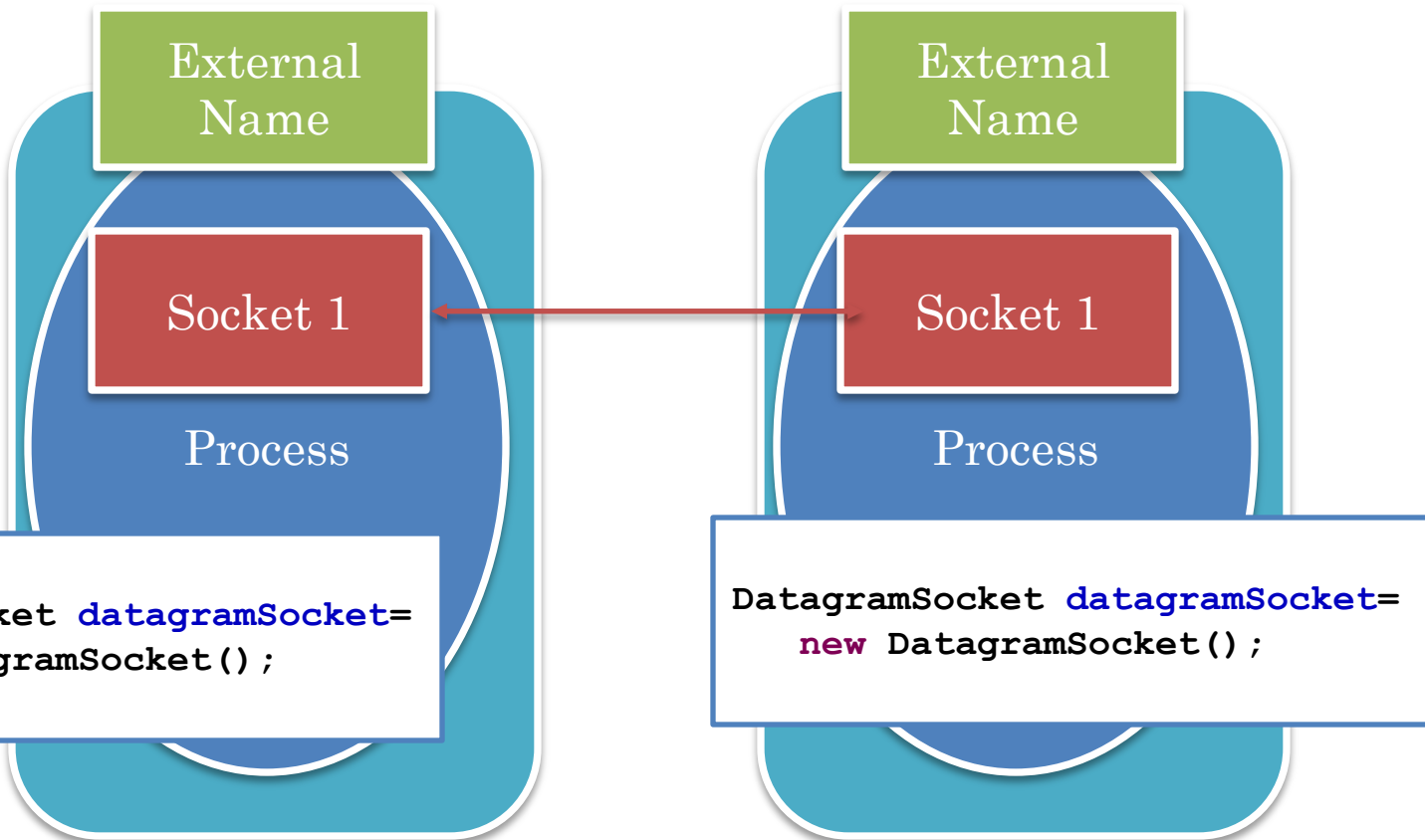
# DATAGRAM SOCKET SHARING



# DATAGRAM SOCKET FOR DATAGRAMS



# DATAGRAMS: DATAGRAM SOCKET AND SPECIAL CALLS

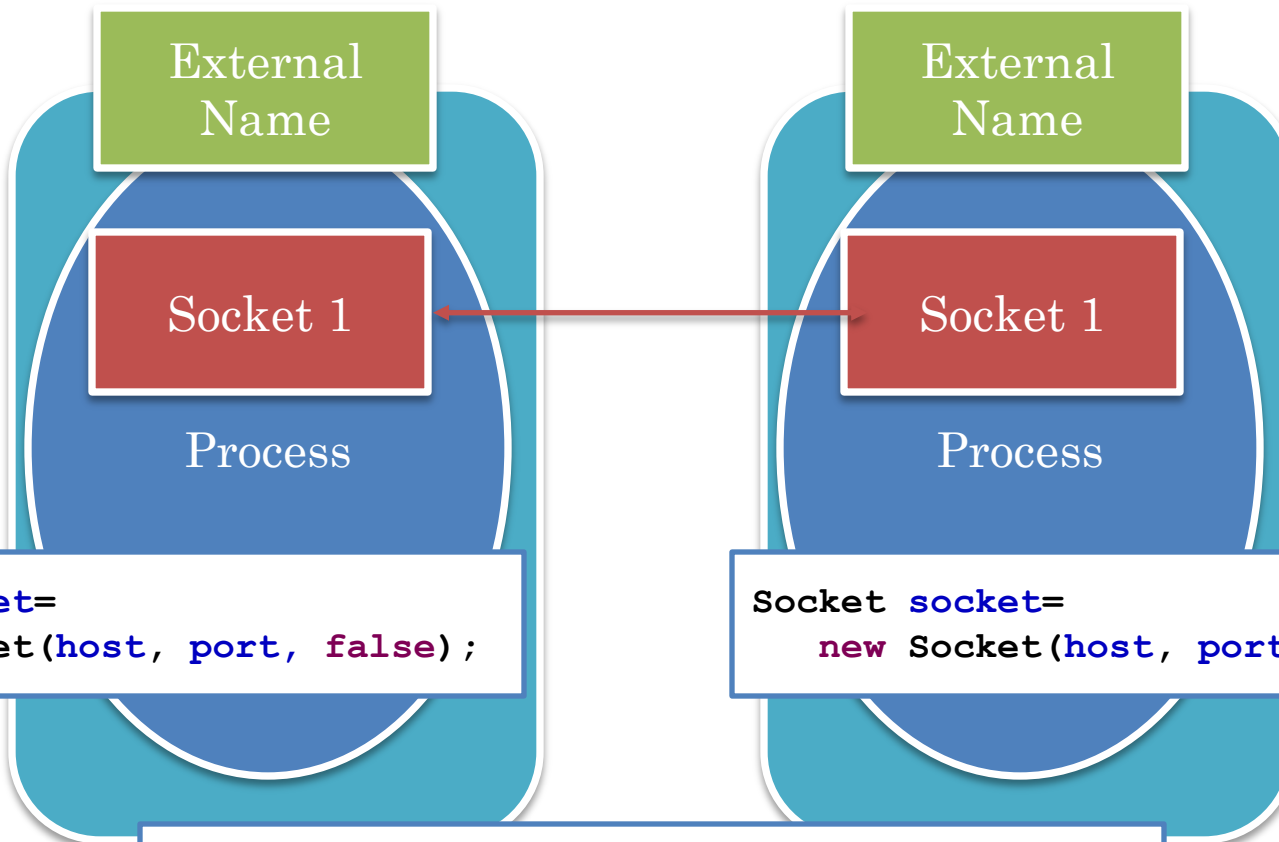


```
DatagramPacket packet= new DatagramPacket(buf, offset,  
length);  
datagramSocket.receive(packet);
```

```
DatagramPacket packet= new DatagramPacket(buf, offset,  
length host, port);  
datagramSocket.send(packet);
```



# DATAGRAMS: SOCKET AND SPECIAL CALLS



```
Socket socket=  
    new Socket(host, port, false);
```

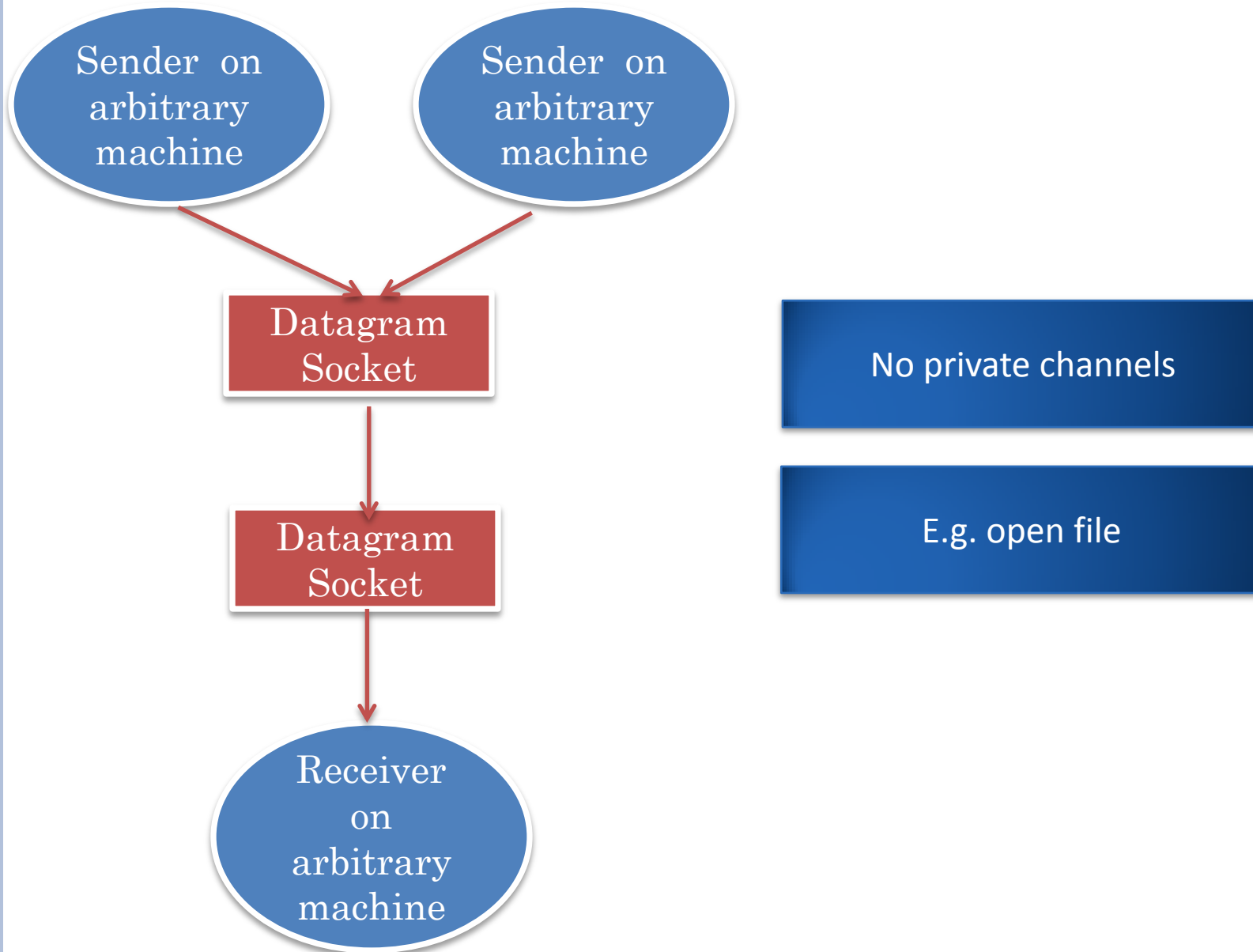
```
Socket socket=  
    new Socket(host, port, false);
```

```
InputStream = socket.getInputStream();  
int retVal =  
    inputStream.read(buf, offset, length);
```

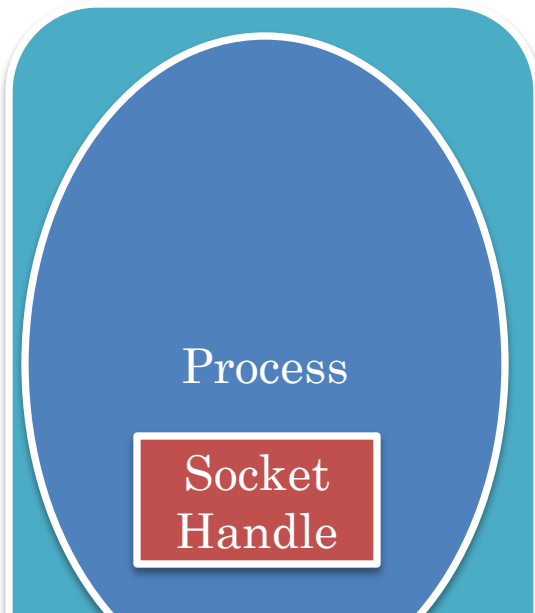
```
OutputStream = socket.getOutputStream();  
OutputStream.write(buf, offset, length);
```



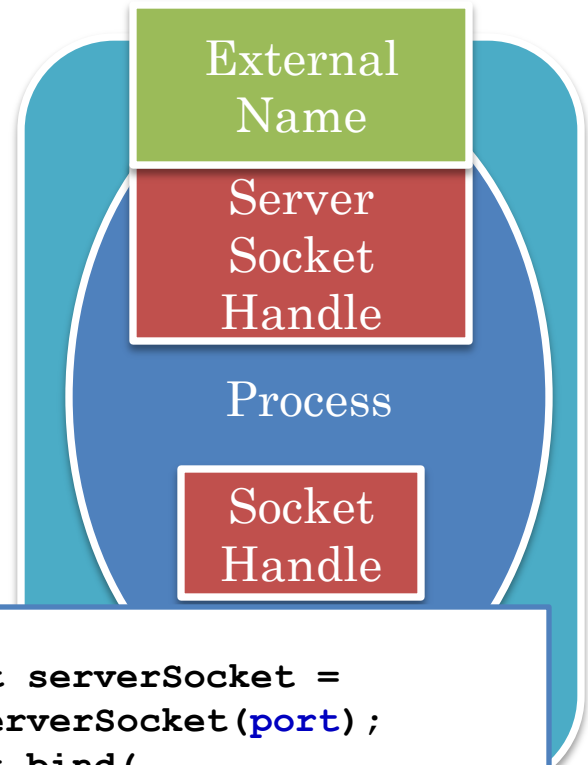
# DATAGRAM SOCKET SHARING



# FROM DATAGRAM TO STREAM SOCKET



```
Socket socket = new Socket();  
socket.connect(  
    new InetSocketAddress(host, port));
```



```
ServerSocket serverSocket =  
    new ServerSocket(port);  
serverSocket.bind(  
    new InetSocketAddress(port));
```

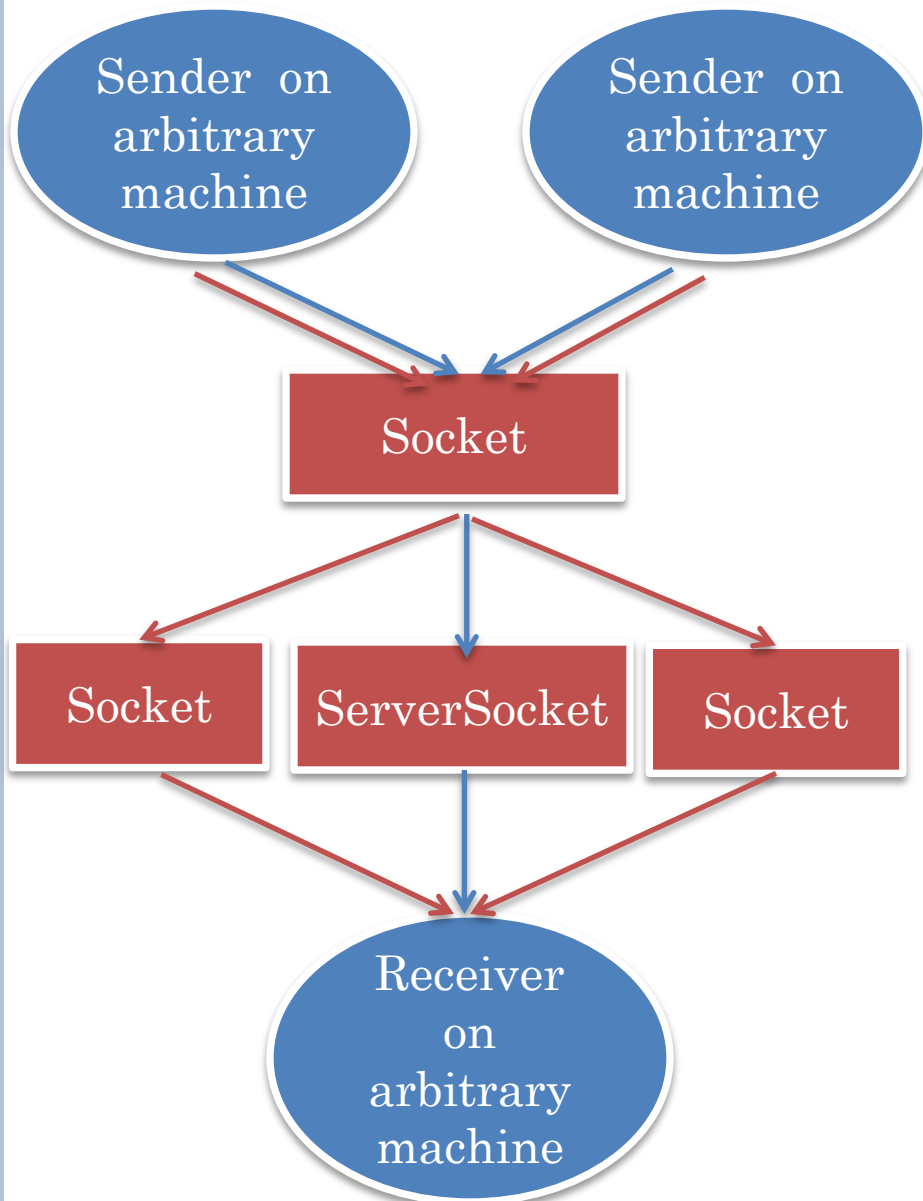
```
Socket socket =  
    serverSocket.accept();
```

```
OutputStream = socket.getOutputStream();  
OutputStream.write(buf, offset, length);
```





# STREAM SOCKET SHARING



Server socket is used to create stream-based socket

Each client connects to it to create a dedicated connection

A data (file) server would create single server socket

“Open” data source operation would connect to server socket

Stream-based socket would represent opened source, which can be read and written



# RELIABILITY AND IN-ORDER?

Datagram sockets: no guarantee,  
built on top of UDP, message size  
limit

Stream sockets: in-order reliable  
on top of TCP/IP

Do not have to change IPC  
mechanism to change guarantee

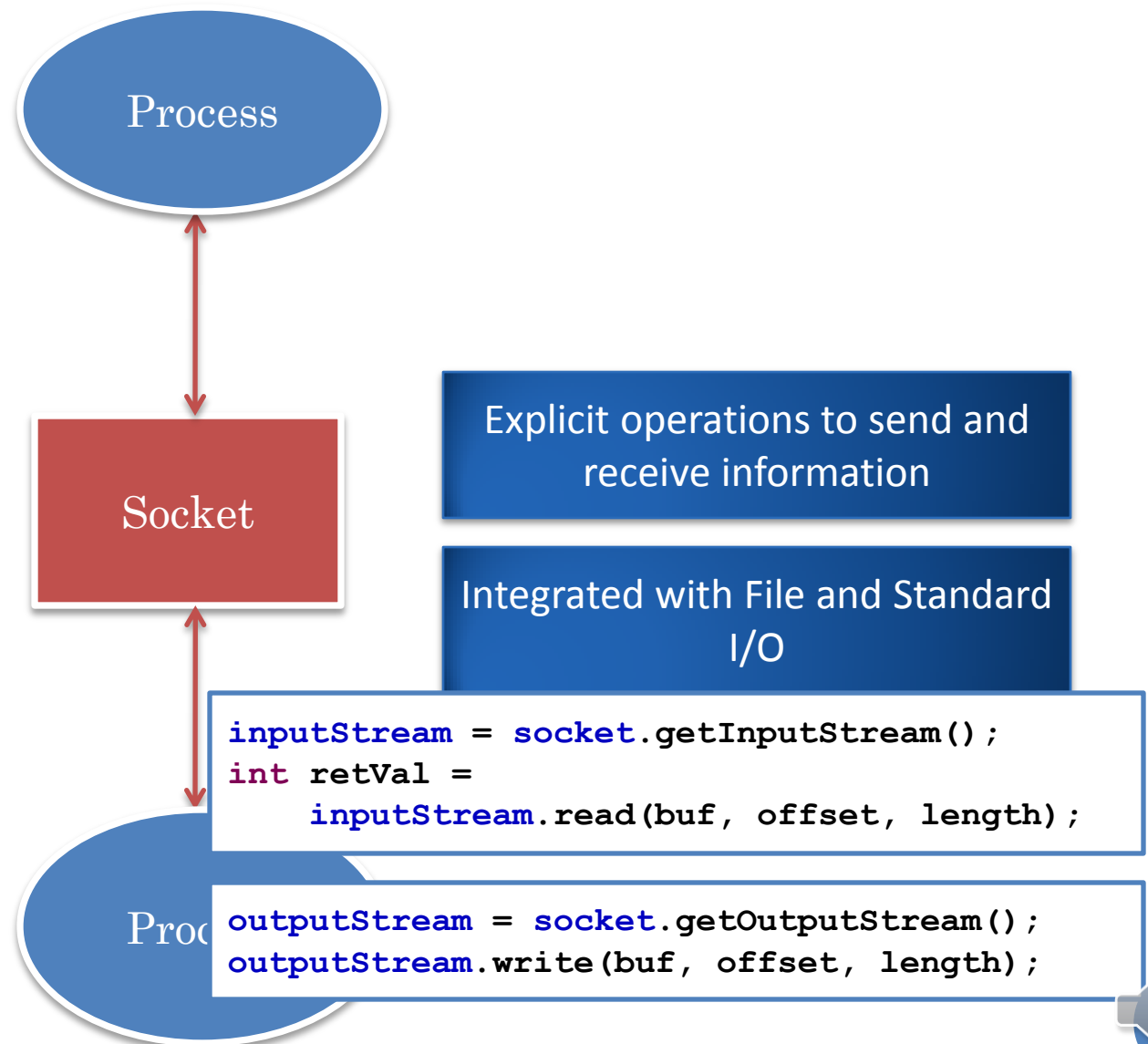
Makes sockets complex, Java  
separating them

```
socket = new Socket(host, port, isStream);
```

```
socket = new DatagramSocket();
```



# READ/WRITE



# IMPLICIT VS NON IMPLICIT

Non connection based ipc can use both implicit and explicit addressing

Connection based ipc uses implicit addressing of destination streams and offset within stream

May also define access rights with which connection was opened



# SEND BLOCKING TIMES?: MESSAGE PIPE LINE

```
OutputStream.write(buf, offset, length);
```

sender buffer

Process  
Thread

Operation started

system buffer

Message in Source  
System Buffer

system buffer

Message in Destination  
System Buffer

Destination  
thread/process starts  
operation

<var>

Process  
Thread

Destination  
thread/process finishes  
operation

Semantics should be like file  
and terminal I/O

In older systems, file I/O  
blocked until data on disk  
(synchronous)

Inefficient, specially if stream  
I/O

Block until in system buffer

If stream socket, then  
message will get through



# RECEIVE BLOCKING TIMES

```
InputStream = socket.getInputStream();  
int retVal =  
    inputStream.read(buf, offset, length);
```

Block until  $\leq \text{length} \geq 1$  bytes received

retVal indicates actual length

Idea is to not block beyond next network message arrival

Give max value so buffer not overwritten

If expecting message of certain size, must loop



# BLOCKING OPERATIONS

```
socket.connect(  
    new InetSocketAddress(host,port));
```

Block until server accepts  
connection to server socket

```
Socket socket =  
    serverSocket.accept();
```

Block until next client tries to  
contact the server socket

```
outputStream = socket.getOutputStream();  
outputStream.write(buf, offset, length);
```

Block until in system buffer

```
inputStream = socket.getInputStream();  
int retVal =  
    inputStream.read(buf, offset, length);
```

Block until  $\leq \text{length} \geq 1$  bytes received



# SOCKET BLOCKING

All operations involve some blocking

What if we want no blocking?

In Java, heavyweight threads can be created

In Unix several primitives for single thread to not block

In Java special NIO layer for blocking and non blocking for sockets (and other I/O resources)





# NIO

NIO (blocking and non blocking)

Sockets (blocking)

Even more flexibility than sockets

How do add non blocking

# BLOCKING OPERATIONS

```
socket.connect(  
    new InetSocketAddress(host,port));
```

Block until server accepts  
connection to server socket

```
Socket socket =  
    serverSocket.accept();
```

Block until next client tries to  
contact the server socket

```
outputStream = socket.getOutputStream();  
outputStream.write(buf, offset, length);
```

Block until in system buffer

```
inputStream = socket.getInputStream();  
int retVal =  
    inputStream.read(buf, offset, length);
```

Block until  $\leq \text{length} \geq 1$  bytes received



# XINU VS ADA

```
int recvclr ()
```

Non blocking, polling, returns either message if it exists, otherwise a special value

```
select
```

```
receive ... <port1> ...
```

```
receive ... <portn> ...
```

```
end
```

Each arm statically registers and interest in an operation on a port, and provides variables and code for completing the operation

Select chooses which of the enabled operations is executed



# COMBINING THE TWO IDEAS

```
int recvclr ()
```

Like rcvclr, no operation will ever block

But will not poll, instead will execute operation that is ready - guaranteed to succeed at least partly

```
select
```

```
receive ... <port1> ...
```

```
receive ... <portn> ...
```

```
end
```

Select-like concept to register interest in receive and other operations

Select will not choose operation to execute, it will tell programmers which operations are ready

Non blocking ready operation such as rcvclr can then be used without blocking



# FORM OF SELECT IN NIO

**select**

**receive** ... **<port<sup>1</sup>>** ...

**receive** ... **<port<sup>n</sup>>** ...

**end**

Objects for arm and selector

Interest in a port and operation, and execution of operation decoupled

Can dynamically register with a selector interest in an operation (e.g. receive, , send, accept) on a resource (e.g. port, file)

Selector blocks and on unblocking tells its user which of the interested (operation, resource )pairs are enabled and thus ready for execution



# SELECTOR

Resource on which async operation may be executed (replaces Socket, ServerSocket, File)

Selectable Channel

```
SelectionKey register(Selector s, int ops )
```

```
configureBlocking(boolean)
```

Single selector for all operations

```
Selector selector = Selector.open();
```

Selector

```
int select()
```

```
Set<SelectionKey> selectedKeys()
```

```
Selector wakeup()
```

Selection Key

```
SelectableChannel channel ()
```

Registers interest in (resource, operation) pair referenced by a key id

```
SelectionKey key =  
selectableChannel.register( selector,  
SelectionKey.OP_ACCEPT );
```

Blocking call waiting until at least one registered pair is enabled by some event, and returning # of such pairs

Keys of actual enabled operations

Unblock select, usually after a new registration

Key to selectable channel, which can be used to execute enabled operation immediately



# SERVERSOCKET AND SOCKET CHANNELS

Selectable  
Channel

IS-A

IS-A

Server  
Socket  
Channel

ServerSocket socket()

SocketChannel accept()

(finish)connect()

Socket  
Channel

Socket socket()

int read(ByteBuffer b)

int write(ByteBuffer b)

HAS-A

Why channel I/O operations?

Operations at channel level are  
non blocking

HAS-A

Server  
Socket

Socket accept()

Socket

InputStream  
getInputStream()

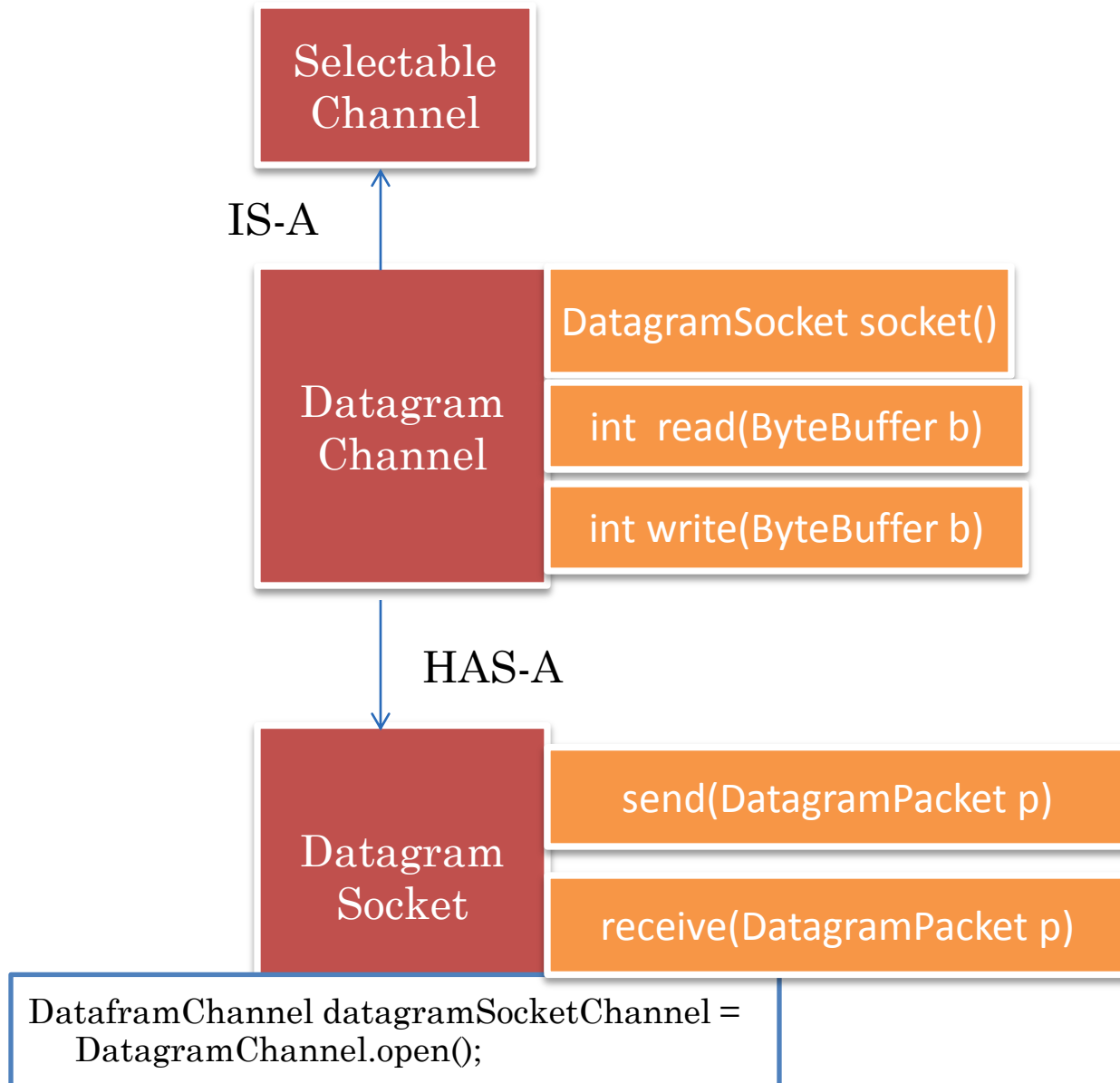
OutputStream  
getOutputStream()

```
ServerSocketChannel serverSocketChannel =  
    ServerSocketChannel.open();
```

```
SocketChannel socketChannel =  
    SocketChannel.open();
```



# DATAGRAM CHANNEL





# CHANNEL VS. STREAM I/O

```
InputStream = socket.getInputStream();  
int retVal = inputStream.read(buf, offset, length);
```

```
OutputStream = socket.getOutputStream();  
ostream.write(buf, offset, length);
```

```
int retVal = socketChannel.read(byteBuffer);
```

```
int retVal = socketChannel.write(byteBuffer)
```

ByteBuffer, like packet, encapsulates buf, offset, and length

Write writes only as many bytes as available in source buffer when in async mode

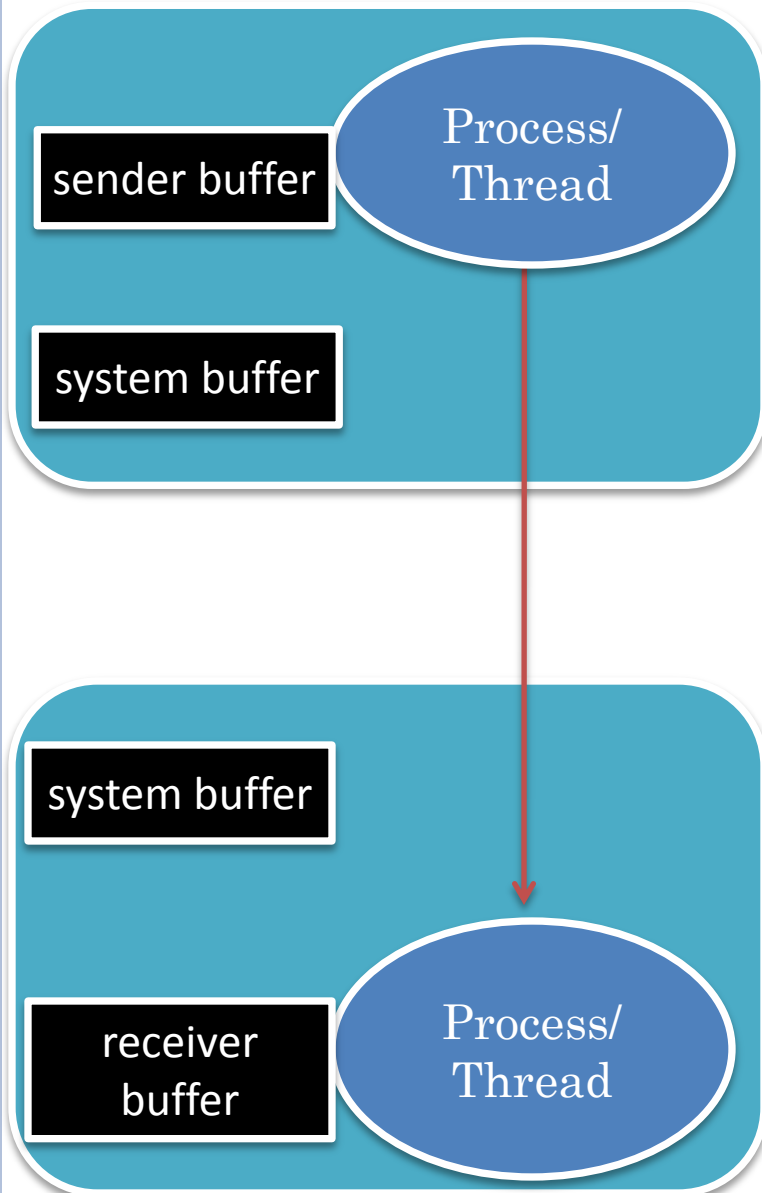
Channel unlike stream and like Unix file/socket can be read and written

IPC mechanism may not complete operation and same buffer may be used for multiple batch operations

System can use the buffer directly instead of creating own source or destination buffer



# DIRECT VS. NON DIRECT



Direct Buffer: System tries to use sender and receiver buffer directly without creating intermediate source or system non direct buffer

Buffer copying is an expensive operation

In synchronous sends safe to avoid copying. In asynchronous, requires careful programming

Direct buffer allocation from kernel space so more costly

Use direct buffer only when performance is an issue and buffer is long lived



# ALLOCATING DIRECT VS. NON DIRECT

Direct

```
ByteBuffer ByteBuffer.allocateDirect(capacity)
```

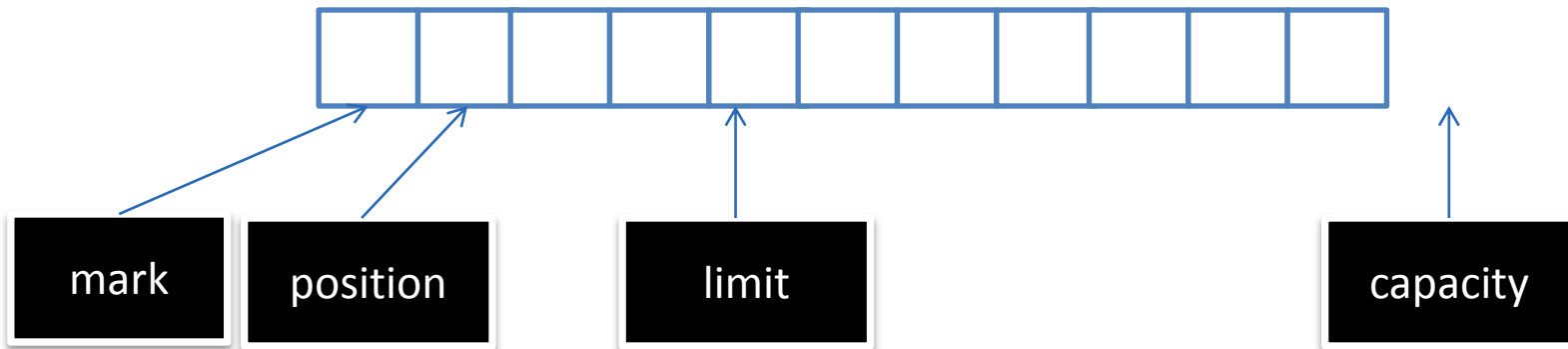
Indirect

```
ByteBuffer ByteBuffer.allocate(capacity)
```

```
ByteBuffer ByteBuffer.wrap(bytes[])
```



# DATA STRUCTURE



Storage for contents

Size of (available) contents

Position of next element to be read or written

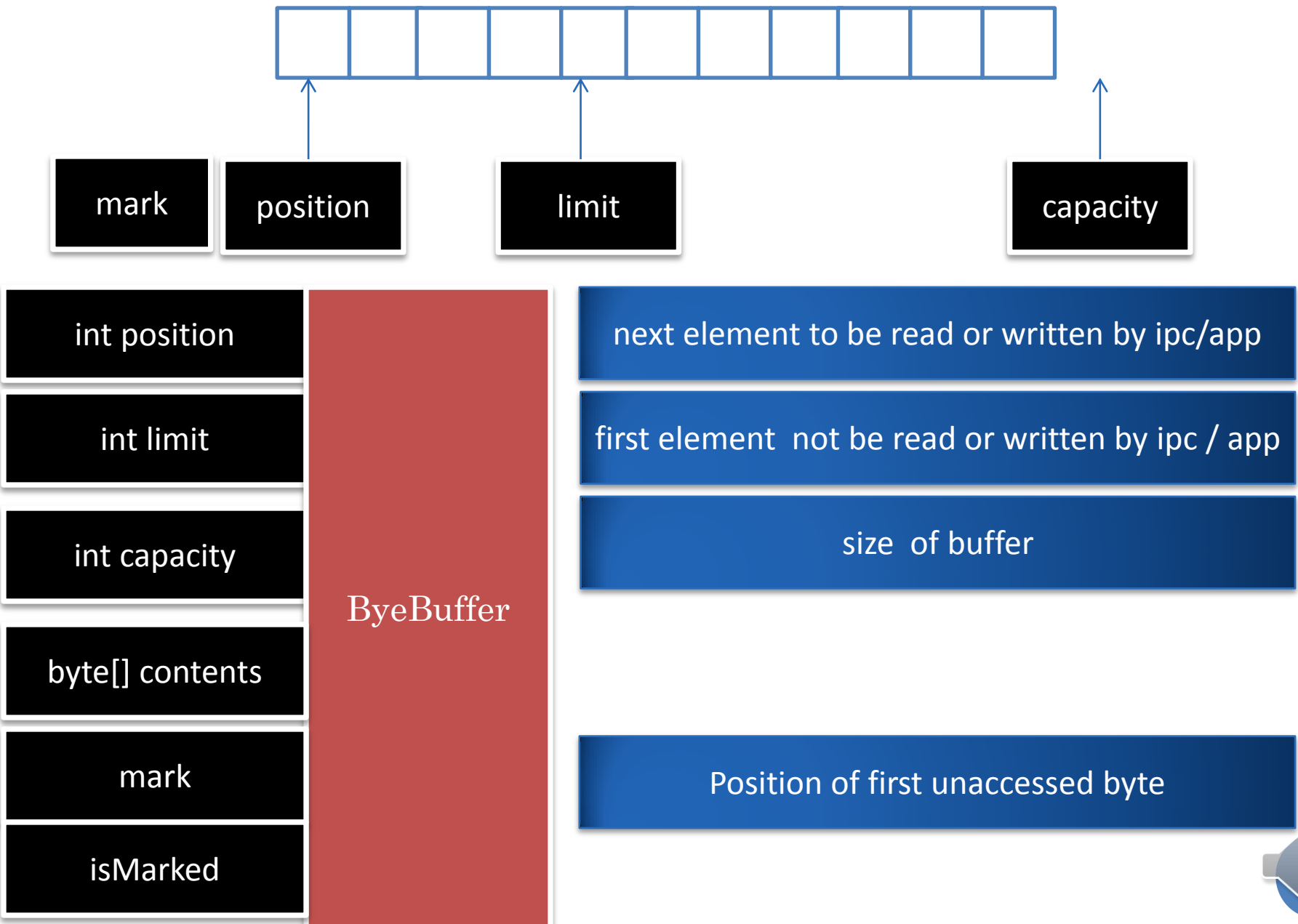
Read may not yield expected bytes, write may not empty all bytes

May use the same buffer for multiple serial operations  
or batch operations, need to mark position of first unconsumed byte

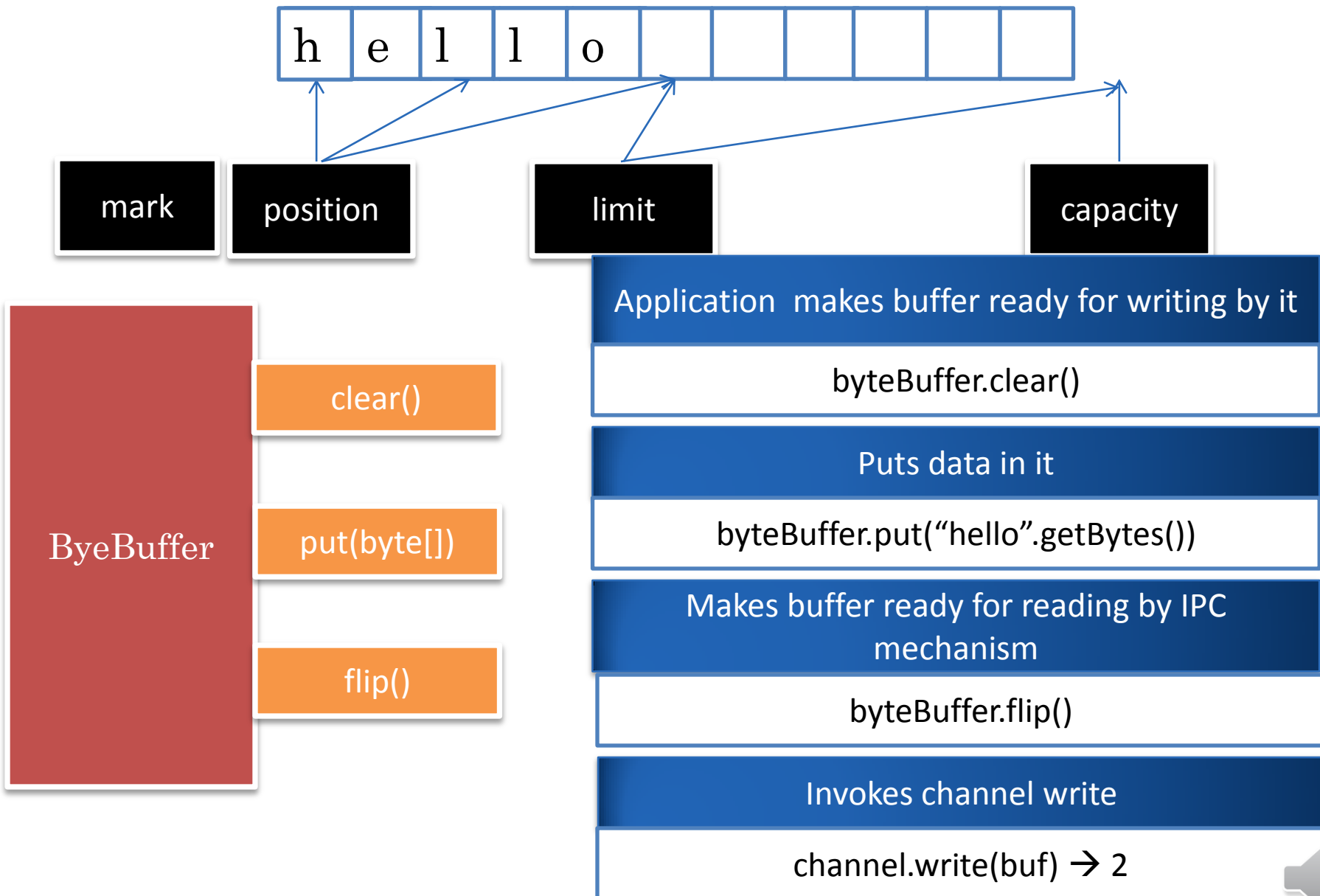
$\text{mark} \leq \text{position} \leq \text{limit} \leq \text{capacity}$



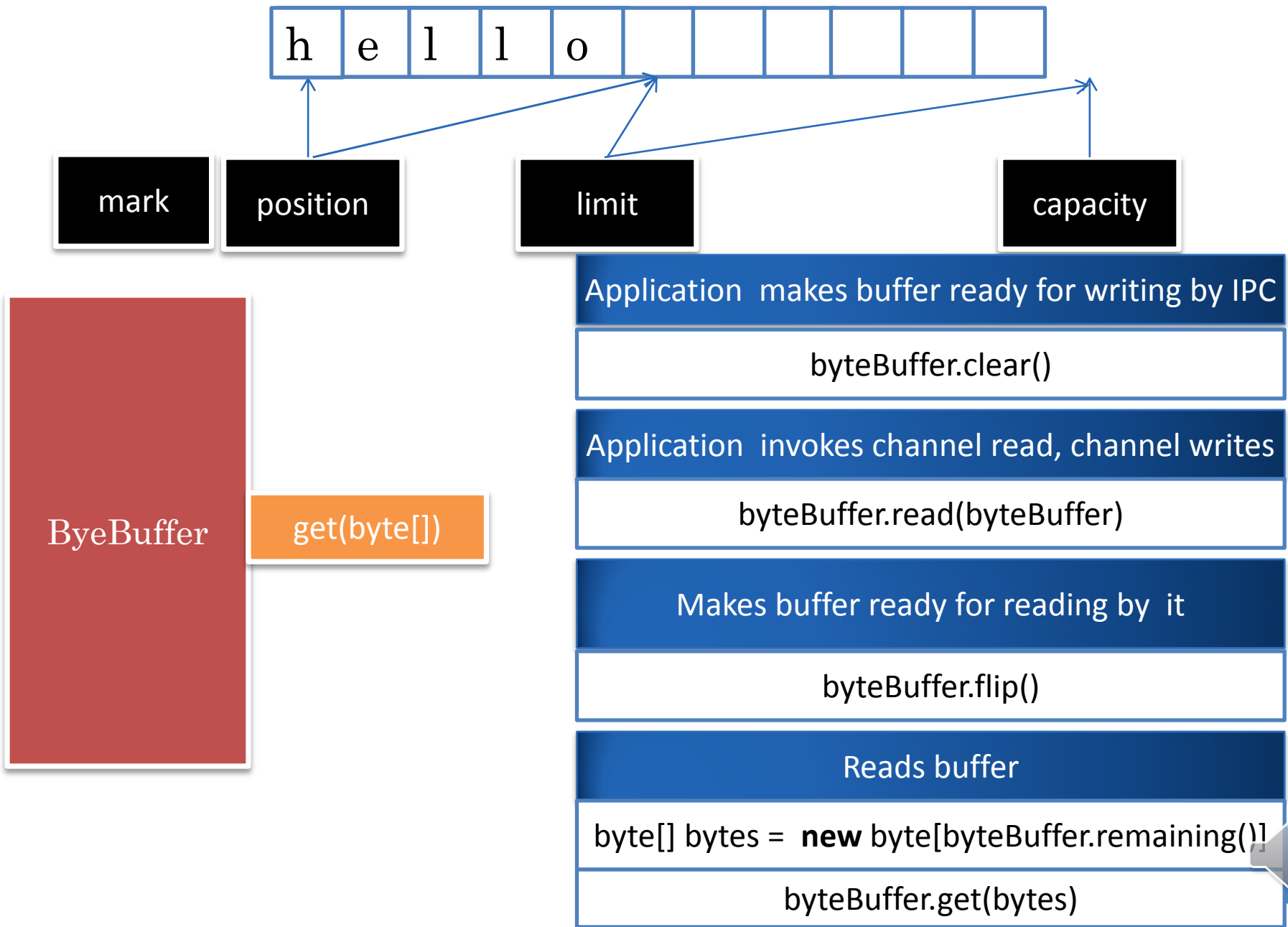
# DATA STRUCTURE



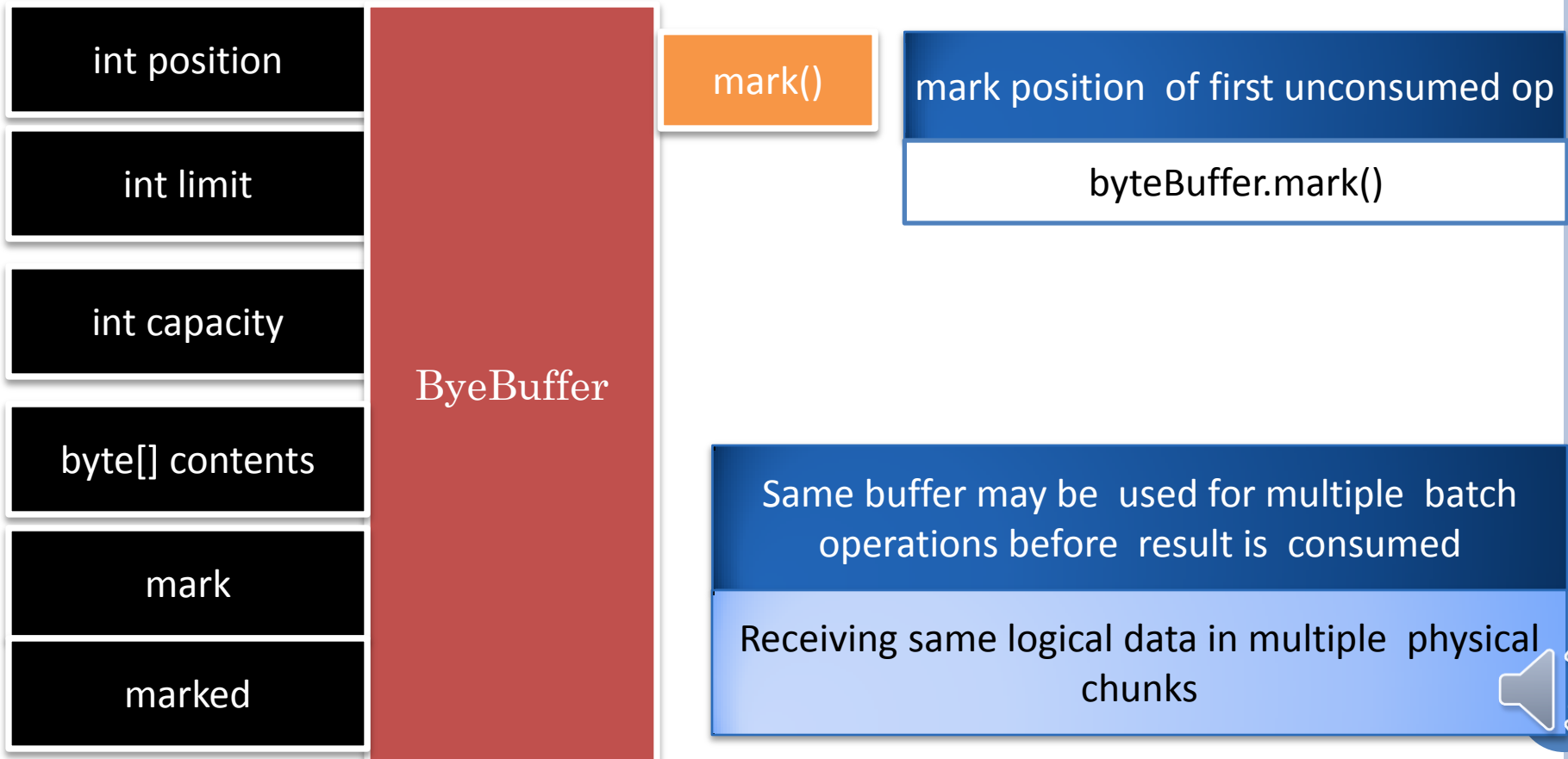
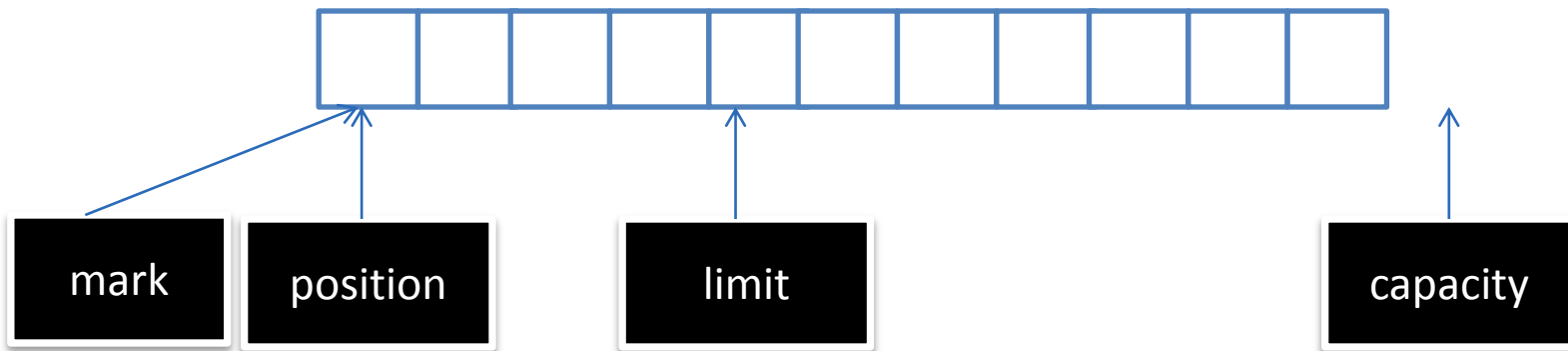
# WRITING



# READING



# MARKING





# OPERATIONS

## ByteBuffer

Int remaining()

flip()

clear()

rewind()

reset()

mark()

put(bytes[] b)

get(bytes[] b)

position(int p)

limit(int l)

limit – position

limit = position, position = 0, marked = false

position = 0, limit = capacity, marked = false

position = 0, marked = false

position = mark, if marked

mark = position, marked = true

For each byte, write next byte at pos, pos++

For each byte, read next byte at pos, and pos++

position = p

limit = l



# NIO EVALUATION

More flexible than even Java sockets

Hence more complex

Even the normal case requires tutorials  
describing normal patterns of user

Selection must be done in a single complex  
thread

