

Todo list

cite ALL images	8
intros for background section	8
ref this	9
ref this	11
is this section right	12
complete this with help from links	14
any more middleboxes - packet cap?	15
Examples of hardware comparisons	23
better word for derive	23
why I need benchmarking	23
show example of optimising java and how different it looks, also different depending on architecture	24
better phrasing off that sentence needed	24
More here	24
get paper	24
explain difference with GC etc	24
get data	24
ref manual	25
cite this	26
should code use dpdk malloc?	26
cite this	28
cite this	29
article on this	30
reasons why JNI is slow	30
ref this	30
ref this	31
expand on this	31

Is this relevant	31
expand on where data is allocated between native and java - use links above	32
Diagram of packets from NIC using c through 'technique' and then processing packets in java and then back	32
ref this	32
do this	32
do this	32
cant use bytebuffers directly with DPDK as dont use hugepages	33
ref this	34
Explain why it has faster performance	35
Nested padding in struct?	35
C struct field always in given order	35
Inconsistencies with datatype length so using uint32t etc	35
Mention this is on 64-bit machine and obviously you don't notice padding and order of elements can pay an important part in this	35
Proof on speed	36
evaluate each technique on ease and number of data copies and results and say why we picked 1 of the others	38
stats for this	38
ref this	39
should these graphs all be together because they look the same	40
cant be done and reference earlier	41
extension to remove packet pointer passing and do it directly from struct array - would have to dynamically pack array though - how????	41
do hardware testing - ie hardware faster than programs (i think) so should aim for through- put and not efficiency	41
ref this	42
ref this	43
get a table of this	44

mention applications below for demonstrations of ease of use of framework, also mention only subset of availability of dpdk, and mention further extensions easy to do like packet fragmentation etc	45
important to check for errors at every stage to stop jvm crashes	45
ref this and say thread can be set to run on multiple cores - need to fix code for this . . .	46
example?	46
Describe 3 objects and fix code to use more than 1	46
ref this	47
elaborate on this, diagram? how did pointers work?	47
ref this	47
get image of setup of memory location, ie number then list of packet pointers	48
ref this	48
put image of gui, also write output to files	49
mention machine specs	50
Talk about pktgen module in kernel - does dpdk version use this	50
how does pktgen work and why we used it, only does 32 packet bursts or less	50
which ones?	51
ref this and ubuntu	51
Figure: Graph of c packet capture, 1 core, different frame sizes, will need to reference transmit speed of pktgen	51
what does the results show	52
Figure: Graph of java packet capture, 1 core, different frame sizes, will need to reference transmit speed of pktgen	52
expand on this	52
Figure: Graph for improvements	53
explain results compared to previous ones and c version	53
Figure: more profiling?	53
image of this	54
Figure: Graph for improvements	54
continue this, where should it go?	56
should I be putting all data in appendix	57

Things to talk about

- Is app portable - it uses unints so why not?, packed structs a problem?
- spell check in excalibur

IMPERIAL COLLEGE LONDON

INDIVIDUAL PROJECT

COMPUTING - BENG

Packet Processing in Java

ASHLEY HEMINGWAY

Supervisor:
Peter PIETZUCH

2nd Marker:
Wayne LUK

June 2015

Abstract

Acknowledgements

Contents

1	Introduction	6
1.1	Motivation	6
1.2	Objectives	6
1.3	Solution Idea	7
1.4	Potential Applications	7
2	Background	8
2.1	Network Components	8
2.1.1	Models	8
2.1.2	Network Packets	10
2.1.3	Packet Handling	12
2.1.4	Network Interface Controller (NIC)	13
2.1.5	Middleboxes	13
2.2	Java	15
2.2.1	JVM	15
2.2.2	Java Native Interface (JNI)	17
2.2.3	Current Java Networking Methods	20
2.2.3.1	Java Sockets	20
2.2.3.2	Remote Method Invocation (RMI)	20
2.2.3.3	Shared Memory Programming	20
2.3	Related Works	20
2.3.1	jVerbs	20
2.3.2	Native I/O API's	22
2.3.3	Packet Shader, infiniband - any more	23
2.3.4	Data Plane Development Kit (DPDK)	23
2.4	Performance Testing Techniques	23
2.4.1	Programming Languages	23

2.4.1.1	Loop Recognition	24
2.4.1.2	Benchmark Game	24
2.4.1.3	Using Economics	24
2.4.1.4	Which is better?	24
2.4.2	Intra-Language Techniques	24
2.4.3	Applications	24
3	DPDK	25
3.1	Environment Abstraction Layer (EAL)	25
3.2	Logical Cores	25
3.3	Huge Pages	25
3.4	Ring Buffer	26
3.5	Memory Usage	26
3.5.1	Allocation	27
3.5.2	Pools	27
3.5.3	Message Buffers	27
3.5.4	NUMA	28
3.6	Poll Mode Driver (PMD)	29
3.7	Models	29
4	Initial language comparison	30
4.1	Benchmarking Algorithm	30
4.2	Results	31
4.3	Further Investigation	31
5	Design Considerations	32
5.1	Data Sharing	32
5.1.1	Objects and JNI	32
5.1.2	ByteBuffers	33

5.1.3	Java Unsafe	33
5.2	Packing Structures	34
5.3	Performance testing	37
5.3.1	Expectations	38
5.3.2	Results	38
5.3.3	Evaluation	41
5.4	Thread affinity	41
5.5	Endianness	43
5.6	Data type conversion	43
5.7	Protocol undertaking	44
5.8	config files	44
6	DPDK Java	45
6.1	Overview (better name?)	45
6.2	Native Libraries	45
6.2.1	Library Compilation Tools	45
6.3	Initialisation	46
6.4	Processing Threads	46
6.5	Packet Data Handling	47
6.6	Packet Polling	48
6.7	Packet Sending	48
6.8	Statistic Profiling	49
6.9	Testing	49
7	Applications	49
7.1	Firewall	49
7.2	Network Address Translator (NAT)	49

8	Evaluation	50
8.1	Packet Generating	50
8.2	Initial Testing	51
8.2.1	Set-up	51
8.2.2	Methods	51
8.2.3	Results	51
8.3	Further Testing	51
8.3.1	Set-up	54
8.3.2	Methods	54
8.3.3	Results	54
8.4	Software Design	54
8.4.1	Portability	55
8.5	Possible Improvement	55
9	Conclusion	56
9.1	Future Work	56
9.2	extrapolate for other non-native languages	56
10	User Guide	57

1 Introduction

1.1 Motivation

As modern computing techniques advance, people are trying to find more generic solutions to problems which have been solved by native applications in the past. A main area of focus has been network middleboxes which are developed to manipulate network packets. Common examples of middleboxes are firewalls, network address translators (NATs) and load balancers, all of which inspect or transform network packets in the middle of a connection between a public and private network. In recent years, people have been developing a number of programmable middleboxes which allow these generic solutions to be used on a wide scale basis.

As middleboxes are mainly used for networking purposes, they are required to process network packets at line rate (i.e. at speeds which allow packets to be processed as they are received) which generally is 10Gbps (gigabits per second), but speeds can reach 40Gbps and even 100Gbps. This requires the application to retrieve the packet from the network line, inspect and transform the packet in the desired way and then insert the packet back onto the network line, all within a time period sufficient enough to not cause a backlog. High performance implementations of such applications are available, but are written in native languages, predominately in C/C++. However, more and more high performance computing projects are being developed in Java and have succeeded in performing at similar speeds to C/C++ applications. These high performance projects typically utilise a distributed system, so it makes sense to have the middleboxes written in Java as well for easy scaling.

The main challenge is actually getting the I/O system for the Java application to run at line rate speeds, due to challenges with how the JVM (Java Virtual Machine) interacts with memory, computer's kernel and the network interface controller/card (NIC). Once this challenge has been overcome, there are no reasons why programmable middleboxes written in non native languages such as Java can't exist within networking systems.

1.2 Objectives

The main objective for this project is to research, develop and test a new application which is implemented in a non-native language such as Java, but can process packets at a high performance level. This also requires that the application can perform I/O operations at line rate in order to pass on data packets without reducing the line rate. The aim is to match current native applications which are generally written in C/C++.

With the main objective stated above, I set out a few initial, smaller objectives in order to divide the project into more manageable objectives:

- Understand similar applications and API's written in C/C++ which process at line speeds and how the implementations can be exploited for Java applications
- Conduct research into Java techniques which can be used in order to increase I/O operations
- Implement basic middlebox applications in Java such as a firewall and a NAT
- Compare Java implementations to those which are pure Java and pure C/C++

1.3 Solution Idea

The idea is to implement a new library in Java using a few techniques in order to speed up network access. Firstly, no networking will be done via the JVM and the operating system kernel. Instead, the Java application will bypass the kernel altogether using a combination of direct memory access and high speed I/O operations via a C library which we interact with the application via the Java Native Interface (JNI).

This eliminates the need for the JVM to interact with the kernel via system calls in order to do the network access, which can be relatively slow compared to direct network access achievable from native applications.

1.4 Potential Applications

2 Background

This section focusses on a number of background concepts which will help to gain further understanding of networking practices and protocols which the solution must adhere to. It then gives an overview of the architecture of the Java language and explores on some of the key techniques which will be exploited. As this report focusses on performance testing throughout, this section will then cover techniques to test custom programs. It finally looks at related works within the research area and existing frameworks which could be exploited.

2.1 Network Components

2.1.1 Models

Generally there are 2 well known network models. The Open System Interconnection (OSI) model in Figure 15 represents an ideal to which all network communication should adhere to while the Transmission Control Protocol/Inter Protocol (TCP/IP) model represents reality in the world. The TCP/IP model combines multiple OSI layers into 1 TCP/IP layer simply because not all systems will go through these exact stages depending on the system implementation.

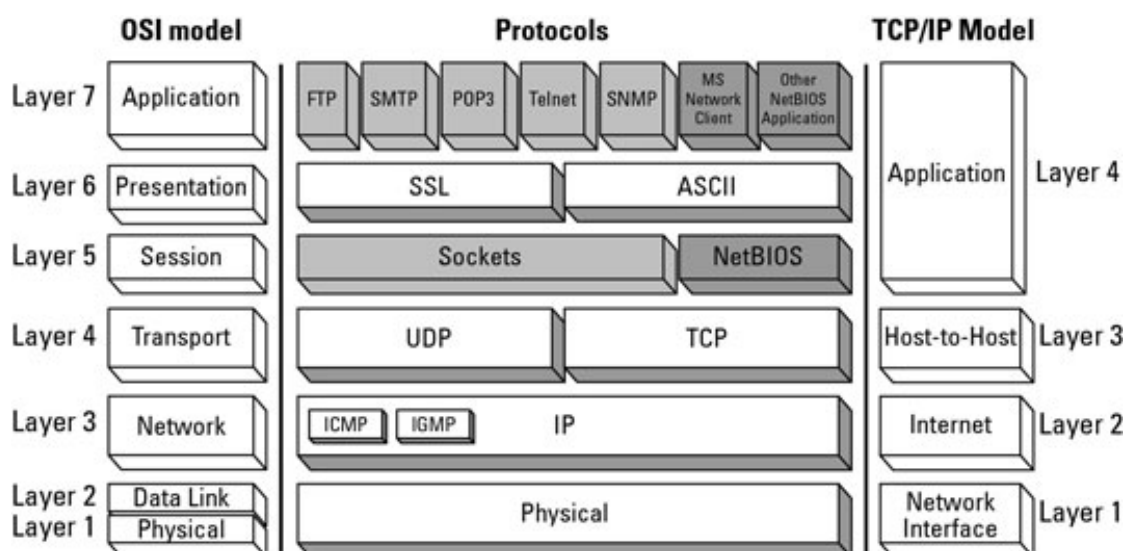


Figure 1: OSI vs TCP/IP Model

The TCP/IP application layer doesn't represent an actual application but instead it's a set of protocols which provides services to applications. These services include HTTP, FTP, SMTP,

POP3 and more. It acts as an interface for software to communicate between systems (e.g. client retrieving data from server via SMTP).

The transport layer is responsible for fragmenting the data into transmission control protocol (TCP) or user datagram protocol (UDP) packets, although other protocols can be used. This layer will attach its own TCP or UDP header to the data which contains information such as source and destination ports, sequence number and acknowledgement data.

The network/internet layer attaches a protocol header for packet addressing and routing. Most commonly this will be an IPv4 or IPv6 header . This layer only provides datagram networking functionality and it's up to the transport layer to handle the packets correctly.

ref this

The network interface or link layer will firstly attach its own ethernet header (or suitable protocol header) to the packets, along with an ethernet trailer. This header will specify the destination and source of the media access control (MAC) address which are specific to network interfaces. The next step is to put the packet onto the physical layer, which may be fibre optic, wireless or standard cables.

This will eventually build a packet of data which include the original raw data along with multiple headers for each layer of the model (Figure 2).

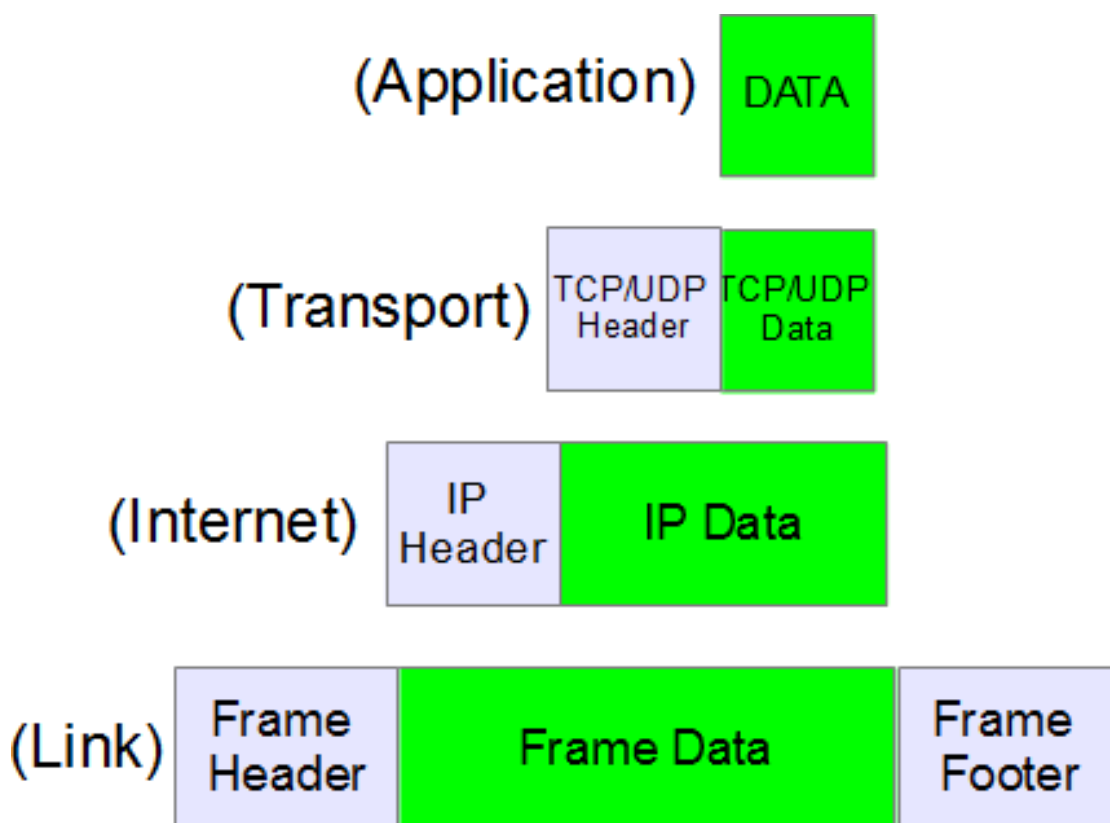


Figure 2: Network model building up of data headers

2.1.2 Network Packets

A network packet is responsible for carrying data from a source to a destination. Packets are routed, fragmented and dropped via information stored within the packet's header. Note: in this report packets and datagrams are interchangeable. Data within the packets are generally input from the application layer, and headers are appended to the front of this data depending on the network level described in Section 2.1.1. Packets are routed to their destination based on a combination of an IP address and MAC address which corresponds to a specific computer located within the network, whether that is a public or private network. In this project we will only be concerned with the Internet Protocol (IP) and therefore IPv4 (Figure 3) and IPv6 (Figure 4) packet headers.

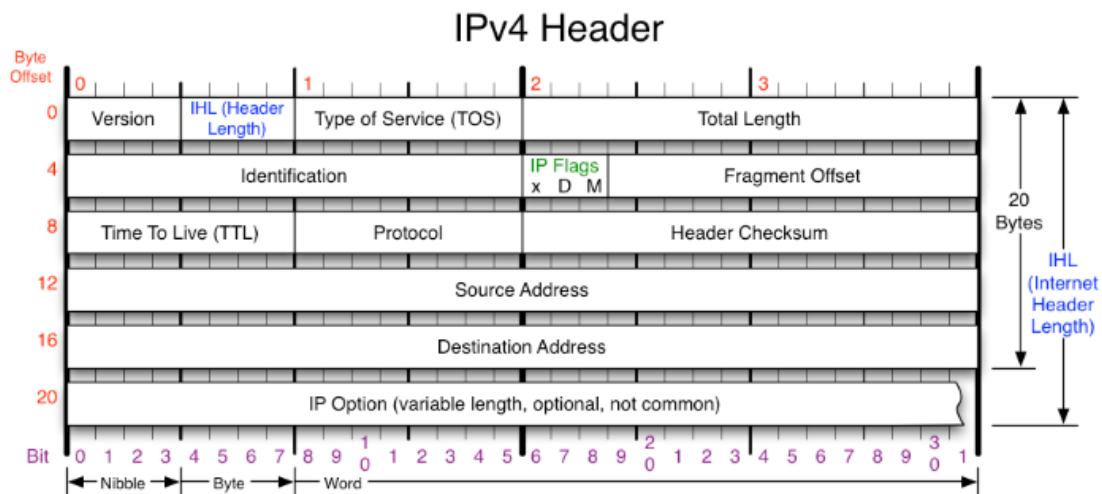


Figure 3: IPv4 Packet Header

- Version - IP version number (set to 4 for IPv4)
- Internet Header Length (IHL) - Specifies the size of the header since a IPv4 header can be of varying length
- Type of Service (TOS) - As of RFC 2474 redefined to be differentiated services code point (DSCP) which is used by real time data streaming services like voice over IP (VoIP) and explicit congestion notification (ECN) which allows end-to-end notification of network congestion without dropping packets
- Total Length - Defines the entire packet size (header + data) in bytes. Min length is 20 bytes and max length is 65,535 bytes, although datagrams may be fragmented.
- Identification - Used for uniquely identifying the group of fragments of a single IP datagram
- X Flag - Reserved, must be zero
- DF Flag - If set, and fragmentation is required to route the packet, then the packet will be dropped. Usually occurs when packet destination doesn't have enough resources to handle incoming packet.

- MF Flag - If packet isn't fragmented, flag is clear. If packet is fragmented and datagram isn't the last fragment of the packet, the flag is set.
- Fragment Offset - Specifies the offset of a particular fragment relative to the beginning of the original unfragmented IP datagram
- Time To Live (TTL) - Limits the datagrams lifetime specified in seconds. In reality, this is actually the hop count which is decremented each time the datagram is routed. This helps to stop circular routing.
- Protocol - Defines the protocol used the data of the datagram
- Header Checksum - Used for to check for errors in the header. Router calculates checksum and compares to this value, discarding if they don't match.
- Source Address - Sender of the packet
- Destination Address - Receiver of the packet
- Options - specifies a number of options which are applicable for each datagram. As this project doesn't concern these it won't be discussed further.

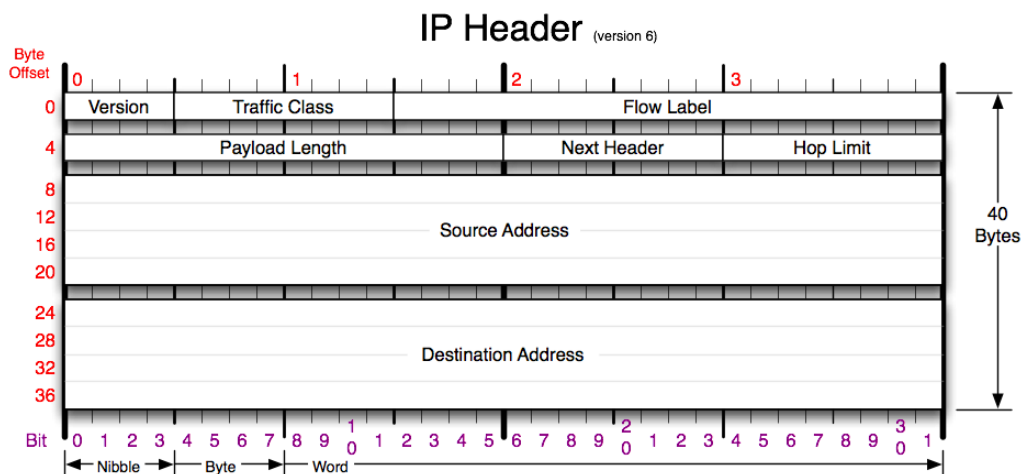


Figure 4: IPv6 Packet Header [?]

- Version - IP version number (set to 6 for IPv6)
- Traffic Class - Used for differentiated services to classify packets and ECN as described in IPv4.
- Flow Label - Used by real-time applications and when set to a non-zero value, it informs routers and switches that these packets should stay on the same path (if multiple paths are available). This is to ensure the packets arrive at destination in the correct order, although other methods for this are available.
- Payload Length - Length of payload following the IPv6 header, including any extension headers. Set to zero when using jumbo payloads for hop-by-hop extensions.

- Next Header - Specifies the transport layer protocol used by packet's payload.
- Hop Limit - Replacement of TTL from IPv4 and uses a hop value decreased by 1 on every hop. When value reaches 0 the packet is discarded.
- Source Address - IPv6 address of sending node
- Destination Address - IPv6 address of destination node(s).
- Extension Headers - IPv6 allows additional internet layer extension headers to be added after the standard IPv6 header. This is to allow more information for features such as fragmentation, authentication and routing information. The transport layer protocol header will then be addressed by this extension header.

2.1.3 Packet Handling

Once the kernel of the given operating system has received data to transmit from a given application, the data is then placed into a packet with the correct header and these packets are placed on a IP stack (Figure 5). Through a few intermediate steps the packets arrive at the driver queue, also known as the transmission queue. This queue is implemented as a ring buffer, therefore has a maximum capacity before it starts to overwrite packets which are still to be transmitted. As long as the queue isn't empty, the network interface card (NIC) will take packets from the queue and place them on the transmission medium. A similar process occurs when receiving packets, but in the opposite direction. For each NIC, there is a receive and transmit queue which are independent of each other allow communication to be bidirectional, although this depends on the kernel and its handling of events associated with packets.

is this section right

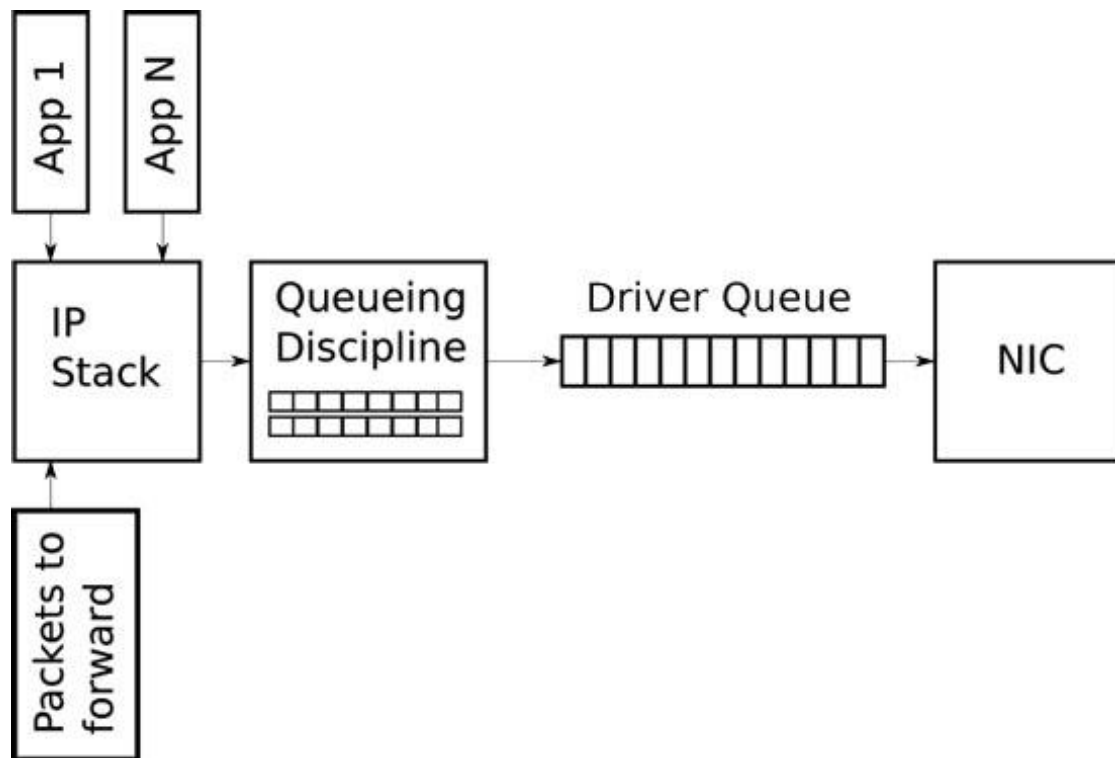


Figure 5: Linux packet handling [?]

2.1.4 Network Interface Controller (NIC)

Also known as a Network Interface Card, they provide the ability for a computer to connect to a network through a variety of mediums such as wireless, ethernet and fibre optics. They provide both the data link and physical layer of the network model (Figure 15), allowing a protocol stack to communicate with other computers on the local area network (LAN) or over a wider area via the IP protocol using IP addresses.

NICs can run at speeds of up to 100Gbps but more commonly run at 10Gbps for servers and 1Gbps for standard computers. The kernel or other applications retrieve packets via the NIC by polling or interrupts. Polling is where the kernel or application will periodically check the NIC for received packets while the use of interrupts allows the NIC to tell the kernel or application that it has received packets. Generally NICs provide the ability for 1 or more receive and transmit queues to be assigned per port, allowing for increased performance by assigning queues to different threads.

2.1.5 Middleboxes

Middleboxes are a device within a network that inspect, alter and forwards packets depending on certain rules and the intended functionality of the middlebox. A few different types are described below and can be combined into 1 single application:

Firewall Firewalls (Figure 6) are generally the major applications which sit between the public and private network of a system. They provide packet filtering which controls which packets can enter the private network via establishing a set of rules which packets have to adhere to. Filtering can be based on a number attributes of the packet such as the source and destination IP address and port and the destination service. Firewalls can also offer a number of other useful features such as NAT's or dynamic host configuration protocol (DHCP) to allow dynamic assignment of IP addresses within a network. As well as providing protection on a network level, application layer firewalls exist which stop certain applications from sending or receiving a packet.

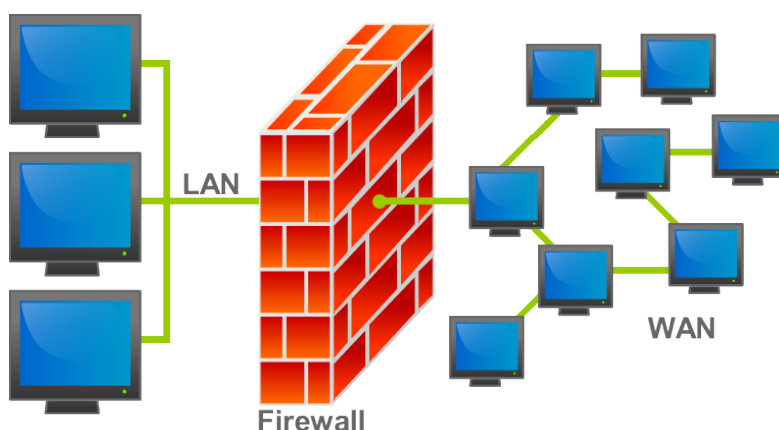


Figure 6: Firewall intercepting packets as a security measure [?]

Network Address Translator (NAT) As a routing device, a NAT is responsible for remapping an IP address to another by altering the IP datagram packet header. NAT's have become extremely important in modern networking systems due to IPv4 address exhaustion, allowing a single public IP address to map to multiple private IP addresses. This is particularly useful in large corporations where only a limited public network connection is required, meaning that all private IP addresses (usually associated with a single machine) are mapped to the same public IP address. A NAT will make use of multiple connection ports to identify which packets are for which private IP address and then re-assign the packet header so the internal routers can forward the packet correctly. As can be seen by Table 1 and Figure 7, each internal address is mapped to via the port number associated with the external address. NAT's are generally implemented as part of a network firewall as they inspect the datagram packets for malicious data and sources.

complete
this with
help from
links

Private IP Address	Public IP Address
10.0.0.1	14.1.23.5:62450
10.0.0.2	14.1.23.5:62451
10.0.0.3	14.1.23.5:62452
10.0.0.4	14.1.23.5:62453

Table 1: Example of public IP address and ports mapping to private IP address

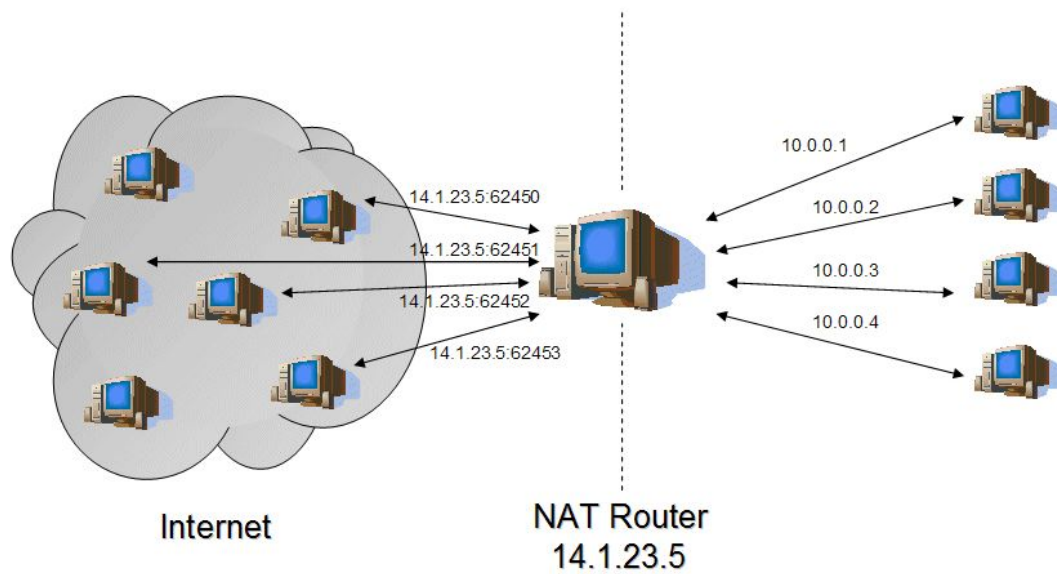


Figure 7: NAT translating public IP addresses into private IP addresses [?]

2.2 Java

2.2.1 JVM

The Java Virtual Machine is an abstract computer that allows Java programs to run on any computer without dependant compilation. This works by all Java source code been compiled down into Java byte code, which is interpreted by the JVM's just in time (JIT) compiler to machine code. However, it does require each computer to have the Java framework installed which is dependant on the OS and architecture. It provides an appealing coding language due to the vast support, frameworks and code optimisations available such as garbage collections and multithreading. Figure 8 shows the basic JVM architecture with the relevant sections explained in the list below.

any more
middle-
boxes -
packet
cap?

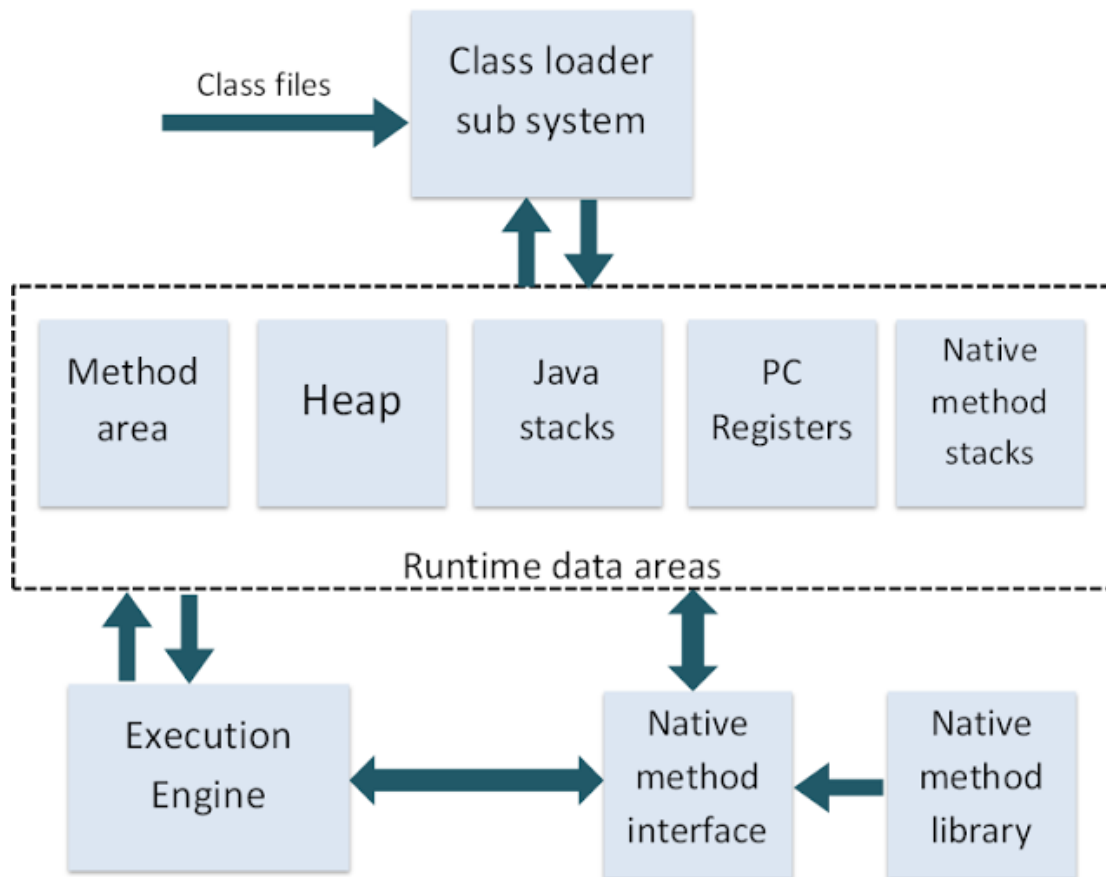


Figure 8: Java Virtual Machine interface [?]

- Class loader sub system - Loads .class files into memory, verifies byte code instructions and allocates memory required for the program
- Method area - stores class code and method code
- Heap - New objects are created on the heap
- Stack - Where the methods are executed and contains frames where each frame executes a separate method
- PC registers - Program counter registers store memory address of the instruction to be executed by the micro processor
- Native method stack - Where the native methods are executed.
- Native method interface - A program that connects the native method libraries with the JVM
- Native method library - holds the native libraries information

- Execution engine - Contains the interpreter and (JIT) compiler. JVM decides which parts to be interpreted and when to use JIT compiler.

Typically, any network communication from a Java application occurs via the JVM and through the operating system. This is because the JVM is still technically an application running on top of the OS and therefore doesn't have any superuser access rights. Any network operation results in a kernel system call, which is then put into a priority queue in order to be executed. This is one of the main reasons why network calls through the JVM and kernel can be seen as 'slow', in relative speeds compared to network line rate speeds.

2.2.2 Java Native Interface (JNI)

Provided by the Java Software Development Kit (SDK), the JNI is a native programming interface that lets Java applications use libraries written in other languages. The JNI also includes the invocation API which allows a JVM to be embedded into native applications. This project and therefore this overview will only focus on Java code using C libraries via the JNI on a linux based system.

In order to call native libraries from Java applications a number of steps have to be undertaken as shown below, which are described in more detail later:

1. Java code - load the shared library, declare the native method to be called and call the method
2. Compile Java code - compile the Java code into bytecode
3. Create C header file - the C header file will declare the native method signature and is created automatically via a terminal call
4. Write C code - write the corresponding C source file
5. Create shared library - create a shared library file (.so) from C code
6. Run Java program - run the Java program which calls the native code

The Java Framework provides a number of typedef's used to have equivalent data types between Java and C, such as jchar, jint and jfloat, in order to improve portability. For use with objects, classes and arrays, Java also provides others such as jclass, jobject and jarray so interactions with Java specific characteristics can be undertaken from native code run within the JVM.

```

1 public class Jni {
2
3     int x = 5;
4
5     public int getX() {
6         return x;
7     }
8
9     static { System.loadLibrary("jni"); }
10
11     public static native void objectCopy(Jni o);

```

```

12
13 public static void main(String args[]) {
14     Jni jni = new Jni();
15     objectCopy(jni);
16 }
17
18 }

```

Code 1: Basic Java class showing native method declaration and calling with shared library loading

Code 1 shows a simple Java program which uses some native C code from a shared library. Line 9 indicates which shared library to load into the application, which is by default lib*.so where * indicates the name of the library identifier. Line 11 is the native method declaration which specifies the name of the method and the parameters which will be passed to the corresponding C method. In this case, the method name is 'objectCopy' and a 'Jni' object is passed as a parameter. Line 15 is where this native method is called.

```

1 $ javac Jni.java

```

Code 2: Compiling basic Java program

Code 2, run from a terminal, compiles the Java class and create a class file which can be executed.

```

1 $ javah -jni Jni

```

Code 3: Generating C header file

In order to generate the C header file the command 'javah' (Code 3) is used with the flag 'jni' which tells Java that a header file is required which is for use with the JNI. It will then produce method signatures which correspond to the native method declared within Jni.java. The auto generated C header file is shown in Code 4.

```

1 /* DO NOT EDIT THIS FILE - it is machine generated */
2 #include <jni.h>
3 /* Header for class Jni */
4
5 #ifndef _Included_Jni
6 #define _Included_Jni
7 #ifdef __cplusplus
8 extern "C" {
9 #endif
10 /*
11  * Class:      Jni
12  * Method:     objectPrint
13  * Signature:  (LJni;)V
14  */
15 JNIEXPORT void JNICALL Java_Jni_objectPrint
16     (JNIEnv *, jclass, jobject);
17
18 #ifdef __cplusplus
19 }
20 #endif
21 #endif

```

Code 4: Auto-generated C header file


```

1 #include "Jni.h"
2
3 JNIEXPORT void JNICALL Java_Jni_objectPrint(JNIEnv *env, jclass class, jobject obj
4 ) {
5     jclass cls = (*env)->FindClass(env, "Jni");
6     jmethodID method = (*env)->GetMethodID(env, cls, "getX", "()I");
7     int i = (*env)->CallIntMethod(env, obj, method);
8     printf("Object x value is %i\n", i);
9 }

```

Code 5: C source file corresponding to auto-generated header file

The C source file implementation is in Code 5. The method signature on line 3 isn't as first declared in the Java source file. The 'env' variable is used to access functions to interact with objects, classes and other Java features. It points to a structure containing all JNI function pointers. Furthermore, the method invocation receives the class from which it was called since it was a static method. If the method had been per instance, this variable would be of type 'jobject'.

Line 4 shows how to find a class identifier by using the class name. In this example, the variables 'class' and 'cls' would actually be equal. In order to call an objects' method, a method id is required as a pointer to this method. Line 5 shows the retrieval of this method id, whose parameters are the jclass variable, the method name and the return type, in this case an integer (represented by an I). Then the method can be called on the object via one of the numerous helper methods (line 6) which differ depending on the return type and static or non-static context.

```

1 $ gcc -shared -fpic -o libjni.so -I/usr/java/include -I/usr/java/include/linux jni
   .c

```

Code 6: Terminal commands to generate shared library file (.so)

The command in Code 6 will create the shared object file called 'libjni.so' from the source file 'jni.c'. This output file is what the Java program uses to find the native code when called. It requires pointers to the location of the Java Framework provided jni.h header file.

```

1 $ java -Djava.library.path="." Jni
2 Object x value is 5

```

Code 7: Output from running Java application calling native C methods

Running the Java application, pointing to the location where Java can find the shared library (if not in a standard location) will output the above in Code 7.

Although the JNI provides a very useful interface to interact with native library code, there are a number of issues that users should be wary of before progressing:

- The Java application that relies on the JNI loses its portability with the JVM as it relies on natively compiled code.
- Errors within the native code can potentially crash the JVM, with certain errors been very difficult to reproduce and debug.
- Anything instantiated with the native code won't be collected by the garbage collector with the JVM, so freeing memory should be a concern.
- If using the JNI on large scale, converting between Java objects and C structs can be difficult and slow

2.2.3 Current Java Networking Methods

<http://haumacher.de/publ/parallel/p086.pdf> <http://www.javacoffeebreak.com/articles/javarmi/javarmi.html>

For high performance computing in Java, a number of existing programming options are available in order for applications to communicate over a network. These can be classified as: (1) Java sockets; and (2) Remote Method Invocation (RMI); (3) shared memory programming. As will be discussed, none of these are capable of truly high performance networking, especially at line rate speeds.

2.2.3.1 Java Sockets

Java sockets are the standard low level communication for applications as most networking protocols have socket implementations. They allow for streams of data to be sent between applications as a socket is one end point for a 2 way communication link, meaning that data can be read from and written to a socket in order to transfer data. Even though sockets are a viable option for networking, both of the Java socket implementations (IO sockets & NIO (new I/O) sockets) are inefficient over high speed networks [?] and therefore lack the performance that is required. As discussed previously, the poor performance is due to the JVM interacting with network cards via the OS kernel.

2.2.3.2 Remote Method Invocation (RMI)

Remote Method Invocation (RMI) is a protocol developed by Java which allows an object running in a JVM to invoke methods on another object running on a different JVM. Although this method provides a relatively easy interface for which JVM's can communicate, its major drawback relates to the speed. Since RMI uses Java sockets as its basic level communication method, it faces the same performance issues as mentioned in section 2.2.3.1.

2.2.3.3 Shared Memory Programming

Shared memory programming provides high performance JVM interaction due to Java's multithread and parallel programming support. This allows different JVM's to communicate via objects within memory which is shared between the JVM's. However, this technique requires the JVM's to be on the same shared memory system, which is a major drawback for distributed systems as scalability options decrease.

Even though these 3 techniques allow for communication between JVM's and other applications, the major issue is that incoming packets are still handled by the kernel and then passed onto the corresponding JVM. This means that packets are destined for certain applications, meaning that generic packets can't be intercepted and checked, which is a requirement for common middlebox software.

2.3 Related Works

2.3.1 jVerbs

Ultra-low latency for Java applications has been partially solved by the jVerbs [?] framework. Using remote direct memory access (RDMA), jVerbs provides an interface for which Java applications can communicate, mainly useful within large scale data centre applications.

RDMA is a technology that allows computers within a network to transfer data between each other via direct memory operations, without involving the processor, cache or operating system of either communicating computer. RDMA implements a transport protocol directly within the network interface card (NIC), allowing for zero copy networking, which allows a computer to read from another computer and write to its own direct main memory without intermediate copies. High throughput and performance is a feature of RDMA due to the lack of kernel involvement, but the major downside is that it requires specific hardware which supports the RDMA protocol, while also requiring the need for specific computer connections set up by sockets.

As jVerbs takes advantage of mapping the network device directly into the JVM, bypassing both the JVM and operating system (Figure 9), it can significantly reduce the latency. In order to have low level interaction with the NIC, jVerbs has a very thin layer of JNI calls which can increase the overhead slightly. However, jVerbs is flawed, mainly because it requires specific hardware to run on, firstly limited by the RDMA protocol reliant hardware and further by the required RDMA wrappers which are implemented by the creators. Also, it can only be used for specific computer to computer connection and not generally packet inspection.

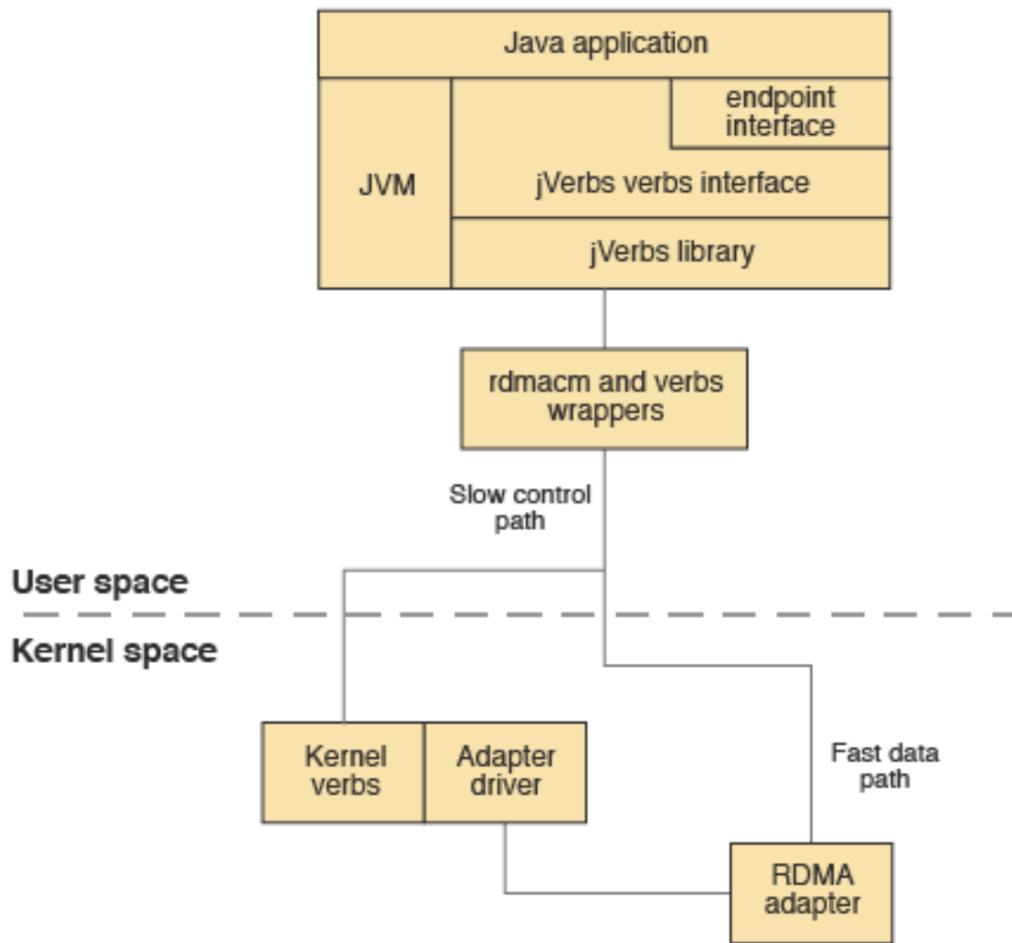


Figure 9: jVerbs architecture - shows how the framework bypasses the kernel and JVM [?]

jVerbs provides a useful example framework which re-emphasises that packet processing in Java is very possible with low latency, while assisting in certain implementation and design choices which can be analysed in more detail.

2.3.2 Native I/O API's

Currently available native networking API's are capable of reading and writing packets to the NIC transmission and receive queues at line rate. This is due to a number of techniques which tend to alter the kernels understand of the underlying NIC, therefore requiring specialist hardware and software to use such tools. DPDK (Section 2.3.4) is one of the tools that is open source and publicly available for use.

2.3.3 Packet Shader, infiniband - any more

2.3.4 Data Plane Development Kit (DPDK)

Data Plane Development Kit (DPDK) [?] is a set of libraries and drivers which enables fast packet processing reaching speeds of 80 Mpps on certain system set ups. Since DPDK is developed by Intel, it only supports Intel x86 CPU's and certain network interface controllers (NIC). DPDK binds the NIC's to new drivers meaning that the operating system doesn't recognise the network cards and can't interact with them. It make use of drivers run in user space allowing it to interact with certain memory locations without permission from the kernel or even involving it in any way.

DPDK makes use of an environment abstraction layer (EAL) which hides the environmental specifics and provides a standard interface which any application can interact with. Due to this, if the system changes in any way, the DPDK library needs to be re-compiled with other features been re-enabled in order to allow applications to run correctly again.

In order to use the DPDK libraries for the intended purpose, data packets have to be written into the correct buffer location so they are inserted onto the network. I similar approach is used when receiving packets on the incoming buffer ring, but instead of the system using interrupts to acknowledge the arrival of a new packet, which is performance costly, it constantly polls the buffer space to check for new packets. DPDK also allows for multiple queues per NIC and can handle multiple NIC's per system, therefore scalability is a major bonus of the libraries.

DPDK is very well documented on a number of levels. Firstly there is a online API which gives in depth details about what the methods, constants and structs do. There are a number of well written guides which give step-by-step details of how to install, set-up and use DPDK on various platforms and finally, there are many sample programs included with the build which give understanding of how the overall library works. DPDK is discussed in much more detail in Section 3.

2.4 Performance Testing Techniques

Benchmarking is a process of testing hardware, individual components or full end to end systems to determine the performance of the application or hardware . Generally, benchmarking should be repeatable under numerous iterations without only minor variations in performance results. This is firstly to allow minor changes to be made to the application/component with re-runs of the benchmark showing the performance changes. Secondly, it allows accurate comparisons to be drawn between similar software or hardware with different implementations in order to derive a better product.

Examples
of hard-
ware com-
parisons

better
word for
derive

why I need
bench-
marking

2.4.1 Programming Languages

It is well known that different programming languages can provide a radical change in execution for a given program. However, direct comparisons can't truly be trusted as certain languages are

suited for specific tasks and finding a benchmarking program to incorporate this is problematic. Other factors can be introduced when deciding on the optimisation level and the compiler of JIT used.

Numerous attempts have been made to compare languages, most noticeably the 'Benchmark Game' and Google.

2.4.1.1 Loop Recognition

Google inducted their own experiment on this problem, testing only C++, Java, Scala and Go on the loop recognition algorithm. Implementations made use of standard looping constructs and memory allocation schemes without the use of non-orthodox optimisation techniques. Selected results of this are shown below:

Benchmark	Time [sec]	Factor
c++	23	1.0x
Java 64-bit	134	5.8x
Java 32-bit	290	12.8x
Java 32-bit GC	106	4.6x
Scala	82	3.6x
Go 6g	161	7.0x

Table 2: Results from Loop Recognition benchmarking

2.4.1.2 Benchmark Game

The Benchmark Game is an online community which aims to find the best programming language by using multiple benchmarking algorithms running on different architecture configurations to determine the outcome. Again, even this community regard the best benchmark application to be your application. A few selected results are shown below between Java and C (those used in this report) for a few different benchmarks.

2.4.1.3 Using Economics

About economics paper

2.4.1.4 Which is better?

2.4.2 Intra-Language Techniques

2.4.3 Applications

show example of optimising java and how different it looks, also different depending on architecture

better phrasing off that sentence needed

More here

get paper

explain difference with GC etc

get data

3 DPDK

This section will go into much more detail about the Data Plane Development Kit (DPDK), focussing on the basic concepts used by the fast packet processing framework for use within custom applications. DPDK does have a number of more advanced features which can be exploited but aren't discussed in this report.

ref manual

3.1 Environment Abstraction Layer (EAL)

The EAL abstracts the specific environment from the user and provides a constant interface in order for applications to be ported to other environments without problems. The EAL compilation may change depending on the architecture (32-bit or 64-bit), operating system, compilers in use or network interface cards. This is done via the creation of a set of libraries that are environment specific and are responsible for the low-level operations gaining access to hardware and memory resources.

On the start-up of an application, the EAL is responsible for finding PCI information about devices and addresses via the `igb_uio` kernel module, performing physical memory allocation using huge pages (section 3.3) and starting user-level threads for the logical cores (section 3.2) (lcores) of the application.

3.2 Logical Cores

Also known as lcores within DPDK, logical cores shouldn't be confused with processor cores. Lcores are threads which allow different applications functions to be run within different threads. This can allow different lcores to control different ports and queues while processing packets as well.

Within DPDK lcores are implemented with POSIX threads (on linux) and make use of processor affinity (CPU pinning). This allows lcores to be only run on certain processing cores which reduces context switching and cache memory swapping and therefore increasing overall performance. However, this only works if the number of lcores is equal or less than the number of available processing cores so the number of lcores which are allowed to be initiated are limited. This feature can also exploit modern processors which allow for hyper-threaded cores but the performance increase isn't as guaranteed as it is with non-hyper-threaded cores.

3.3 Huge Pages

Huge pages are a way of increasing performance when dealing with large amounts of memory. Normally, memory is managed in 4096 byte (4KB) blocks known as pages, which are listed within the CPU memory management unit (MMU). However, if a system uses a large amount of memory, increasing the number of standard pages is expensive on the CPU as the MMU can only handle thousands of pages and not millions.

The solution is to increase the page size from 4KB to 2MB or 1GB (if supported) which keeps the number of pages references small in the MMU but increases the overall memory for the system. Huge pages should be assigned at boot time for better overall management, but can be manual assigned on initial boot up if required.

DPDK makes uses of huge pages simply for increased performance due to the large amount of packets in memory. This is even the case if the system memory size is relatively small and the application isn't processing an extreme number of packets.

3.4 Ring Buffer

A ring buffer is used by DPDK to manage transmit and receive queues for each network port on the system. This fixed sized first-in-first-out queue allows multiple objects to be enqueued and dequeued from the ring at the same time from any number of consumers or producers. A major disadvantage of this is that once the ring is full (more of a concern in receive queues), it allows no more objects can be added to the ring, resulting in dropped packets or packet caches. It is therefore imperative that applications can processes packets at the required rate.

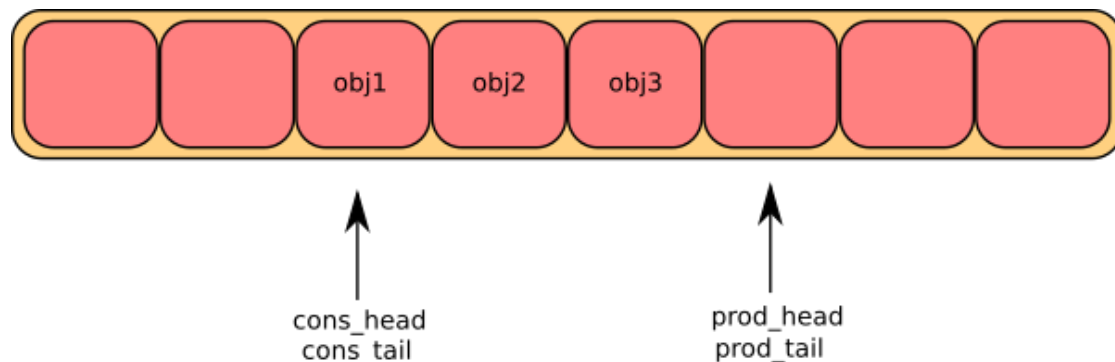


Figure 10: Ring Buffer

[cite this](#)

3.5 Memory Usage

DPDK tries to regulate memory usage in order to be more efficient and therefore handles its own memory management. However, it does allow application to allocate memory blocks for specific port queues and for data sharing between lcores. The major memory techniques used are described below.

should code use dpdk malloc?

3.5.1 Allocation

Since DPDK make use of huge pages, it provides its own memory allocation (malloc) library to allow memory blocks of any size, which also improves the portability of applications. However, even the DPDK malloc library is slow compared to pool memory access due to synchronisation constraints. Generally this library should only be used at initialisation time but does support NUMA for specific socket memory access and memory alignment.

3.5.2 Pools

Memory pools allow for fixed size allocation of memory which uses a ring to store fixed sized objects. These are almost guaranteed to be used for message buffer storage for receive and transmit queues for ports. They are initialised with a number of parameters to increase performance such as cache sizes and NUMA socket identification as well as a name identifier.

Pools offer increased performance over the standard memory allocation since the object padding is optimised so each object starts on a different memory channel and rank. Furthermore, a per core cache can be enabled at initialisation. This offers performance advantages, as without caching per core locks are required for every pool access. Caches offer cores lock free access to data, while bulk requested can be carried out on the pool to reduce locking on the pool.

3.5.3 Message Buffers

Message Buffers (mbufs) are stored within a specified memory pool and are used to carry data between different processes within the application and are primarily used for carrying network packets through the application. Mbufs also contain metadata about the information it is carrying which includes the data length, message type and offsets for the start of the data. The headroom shown in figure 11 shows empty bytes between the metadata and start of the data which allows the data to be memory aligned for quicker access. Mbufs can also be chained together to allow for longer data, more commonly jumbo packets.

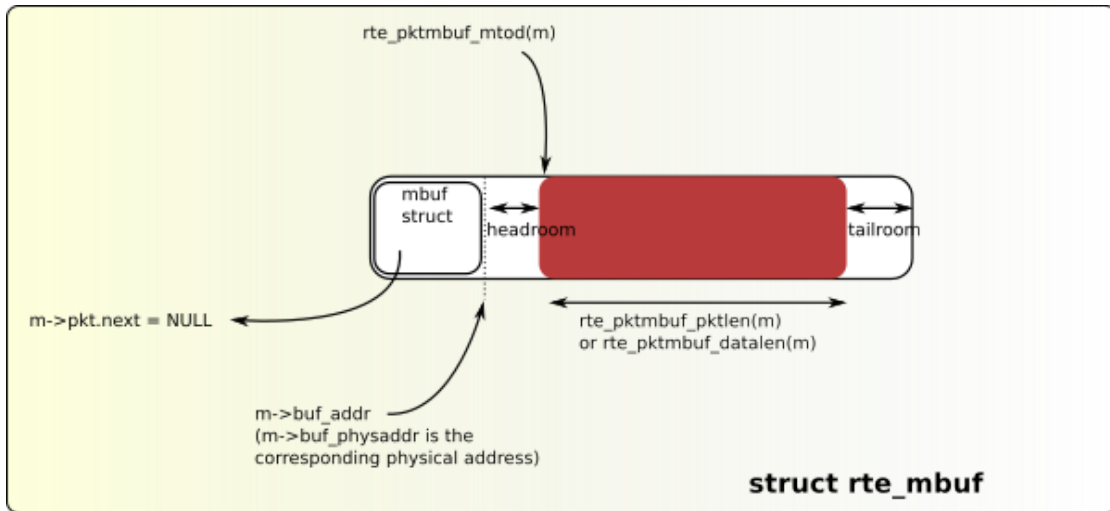


Figure 11: Message Buffer

cite this

3.5.4 NUMA

DPDK can make use of non-uniform memory access (NUMA) if the system supports it. NUMA is a method for speeding up memory access when multiple processors are trying to access the same memory and therefore reduces the processors waiting time. Each processor will receive its own bank of memory which is faster to access as it doesn't have to wait. As applications become more extensive, processors may need to share memory, which is possible via moving the data between memory banks. This somewhat negates the need for NUMA, but NUMA can be very effective depending on the application. DPDK can make extensive use of NUMA as each logical core is generally responsible for its own queues, and since queues can't be shared between logical cores (although dedicated ring buffers can), data sharing is rare.

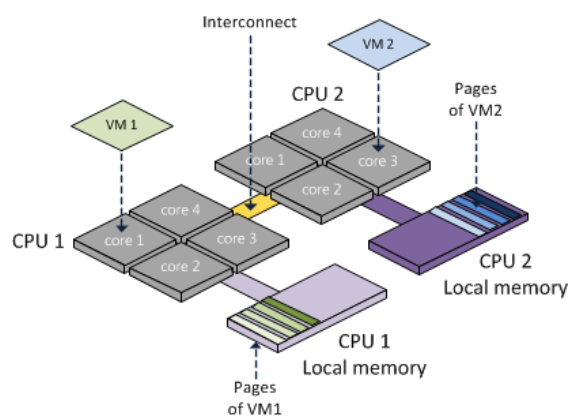


Figure 12: NUMA

3.6 Poll Mode Driver (PMD)

A Poll Mode Driver (PMD) is a user space device which allows configuration of network interface cards/controllers and their associated queues. PMD's work without interrupts which allow for quicker receiving, processing and transmitting of packets and are therefore lock free. Generally anything which can be achieved by interrupts can also be achieved via using rings and continuous polling of the rings. This means that 2 lcores running in parallel can't receive from the same queue on the same port. However, 2 parallel cores can receive from the same port on different queues.

There are a number of considerations for application design depending on the hardware in use. Optimal performance can be achieved by carefully considering the hardware properties such as caches, bus speed and bandwidth along side the software design choices. For example, NIC's are more efficient at transmitting multiple packets in a burst rather than individually but consequently overall throughput may be reduced.

3.7 Models

DPDK supports 2 methods for packet processing applications:

Pipe-line This is an asynchronous model where lcores are designated to perform certain tasks. Generally certain lcores will receive packets via the PMD api and simply pass those packets to other lcores via the use of a ring. These other lcores will then process the packets depending on the requirements and then either forward the packets to the PMD for transmitting or pass them onto other lcores via a ring.

Run-to-completion This is a synchronous model where each lcore will retrieve the packets, process them individually and then output them for transmission. Each lcore should be assigned its own receive and transmit queue on a given port in order to negate the need for locks. This report will focus on the run-to-completion model.

4 Initial language comparison

Before any implementation or specific design considerations were undertaken, an evaluation of the performance of C, Java and the Java Native Interface (JNI) was carried out. Although data from existing articles and websites could be used for Java and C, there was no existing direct comparisons between them and the JNI, therefore custom tests were carried out.

The JNI is inherently seen as a bottleneck of an application (even after its vast update in Java 7).

article on
this

As this application would be forced to use the JNI, numeric values of its performance was helpful to evaluate the bridge in speed required to be overcome.

reasons
why JNI is
slow

4.1 Benchmarking Algorithm

As discussed previously, there are always advantages and disadvantages of any algorithm used for benchmarking. In order to minimise the disadvantages, an algorithm was used which tried to mimic the procedures which would be used in the real application, just without the complications. Algorithm 1 shows that the program basically creates 100,000 packets individually and populates their fields with random data, which is then processed and returned in the 'result' field. This simulates retrieving low-level packet data, interpreting and acting upon the data, and then setting data within the raw packet.

ref this

Algorithm 1 Language Benchmark Algorithm

```
1: function MAIN
2:   for i = 1 to 100000 do
3:      $p \leftarrow \text{Initialise Packet}$ 
4:     POPPACKET(p)
5:     PROPACKET(p)

6: function POPPACKET(Packet p)                                ▷ Set data in a packet
7:    $p.a \leftarrow \text{randomInt}()$ 
8:    $p.b \leftarrow \text{randomInt}()$ 
9:    $p.c \leftarrow \text{randomInt}()$ 
10:   $p.d \leftarrow \text{randomInt}()$ 
11:   $p.e \leftarrow \text{randomInt}()$ 

12: function PROPACKET(Packet p)                                ▷ Process a packet
13:    $res \leftarrow p.a * p.b * p.c * p.d * p.e$ 
14:    $p.result \leftarrow res$ 
```

For the JNI version, the same algorithm was used, however, the *PopPacket* method was carried out on the native side to simulate retrieving raw packet data. The *ProPacket* method was executed on the Java side with the result been passed back to the native side to be entered back into the packet structure.

Timing within the algorithm for all variations was carried out between each iteration. This firstly eliminated any initial start-up time associated with the application which is common with the JVM. Secondly, any calls for time stamps to the system would be minimised as 100,000 iterations would occur in-between them.

4.2 Results

Each language had the algorithm run 1,000 times in order to minimise any variations due to external factors. Figures show that C was considerably quicker than Java, while Java using the JNI was extremely slow.

ref this

expand on
this

4.3 Further Investigation

Due to the very poor performance of the JNI compared to other languages, further investigations were carried out to find more specific results surrounding the JNI.

Is this rel-
evant

5 Design Considerations

5.1 Data Sharing

The proposed application will be sharing data between the DPDK code written in C (low level) and the Java (medium level) side used for the highly abstracted part of the application. This requires a large amount of data, most noticeably packets, to be transferred between 'sides' in a small amount of time. This section will explore possible techniques for this data sharing, discuss the advantages and disadvantages and finally reason about a decision.

expand
on where
data is
allocated
between
native and
java - use
links above

Diagram of packets from NIC using c through 'technique' and then processing packets in java and then back

5.1.1 Objects and JNI

By far the simplest technique available is using the Java Native Interface (JNI) in order to interact with native code and then retrieve the required data via object method accessing provided by the JNI environment pointer. This can be done 2 ways, either by creating the object and passing it as a parameter to the native methods or creating an object on the native side. Both ways require the population of the fields to be done on the native side. From then on, any data manipulation and processing could be done on the Java side. Unfortunately, this does require all data to be taken from the object and placed back into the structs before packets can be forwarded. Obviously this results in a lot of unneeded data copying, while the actual JNI calls can significantly reduce the speed of the application as shown in .

ref this

There are a few frameworks which aim to solve the problem of mapping structs to objects while trying to increase the overall speed of the JNI using different techniques.

Java Native Access (JNA) The JNA is a community developed framework which aims to abstract the JNI and native code away from developers who want to use native shared libraries. It uses its own data mapping between native and Java datatypes with automatic conversion of string and object to character arrays and structs. Although this framework provides much of the functionality within this project, it focusses on correctness and ease of use and therefore neglects performance.

Javolution Javolution is similar to the JNA but focusses more on real-time performance within Java. It provides a number of high performance utilities for this purpose while allowing access to native libraries. For this purpose, struct and union classes are created automatically for mapping between objects. Javolution provides the functionality required for this project but it requires a lot initial set-up which can be heavy, especially for small applications like middleboxes.

Swig

do this

Preon

do this

5.1.2 ByteBuffers

ByteBuffers are a Java class which allow for memory to be allocated on the Java heap (non direct) or outside of the JVM off heap (direct), while abstracting pointer arithmetic and boundary checks away from the application. Non direct ByteBuffer's are simply a wrapper for a byte array on the heap and are generally used as they allow easier access to the data as multiple byte reading and datatype conversion is handled by the class. Direct ByteBuffers allocate memory outside of the JVM in native memory which means that the only limit on the size of bte buffers is memory itself.

cant use
byte-
buffers di-
rectly with
DPDK as
dont use
hugepages

The performance characteristics of direct and non-direct byte buffers are very similar. Direct byte buffers can be improved by aligning the bytes with the native endian (normally little endian) format since the class makes use of the Java Unsafe native access which allows for the methods to be inlined with machine code.

The Java garbage collector doesn't have access to native memory and therefore direct byte buffers allow the application to manage memory and reduce memory usage. Garbage collection bottlenecks within high performance system are therefore reduced, while allowing direct referencing of objects stored natively as the garbage collector doesn't compact the heap and move objects around.

Direct byte buffers would be an ideal choice for data sharing if it was possible to redirect DPDK mbuf structs onto the byte buffer directly, consequently allowing the Java side easy access to the data using the api. However, since DPDK make use of huge pages and handles its own memory management, it's not possible to use byte buffers directly as the memory for the NIC queues. Instead, data would have to be copied from structs and extracted into objects and vice versa.

5.1.3 Java Unsafe

The Java Unsafe class is actually only used internally by Java for its memory management. It generally shouldn't be used within Java since it makes a safe programming language like Java an unsafe language (hence the name) since memory access exceptions can be thrown which will ultimately crash the JVM. Even so, it has a number of uses such as:

Object initialisation skipping This is where any instance of an object can be created from the class, but no constructors are used meaning that the object created without any of the fields set. This has a number of uses including security check bypassing, creating instances of objects which don't have a public constructor and allowing multiple objects of a singleton class.

Intentional memory corruption This allows the setting of private fields of any object. It is a common way of bypassing security features as private fields to allow access to certain situation can be overwritten to gain access.

Nullifying unwanted objects This has a common use of nullifying passwords after they have been stored as strings. If a password is stored as a string in Java, even setting the field to (null)

will only dereference it. The original string will still be in memory after the dereferencing up until it is garbage collected. Even rewriting the field with a new field won't work as strings are immutable in Java. This makes it susceptible to a timing attack to retrieve the password. Using unsafe allows for the actual memory location to be overwritten with random values to prevent this.

Multiple inheritance Java doesn't allow multiple inheritance within its class declaration of casting. However, using Unsafe any object can be cast to any other object without a compiler of run-time error. This obviously only works if data fields are compliant with each other and any method invocations are referenced.

Very fast serialization The Java *Serializable* abstraction is well known to be slow, which can be a major bottleneck in fast processing applications over a network. Even the *Externalizable* functionality isn't much factor and that requires a class schema to be defined. However, custom serializing can be extremely fast using the Unsafe Class. Basically allocating memory and then putting/getting data from the memory requires little JVM usage and can be done with machine instructions.

Obviously without proper precautions any of these actions can be dangerous and can result in crashing the full JVM. This is why the Unsafe class has a private constructor and calling the static `Unsafe.getUnsafe()` will throw a security exception for untrusted code which is hard to bypass. Fortunately, Unsafe has its own instance called 'theUnsafe' which can be accessed by using Java reflection :

ref this

```
1 Field f = Unsafe.class.getDeclaredField("theUnsafe");
2 f.setAccessible(true);
3 Unsafe unsafe = (Unsafe) f.get(null);
```

Code 8: Accessing Java Unsafe

Using Unsafe then allows direct native memory access (off heap) to retrieve data in any of the primitive data formats. Custom objects with a set structure can then be created, accessed and altered using Unsafe which provides a vast increase in performance over traditional objects stored on the heap. This is mainly thanks to the JIT compiler which can use machine code more efficiently by inlining certain memory access directly with assembly code. This also removes the need for copying of data between memory locations, structs and objects, therefore meaning it is zero-copy.

5.2 Packing Structures

Structures (structs) are a way of defining complex data into a grouped set in order to make this data easier to access and reference as shown in Code 9. They are heavily used with C applications and can be seen as Java object without the associated methods.

```
1 struct example {
2     char *p;
3     char c;
4     long x;
5     char y[50];
```



```

6     int z;
7 };

```

Code 9: Example C Struct

On modern processors all commercially available C compilers will arrange basic C datatypes in a constrained order to make memory access faster. This has 2 effects on the program. Firstly, all structs will actually have a memory size larger than the combined size of the datatypes in the struct as a result of padding. However, this generally is a benefit to most consumers as this memory alignment results in a faster performance when accessing the data.

Explain why it has faster performance

Nested padding in struct?

C struct field always in given order

Inconsistencies with datatype length so using uint32t etc

Code 12 shows a struct which has compiler inserted padding. Any user wouldn't know the padding was there and wouldn't be able to access the data in the bits of the padding through conventional C dereferencing paradigm (only via pointer arithmetic). This example does assume use on a 64-bit machine with 8 byte alignment, but 32-bit machines or a different compiler may have different alignment rules.

```

1 struct example {
2     char *p;           // 8 bytes
3     char c;            // 1 byte
4     char pad[7];       // 7 byte padding
5     short x;           // 2 bytes
6     char pad[6];       // 6 byte padding
7     char y[50];        // 50 bytes
8     int z;             // 4 bytes
9 };

```

Code 10: Example C Struct with compiler inserted padding

Mention this is on 64-bit machine and obviously you don't notice padding and order of elements can pay an important part in this

```

1 struct __attribute__((__packed__)) example {
2     char *p;           // 8 bytes
3     char c;            // 1 byte
4     short x;           // 2 bytes
5     char y[50];        // 50 bytes
6     int z;             // 4 bytes
7 };

```

Code 11: Example C Struct stopping padding

Since the proposed application in this report requires high throughput of data, the initial thought would be that this optimisation is a benefit to the program. Generally this is the case, but for data which is likely to be shared between the C side and Java side a large amount, data accessing

is far quicker on the Java side if the struct is packed (no padding). This results in certain structs been forced to be packed when compiled, more noticeably, those used for packet and protocol headers.

Proof on
speed

Packed structures mean there are no gaps between elements, required alignment is set to 1 byte. Also `__attribute__((packed))` definition means that compiler will deal with accessing members which may get misaligned due to 1 byte alignment and packing so reading and writing is correct. However, compilers will only deal with this misalignment if structs are accessed via direct access. Using a pointer to a packed struct member (and therefore pointer arithmetic) can result in the wrong value for the dereferenced pointer. This is since certain members may not be aligned to 1 byte. In the below example, `uint32` is 4 byte aligned and therefore it is possible for a pointer to it to expect 4 byte alignment therefore resulting in the wrong results.

```
1 #include <stdio.h>
2 #include <inttypes.h>
3 #include <arpa/inet.h>
4
5 struct packet {
6     uint8_t x;
7     uint32_t y;
8 } __attribute__((packed));
9
10 int main ()
11 {
12     uint8_t bytes[5] = {1, 0, 0, 0, 2};
13     struct packet *p = (struct packet *)bytes;
14
15     // compiler handles misalignment because it knows that
16     // "struct packet" is packed
17     printf("y=%"PRIx32", ", ntohl(p->y));
18
19     // compiler does not handle misalignment - py does not inherit
20     // the packed attribute
21     uint32_t *py = &p->y;
22     printf("py=%"PRIx32"\n", ntohl(*py));
23     return 0;
24 }
```

Code 12: Example C Struct with compiler inserted padding

On an x86 system (which does not enforce memory access alignment), this will give `y=2` and `*py=2` which is as expected. Conversely, other systems using the SPARC architecture or similar will give `y=2` and `*py=1` which isn't what the user would expect.

However, since a packed struct is much easier to traverse from Java than a padded struct, the decision was made to make certain structs packed within the DPDK framework and then recompile the libraries. This decision could be made since other structs within the DPDK framework were also packed and therefore consideration of this was already made.

Note that if a struct contains another struct, that struct should be packed recursively as-well to ensure the first struct has no padding at all.

Char doesn't have alignment and can start on any address. But 2-byte shorts must start on an even address, 4-byte ints or floats must start on an address divisible by 4, and 8-byte longs or doubles must start on an address divisible by 8. Signed or unsigned makes no difference.

Self-alignment makes access faster because it facilitates generating single-instruction fetches and puts of the typed data. Without alignment constraints, on the other hand, the code might end up having to do two or more accesses spanning machine-word boundaries. Characters are a special case; they're equally expensive from anywhere they live inside a single machine word. That's why they don't have a preferred alignment.

Casting to an odd pointer will slow down code and could work. Other architectures will take the word which the pointer points to and therefore the problem occurs above.

5.3 Performance testing

In order to evaluate the most suitable data sharing technique, performance testing on 4 different implementation options for sharing data between Java and native memory were considered. Since the ultimate aim of the implementation is to maximise throughput of packets, the performance test tried to mimic this by processing data on 1 million packets per iteration. This processing involved retrieving data from the native packet struct, loading that data into a Java object, changing the data and then setting the data back into the original struct. Various techniques to do this were used to try and find the best performance possible. All of the techniques made use of a static native struct which acted like a new packet been received. The changed data was then set back into this struct.

Algorithm 2 Data Sharing Performance Test Algorithm

```

1: function MAIN
2:   for i = 1 to 10 do
3:     startTimer
4:     PERFORMTEST( )
5:     stopTimer
6:     outputTime

7: function PERFORMTEST
8:   for i = 1 to 1000000 do
9:     retrieveData
10:    setNewData
11:    saveData

```

Considering there were 4 different data sharing techniques tested, the *retrieveData*, *setNewData* and *saveData* methods were different and are described below:

Object Passing The object technique involved creating a packet and sending its reference through the JNI to the native code. From there, this packet could be populated with data from a given struct through the Java environment pointer. For each setter method called, the method id of that method must be retrieved for the given class so the combination of that and the packet reference could set the data.

From there, new data is input into the packet from the Java side and passed back through to JNI so the struct can be set with the new data. This is done using getters for the objects' fields via the Java environment pointer.

Byte Buffer The technique involved declaring a direct byte buffer to assign off heap memory of the size of the packet struct. From there, the pointer to this memory location was sent to the native code, where data from the struct was directly copied into the byte buffer, therefore populating the byte buffer with duplicate data. The Java code then pulled the data from this via the byte buffer's inbuilt methods and set the data into a new packet object. From there the packet object could be used whenever desired.

To save the data back into the original struct, the data was copied from the packet object back into the byte buffer in the order of the members of the struct. The data was then pulled from the memory of the byte buffer and set back into the original struct.

Unsafe Using the Unsafe class allows for direct access to members of the struct. To take account for this this technique first allocates off heap memory for the pointer to the struct to be stored. This pointer is put into the memory location natively and accessed via Unsafe methods. From this, data can be accessed directly from the struct and input to a new packet object for later use. New data is then set in the object.

To set the data in the struct, the data is removed from the object and directly put into the struct. This is done using the pointer and the correct byte offsets depending on the previous data type inputted.

Direct Direct accessing relates to not storing the struct members in Java at all. Instead a different type of Packet object is used which just stores pointers to the struct. From there any accessing and setting of data is simply done using the Unsafe class to directly get or set the values within the struct. This different packet object also contains offset information for the struct so the correct values are accessed.

5.3.1 Expectations

Of the 4 techniques, it is expected that the object method would be by far the slowest, mainly due to the excessive number of JNI calls used which significantly slows an application down. The Byte Buffer method will most likely be the slowest of the other 2 techniques due to the large number of data copying which goes on. The other 2 techniques (unsafe and direct) will likely be very close in performance, mainly because they both use direct accessing into the struct. The unsafe technique should be slightly slower however due to the setting and getting from the packet object.

5.3.2 Results

The graphs below show the results from the different techniques run on different systems. Considering there were 10 iterations of 1 million packets per technique the results show the averages of the times. However, certain times were disregarded and seen as anomalies since they were well above the average. These were generally at the start of the new technique and can be either be

evaluate each technique on ease and number of data copies and results and say why we picked 1 of the others

stats for this

related to the just in time compiler warming up after switching to a different class or garbage collection on the previous technique.

The graphs show the time in nanoseconds which it took to process an individual packet. The scale is logarithmic due to the excessive size of the object technique. For easier reading, the numbers at the top of the columns show the times factor for each technique compared to the fastest. For example, on Figure [the byte buffer takes 2.74 times the direct technique to process the packet](#).

ref this

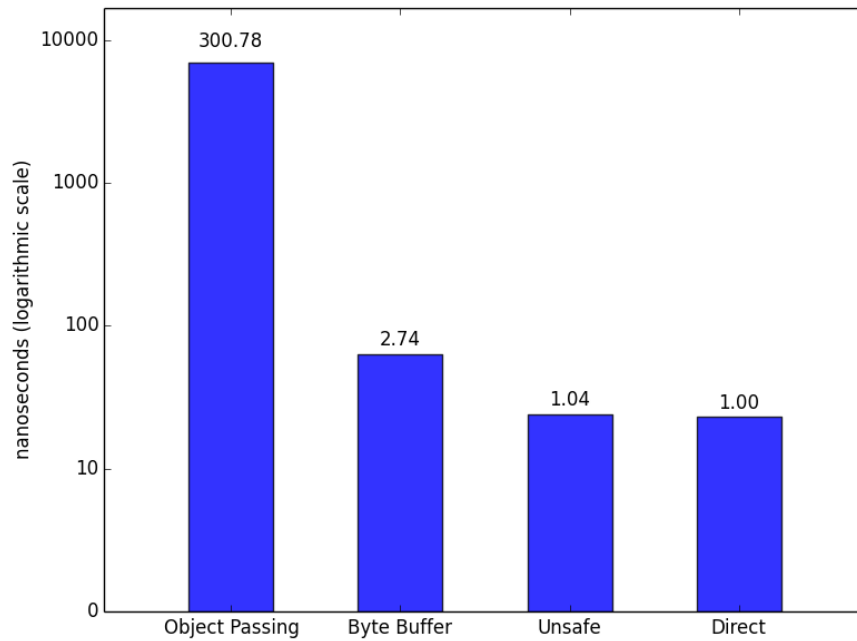


Figure 13: Server

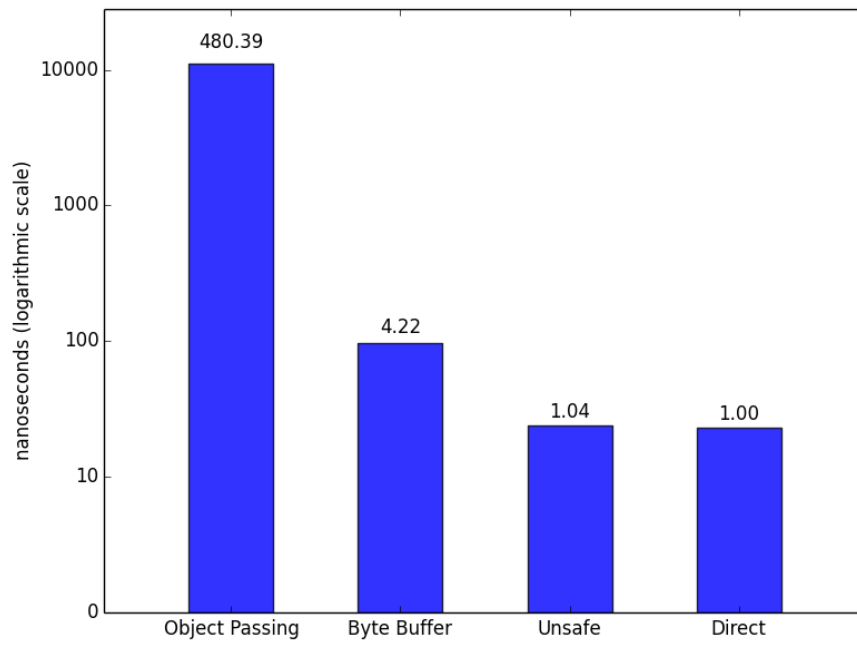


Figure 14: Mac

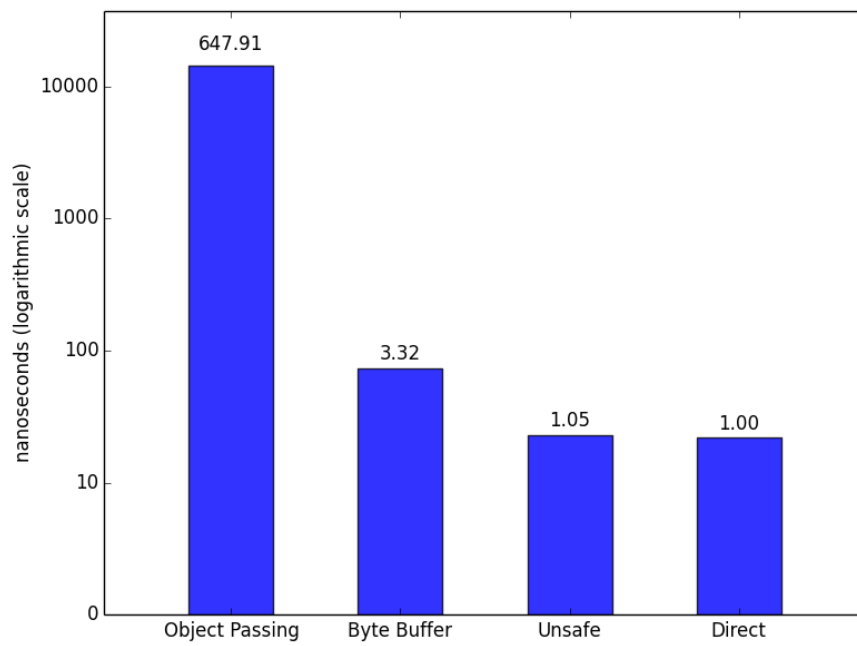


Figure 15: VM

should these graphs all be together because they look the same

5.3.3 Evaluation

As the graphs show, the exceptions declared earlier have been met. The object technique should be completely disregarded, even though it offers the easiest solution. Due to its abysmal performance, been 300 - 650 times slower than the faster technique, its not going to help with packet throughput.

For similar reasons Byte buffer should be disregarded as well simply because its too slow compared to the other methods. The amount of data copying with this technique is the downfall, although this could have potentially worked if structs could be directly written into the buffers.

cant be
done and
reference
earlier

The difference between the direct and unsafe methods is minimal on all 3 machines. The direct technique is slightly faster since it doesn't copy the data into the object although Java is surprisingly fast at this as proven by the unsafe results. Furthermore, in general middleware software, there isn't the requirement to access all data fields of the packets generally, therefore storing them all in objects can be seen as needless copying. This is where the direct access technique excels since it only accesses data which it needs to be zero-copy. Direct access was therefore chosen as the technique to use for the implementation which goes into further details about how this was done.

extension to remove packet pointer passing and do it directly from struct array - would have to dynamically pack array though - how????

do hardware testing - ie hardware faster than programs (i think) so should aim for throughput and not efficiency

5.4 Thread affinity

`Thread.currentThread().getId();` just gets id of thread relative to jvm.

It keeps a process limited to certain a certain core or cores. Process will still be taken out of use and switched back in but without the problem of moving cache between cores.

Normally as a thread gets a time slice (a period in which to use the core), it is granted whichever core [CPU] is determined to be most free by the operating system's scheduler. Yes, this is in contrast to the popular fallacy that the single thread would stay on a single core. This means that the actual thread(s) of an application might get swapped around to non-overclocked cores, and even underclocked cores in some cases. As you can see, changing the affinity and forcing a single-threaded CPU to stay on a single CPU makes a big difference in such scenarios. The scaling up of a core does not happen instantly, not by a long shot in CPU time.

Therefore, for primarily single (or limited) thread applications, it is sometimes best to set the CPU affinity to a specific core, or subset of cores. This will allow the 'Turbo' processor frequency scaling to kick in and be sustained (instead of skipping around to various cores that may not be scaled up, and could even be scaled down).

core thrashing - ust by the name, you know this is a bad thing. You lose performance when a thread is swapped to a different core, due to the CPU cache being 'lost' each time. In general, the *least* switching of cores the better. One would hope the OS would try to avoid this, but it doesn't seem to at all in quick tests under Windows 7. Therefore, it is recommended you manually adjust the CPU affinity of certain applications to achieve better performance.

Another important issue is avoiding placing a load on a HyperThreaded (non-physical) core. These cores offer a small fraction of the performance of a real core. The Windows scheduler is aware of this and will swap to them only if needed. As of mid Jan 2012 the Windows 7 and Windows 2008 R2 schedulers have a hotfix for AMD Bulldozer CPUs that see them as HyperThreaded, cutting them down from 8 physical cores to 4 physical cores, 8 logical cores. This is for two reasons: The AMD Bulldozer platform uses pairs of cores called Bulldozer Modules. Each pair shares some computation units, such as an L2 cache and FPU. To spread out the load and prevent too much load being placed on two cores that have shared computational units, the Windows patch was released, boosting performance in lightly threaded scenarios.

Processor affinity takes advantage of the fact that some remnants of a process that was run on a given processor may remain in that processor's memory state (for example, data in the CPU cache) after another process is run on that CPU. Scheduling that process to execute on the same processor could result in an efficient use of process by reducing performance-degrading situations such as cache misses. A practical example of processor affinity is executing multiple instances of a non-threaded application, such as some graphics-rendering software.

```
1 cpu_set_t cpuset; \\ structure used to manage cpu affinity settings
2 pthread_t thread = pthread_self(); \\ get own system wide thread id
3 CPU_ZERO(&cpuset); \\ zero cpuset structure
4 CPU_SET(5, &cpuset); \\ set cpuset structure use with the 6th core (0-5)
5 int res = pthread_setaffinity_np(thread, sizeof(cpu_set_t), &cpuset); \\ set
    affinity of thread
6 \\ error handling here of res
7 res = pthread_getaffinity_np(thread, sizeof(cpu_set_t), &cpuset); \\ check
    affinity of thread
8 \\ error handling here of res
```

Code 13: Example of setting thread affinity

[ref this](#)

In Linux, Java thread uses the native thread(i.e, thread provided by Linux). This means the JVM creates a new native thread when the Java code creates a new Java thread. So, the Java threads can be organised in any way the native threads can be organised.

A native thread can be bound to a core through the `pthread_setaffinity_np()` function. So, a Java thread can be bound to a core. If Java standard library does not provide a function to do so, then this function need to be provided through JNI.

In Linux, multi-threading is same as parallel threading. Linux kernel distribute threads among processors to balance the cpu load. However the individual threads can be bound with any core as wished. So, in Linux Java multi-threading is same as parallel threading.

5.5 Endianness

This describes the order in which bytes of data types are stored in memory relative to pointers. In big-endian systems, the most significant bytes are stored at the pointer with every successive data bytes stored in successive memory locations. Conversely, in little-endian systems, the least significant byte is stored at the pointer. There is no advantage or disadvantage to either endian types, its simply a matter on convention for certain systems. Figure 16 shows how different endian systems store the value of the hexadecimal value `0x0A0B0C0D`.

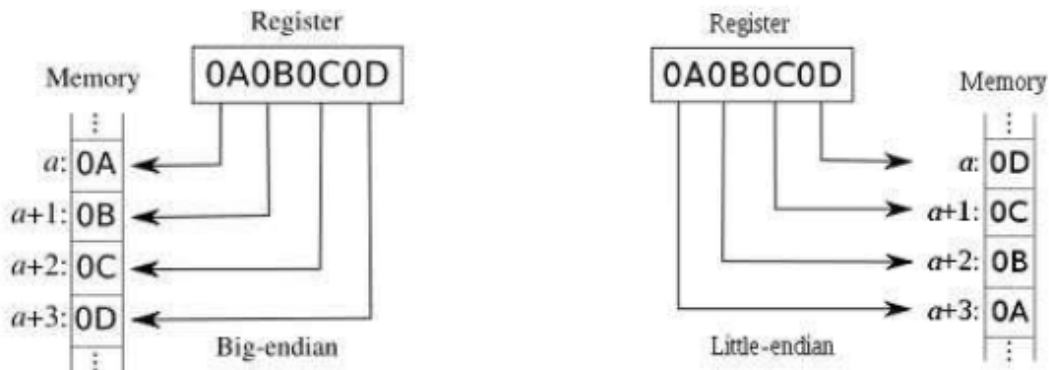


Figure 16: Difference between little and big endian

This becomes a problem when using the Java Native Interface with native code code on Intel architectures since the Java Virtual Machine uses big-endian format while Intel uses little-endian format. Normally the JNI environment would handle this byte ordering conversion but since the JNI has been proven to be too slow to meet the requirement, this needs to be handled by the application.

ref this

5.6 Data type conversion

Between different languages the same data type can be represented by varying lengths in bytes and whether unsigned or signed (in the case of numerical values). It can also be the case that data type lengths in native languages differ depending on the architecture and whether it is 32-bit or 64-bit.

This becomes more of an issue when sharing data between C and Java. Java always uses a signed integer representation with the most significant bit representing whether the number is negative or positive. C uses both signed and unsigned representation depending on the requirements with varying byte length. Again, any conversion between the differences is normally handled via the JNI, but since the implementation aims to bypass the JNI as much as possible, data type conversion will have to be handled.

Lets take the integer representation as an example. In C, and integer can vary between 2, 4 and 8 bytes in length depending on the architecture and compiler. To solve this, standard integer types

are used (e.g. `uint8_t` & `int32_t`) which guarantee that the representation is at least the length of the type definition stated. In Java, an integer is guaranteed to be 4 bytes long regardless of the system. However, DPDK uses unsigned integers while Java uses signed integers. This requires conversion between the unsigned and signed, but since unsigned has a higher upper bound on the value it can store due to the extra bit (MSB) there can be an overflow error when converting to Java. This requires that Java uses a long (8 byte) representation to hold the C 4 byte unsigned integer.

Conversion from Java to C then could then result in an underflow error if a value which can be represented by a Java long can't be represented by a C integer. This means bound checking is required on the Java side for any number conversions.

get a table
of this

5.7 Protocol undertaking

just IP (4 and) supporting udp and tcp (i think), no ARP etc

5.8 config files

6 DPDK Java

Since the Data Plane Development Kit is a framework, it was decided that instead of implementing separate middleboxes, a Java framework would be implemented. This would allow any number of other applications to be designed with ease.

mention applications below for demonstrations of ease of use of framework, also mention only subset of availability of dpdk, and mention further extensions easy to do like packet fragmentation etc

important to check for errors at every stage to stop jvm crashes

6.1 Overview (better name?)

As proven earlier in the report, for fast packet processing in Java the number of JNI calls should be minimised as possible due to their large overhead. However, since DPDK is a native library some JNI calls were obviously mandatory, mainly those initialising the application and packet interactions. This was minimised by making the majority of the JNI calls at the start of the application, before individual processing threads had been started. From then on, JNI calls would only be made at vital times within the application.

6.2 Native Libraries

For the Java framework to work, it requires a native shared library which includes all of the native methods to be called. This shared library can be placed anywhere as long as the correct Java path is set to its location, however, the common place is to install it within the users library directory (/usr/lib/) in Linux.

This shared library must also be dynamically linked to other shared libraries which contain the other non-JNI methods called, whether this is system calls, standard system libraries or the actually compiled DPDK libraries. For this, there a number of tools which can be used, all with their own advantages and disadvantages.

6.2.1 Library Compilation Tools

GNU Linker The GNU linker on linux is a basic linker which can be used with any number of shared libraries, although as the number increases so does the number of number of commands needed, as well as the number of flags needed. Considering the large amount of shared libraries which DPDK uses, this option is less attractive.

GNU Libtool The GNU Libtool uses the GNU linker underneath but abstracts away the compilations into a build unit, similar the makefiles used. It also has the advantage of only recompiling those libraries required if certain source files are updated. However, it does require a significant amount of set-up, and problems can occur when migrating to other machines.

DPDK Makefile By far the best option is to use the DPDK makefile. It is generally used when the application is written in C and uses the DPDK libraries which is why a lot of documentation is provided on this. However, when wanting to compile shared libraries (which don't include a main method) there is very limited documentation on this and certain flags and processes can only be found by delving into the build process of DPDK. Even so, this is the easiest option to use on a regular basis and ports perfectly well to other systems as it installs the shared library for you as well.

6.3 Initialisation

DPDK requires a number of initialising procedures to create the environment abstraction layer, start ports, allocate memory pools and setup the port specific queues. Each of these procedures is directly mapped to a Java implementation of them, which also handle error checking. Generally any native applications using DPDK create a number of threads depending on the number of available cores to split the work and maximise throughput. However, since the Java side creates its own affinity threads for this, DPDK is always initialised with only 1 thread to be used for the set-up.

When writing custom applications, all initialising should be complete before starting the threads. This is to stop any DPDK errors when trying to access uninitialised ports, memory pools or queues. For this reason, the last thing which should be called is 'startAll()' which starts all threads and affinity threads one by one in quick succession to allow the actual packet processing to take place.

6.4 Processing Threads

The framework supports multiple processing procedures which can run simultaneously and even with different implementations depending on the requirements. This allows an application to be built where certain processing objects do the polling of packets, and then pass these objects to be inspected, and then onto another object for sending and freeing of the packets. This basically creates a pipelined application, although a more simple processing object would receiver, process and forward packets in the same thread. This is instead of multiple processing units doing exactly the same job.

Each of the processing units are automatically threaded by the framework and put into an 'Affinity Thread'. This follows the same logic of the DPDK framework where each thread is set to only run on 1 core of the machine's processor using affinity cores . This provided a few complications as Java and the JVM doesn't provide functionality for assignment of threads to cores. This is mainly as because JVM abstracts away the complications of this and allows the kernel to do its own thread scheduling. However, on Linux, Java does utilise the native POSIX Threads (pthread) and assigns a Java thread to 1 pthread. This meant via a few JNI calls and native system calls , each thread could in fact be associated with certain cores. As with DPDK, this limits the number of threads to be equal or less than the number of available cores on the machine (of hyper-threaded cores) as to fully maximise the application speed.

Each of the processing units must be an extension of the 'PacketProccesor' class which provides an underlying abstraction so statistics of packet data can be retrieved as mentioned below.

ref this
and say
thread can
be set to
run on
multiple
cores -
need to
fix code
for this

example?

Describe 3 objects and fix code to use more than 1

6.5 Packet Data Handling

DPDK and therefore the Java framework version is primarily used for packet processing, meaning that efficient handling of the packet's data and header information are a necessity. Investigations outlined in section [supported that copying packet data from the native side to Java and back](#) again was a very poor choice in terms of processing speed. To solve this, all packet information is left in native memory (not copied to the Java heap) and accessed directly from the Java application in order to read and write data to/from specific packets.

ref this

In order to do this, the Java Unsafe class was used extensively as it provided then functionality. However, as discussed previously, the Unsafe class directly accesses native memory and is therefore inherently unsafe to use, as opposed to Java itself. This memory accessing was therefore abstracted away into the 'UnsafeMemory' class which handled the type conversion between native unsigned and Java signed, pointer arithmetic and big/little endian conversion. It also checks whether values are going to be out of data type number representation range when converted from signed to unsigned, since unsigned values are represented as larger data types to account for the signed bit.

Since little-endian to big-endian conversion and vice-versa was required to pass data between the JVM and native memory it was decided that all handling of this would be done on the native side. This firstly abstracts the certain complications of this away from the Java side to provide a cleaner interface and secondly the native implementation provides faster byte shifting. To accomplish this, whenever data was read into a memory location assigned from Java Unsafe, it would flip the data's byte order and put the data into big-endian format for the JVM to interpret correctly. Whenever data was been passed from Java to the native side, the opposite would occur .

elaborate on this, diagram? how did pointers work?

Each individual network packet is assigned its own Packet object representation. Even though object usage on the Java heap can be relatively slow compared with native structure handling due to various reasons, it was required or else the application would be wondering away from the object orientated side of Java. However, the Packet object only tracks 3 fields in order to minimise data copying:

- protected long mbuf_pointer - points to memory location of the start of mbuf header for the packet
- protected long packet_pointer - points to memory location of the start of the packet header (either IPv4 or IPv6)
- protected UnsafeAccess ua - packets own unsafe memory accessing object for secure navigation around packet data

Since the native mbuf and ipv4/ipv6 structures are forced to be packed, [this means that all fields can be traversed and therefore read from and written to via pointer arithmetic](#). The packet_pointer allows for simple getters and setters which relate to the standard IP packet headers, even if the actual fields don't exist. The mbuf_pointer is generally used to gain access to the raw packet data and used later for freeing and forwarding the packet.

ref this

6.6 Packet Polling

Packet polling is the act of receiving data from the memory buffers used to store packets received by the network interface card/controller. When an application requests new packets from the received queue of a specific port, it does so using native methods via the JNI. These are the only JNI calls used after the initial start up (i.e. used in processing threads) in order to increase performance. The native methods' parameters specify the port id, queue id and a memory location pointer, while defaulting to fetching the default number of packets set on the application initialisation. Before this call happens, Java DPDK allocates a set number of bytes within the native memory and holds its pointer in memory in long format. This pointer location is then passed as one of the parameters to the native method.

get image of setup of memory location, ie number then list of packet pointers

Within the native method, the number of received packets is put as a short data type into the memory location of the pointer. Every subsequent value inserted into the memory location is first a pointer to the mbuf location of a given packet and then a pointer to the location of the packet header. This happens for all of the packets received by the call.

Within the Java side, using the UnsafeMemory class, the number of packets is pulled from the memory location. From there, the pointers can be pulled within a loop and added to their own packet class, depending if they are IPv4 or IPv6 headers. A list of packets is then returned from the method.

6.7 Packet Sending

Packet sending involves taking processed packets and putting them on the port queues ready to be sent by the network interface card/controller. When ever a packet needs to be sent, it is passed to a PacketSender object which initiates the process. Within PacketSender, packets are stored within a list ready to be sent. Since packet sending within DPDK is more efficient when sending multiple packets at the same time, whenever the list reaches the default value for a sending burst (n), the first n packets in the list are sent. However, there is also a time-out functionality which is used whenever the receiving of packets is slow. Since the program doesn't want to wait for n packets which could take a long time, the time-out allows any number of packets to be sent if a time reaches a user set time period away from the last send burst.

Sending of packets is very similar to the polling of packets but in the opposite direction. Firstly a memory location of the required size (depending on the number of packet to send) is allocated in native memory, and that pointer, along with port id and queue id are passed as parameters to the native code. Each packets' mbuf pointer is then added to the memory location in increasing order. This pointer is then directly passed as a parameter to the DPDK packet sending methods as the memory is already set up in the correct format.

ref this

Once the packets have been put onto the queue, a loop is made through the pointers to the mbuf's to free the memory.

6.8 Statistic Profiling

The statistic profiling is an optional part of the framework which allows statistics to be gathered on the number and sizes of packets received and packets sent. It gathers data every second by default, but can be user set to find information about the number of packets sent/received for a given time unit. It's up to the user to set which ReceivePoller and PacketSender objects to gather information on. In order to be most efficient for the application, it runs on its own thread, although not an affinity thread, so it can be context switched into any of the cores.

It also has the option of running a graphical user interface (GUI) which displays the data in better format instead of constant prints to the console. It also redirects other console information on the running of the application to a user console on the GUI for easier use as DPDK outputs a lot of information which isn't that useful.

put im-
age of gui,
also write
output to
files

6.9 Testing

7 Applications

Since a Java DPDK framework was designed, any applications are relatively easy to create and only use Java code so no interaction with the native DPDK is necessary from a users point of view. Below are 2 simple applications created using this framework, which both show the small amount of code required to make a working application. However, either of these could be heavily extended to use multiple threads, different processing objects, shared ports in order to design more complex application such as packet fragmentation, reassembly, multi-casting or ordering.

7.1 Firewall

7.2 Network Address Translator (NAT)

8 Evaluation

This section focusses on the comparison testing which was carried out between the C and Java implementations of basic middleboxes. Evaluation of the results will then be discussed and further improvements to the Java solutions will be considered.

mention
machine
specs

8.1 Packet Generating

Packet generating is the act of creating packets with random payloads to be sent to certain MAC addresses on the network. This can either be done via the use of specialised hardware or using software. They are used for load testing of packet processing applications to test the amount of data which applications can process per second. This can reveal whether limitations on a system is software or hardware based.

Talk about pktgen module in kernel - does dpdk version use this

Pktgen is open source software tool, maintained by Intel, which aims to generate packets using the DPDK framework. It can generate up to 10Gbits of data per second with varying frame sizes, and send the data in the form of packets across a compatible network interface card/controller. It has a number of benefits which include:

- Real time packet configuration and port control
- Real time metrics on packets sent and received
- Handles UDP, TCP, ARP and more packet headers
- Can be commanded via a Lua ?? script

how does pktgen work and why we used it, only does 32 packet bursts or less

Pktgen's transmitting performance can be reduced when increasing the frame size to those greater than 64 bytes. This can be negated by running the application with multiple lcores on the same port with different queues, which allow the transmitting speed to match those of the NIC. With the packet size varying, obviously the number of packets transmitted per second is depends on this. The differences are shown in Table 3.

Packet Size	Packet/s (millions)	MBit/s
64 (min)	14.9	9999
128	8.4	9999
256	4.5	9999
512	2.3	9999
1024	1.2	9999
1518 (max)	0.8	9999

Table 3: Pktgen speed at 32 packet burst, used 2 processors to get meet NIC speed

8.2 Initial Testing

The initial testing of applications was carried out on a local Mac OS X machine running Ubuntu 14.04 LTS 64-bit on a VirtualBox virtual machine. Although this set-up didn't provide the ability to load test on very high speeds (anything above 1Gbit/s), it allowed for basic testing to check that the application was running as expected. Load testing of speeds up to roughly 700Mbit/s were also possible which have a basic testing platform without the need to move code to servers.

which ones?

ref this and ubuntu

8.2.1 Set-up

Testing could be carried out using the 2 available 1Gbit NICs of the machine via a bridged network from the host to guest machine which severely reduced transmission speed. This allowed an ethernet cable to be looped back and connected between the ports, meaning anything transmitted via 1 port was guaranteed to be received by the other port.

Pktgen and the custom application were booted up simultaneously running in parallel. Careful memory allocation, port addressing and processor core assignment had to be carried out to stop shared resources impacting the overall performance of either application. This allowed Pktgen to send packets and the application to receive and process them.

8.2.2 Methods

8.2.3 Results

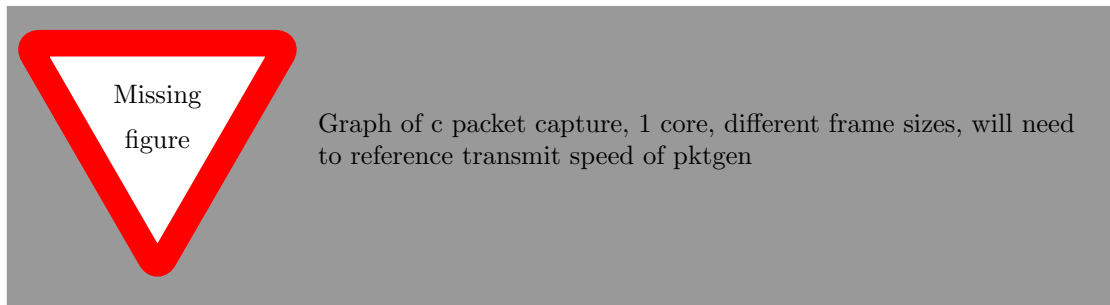
8.3 Further Testing

In order to fully understand the capabilities of the middle boxes, testing was carried out on Imperial College's Large Scale Distributed Systems (LSDS) test-bed. Although this system consists of numerous machines, tests were carried out using just 2. The first machine was used to host the middlebox application and receive the packets. The other machine was used as the client and ran the pktgen software allowing it to generate packets at up to 10Gbps in order to take advantage of the machines network interface controllers.

The first test on the LSDS test-bed involved comparing the C and Java implementations of a packet capturing application which simply received the packets and freed them straight away without forwarding them. The C implementation is used to give the optimal readings possible from this and further tests, since very limited processing is carried out between receiving and dropping the packet. The figure below shows the performance of the application at varying packet sizes.

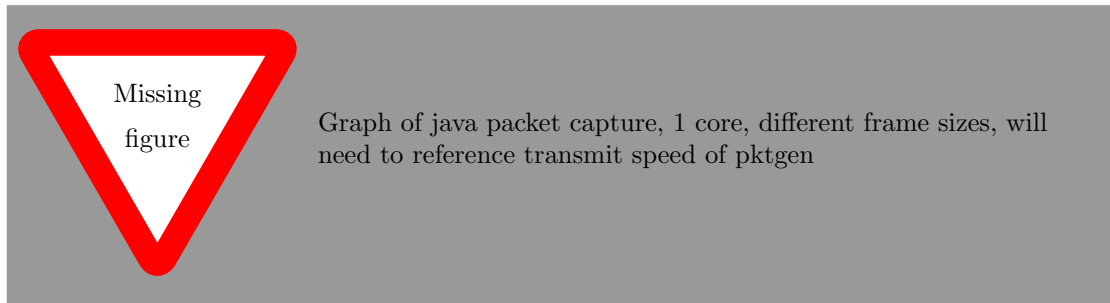


Figure 17: sdfd



The same testing was carried out with an identical packet capture algorithm which was instead implemented in Java using the DPDKJava framework described previously. These results are shown below.

what does the results show



The results above are very low to what was expected for such a simple application.

expand on this

To check where the problems lied within the code, a Java profiler (JProfiler) was used to check multiple parts of the code including memory usage, cpu usage and the number of method calls and the average time per method call. This provided invaluable analysis of where the problems were, although the profiler significantly reduced the performance of the application since it connects to the JVM and reads the data itself.

Figure ?? shows some of the output from the profiler which indicated where the main problems were in terms of memory usage and performance. From this, a number of performance upgrades were implemented:

Capture Processor fields The first improvement were to the Packet Capture application itself. The associated processor originally stored the ReceivePoller and PacketFreeer objects within a list to allow for easy iteration if there were numerous objects. Since the Packet Capture application only used 1 of each, seperate fields could be used for the objects which removed the list accessing times and iterations.

Object lifespan The other improvements were made to the actual DPDKJava framework. These involved utilising objects throughout the application lifespan instead of creating new ones on every loop. This dramatically reduced the number of initializing methods invoked for the objects and reduced the memory usage in the heap, which reduced the number of times the Java garbage collector was invoked. The class which caused the most problem with this was the ArrayList, which were created on every loop to pass packets through the Java system. Since the framework uses threads without the need to synchronise objects, an ArrayList could be created on initialisation and simply cleared before been used again.

Off heap allocation Finally, instead of allocating new off heap memory to receive the packet pointers through on every loop, a memory bank was allocated on initialisation and the same memory was simply overwritten on every loop iteration. These improvements resulted in the graph below.

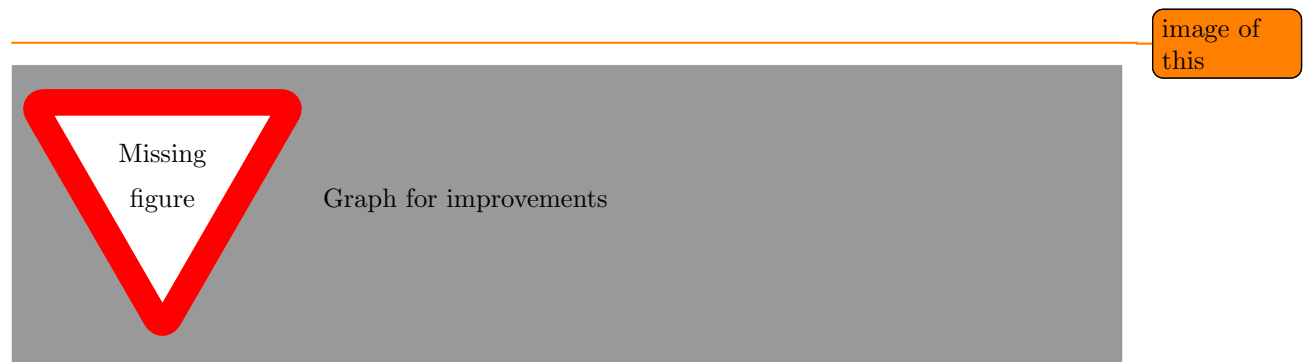
Again, the same tests were carried out to check the performance increase of the framework implementation.



explain re-
sults com-
pared to
previous
ones and c
version

Packet creation Of course, for every packet entering the application a packet object is created for easy referencing further down the application pipe-line. However, not using packet objects would significantly reduce the usability and scalability of applications. It was decided not to alter this. However, for each packet initialization a new UnsafeAccess object was been created for use with accessing packet header information. However, since each thread controls its own packets and therefore can only process 1 packet at a time, 1 UnsafeAccess object could be shared between all packets which significantly reduced the number of objects on the heap.

Send and Free lists Further problems were identified with the lists of packets awaiting to be sent and freed. This list was been iterated over with the Packet's mbuf pointer then been stored in an off heap memory bank waiting to be freed. This was pointless since the packets mbuf pointer could directly be put into the off heap memory, therefore eliminating the need for the list while also allowing the packet objects to be dereference quicker.



The results show that the performance was further improved and pretty close to that of the C implementation. After further profiling, there was only 1 obvious improvement which could be made, This would be to replace the list storing packets received from the poller. However, replacing this with a custom implementation using off heap memory would firstly reduce the usability of the code and would also mean that the code was drifting away from the Java language. It was decided not to implement this fix and instead to continue with testing of other applications, assuming that the would be negated by applications. This initial series of testing also proved that a dramatic increase in performance can be achieved simply by programming the applications efficiently, by trying to reduce the number of objects created.

8.3.1 Set-up

8.3.2 Methods

8.3.3 Results

8.4 Software Design

Mention somewhere about the limitations of pktgen

8.4.1 Portability

8.5 Possible Improvement

9 Conclusion

9.1 Future Work

9.2 extrapolate for other non-native languages

A number of other programming languages such as Visual Basic (.NET) and C# use similar a virtual machine just like Java to make them portable to other machines. This gives the opportunity to extrapolate certain high-level implementation details and apply them to other non-native programming languages, allowing for other application not written in Java to take advantage of fast non-native packet processing.

continue
this, where
should it
go?

10 User Guide

should I be putting all data in appendix

<http://plvision.eu/blog/deploying-intel-dpdk-in-oracle-virtualbox/> <https://www.virtualbox.org/manual/ch08.html>