



## Accessing InfiniBand hardware from Java using the Foreign Function & Memory API

---

Filip Krakowski, Fabian Ruhland, Michael Schöttner

Department of Computer Science  
Heinrich Heine University Düsseldorf

18/04/2023

# Motivation

---

Many data processing and computation frameworks are written in the Java programming language or run on the Java Virtual Machine.

**Spark™**

**Flink**

**Storm**

**Samza**

**Mahout**

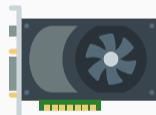
Using traditional sockets these frameworks produce computational overhead regarding the transport layer due to **context switches** / **system calls** and **buffer copies**.

**InfiniBand** eliminates these problems by allowing applications to **bypass the kernel** and read/write remote memory directly in a zero-copy fashion.

Remote direct memory access (**RDMA**) is not limited to Random Access Memory

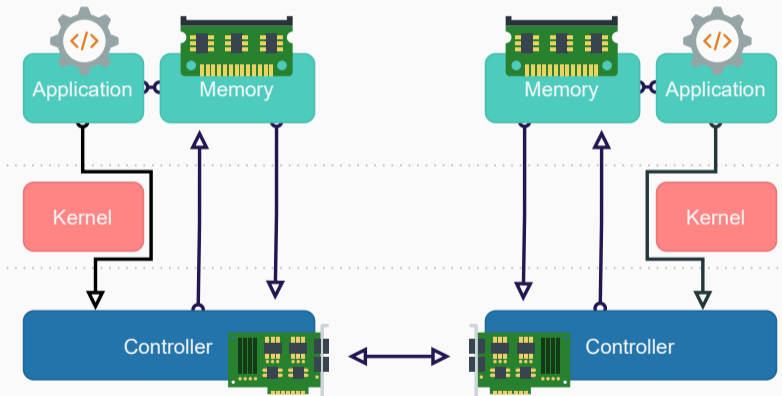


NVMe



GPU

Unlike Ethernet, the InfiniBand protocol stack is implemented within hardware



# Kernel Bypass

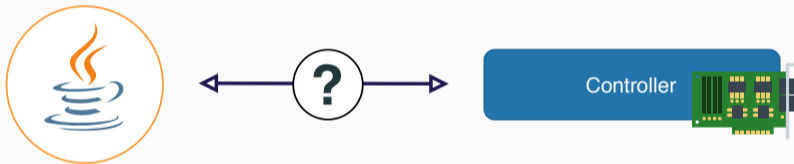


**Low latencies**  
(  $< 0.7 \mu\text{s}$  )



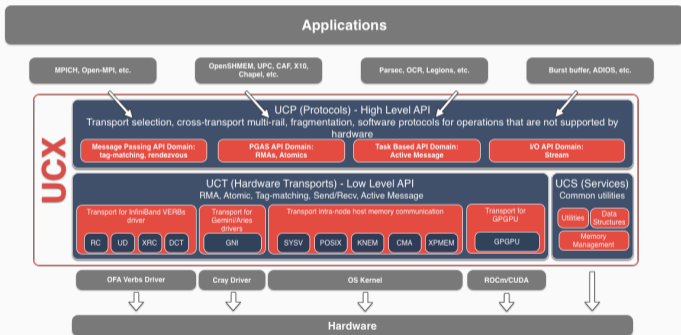
**High throughput**  
( 400 Gb/s )

Current JDK implementations do not provide support for InfiniBand hardware



A library is required for interfacing with InfiniBand hardware

The *Unified Communication X (UCX)*<sup>1</sup> Library offers **simple access** to InfiniBand hardware, provides **easy-to-use abstractions** and is **written in C**.



<sup>1</sup><https://github.com/openucx/ucx>



# Native Interface

---

Two options for interfacing with native code from Java

## Java Native Interface <sup>2</sup>

(since JDK 1.1 - 1997)

## Project Panama <sup>3</sup>

(since JDK 16 - 2021)

*“ We are improving and enriching the connections between the Java virtual machine and well-defined but "foreign" (non-Java) APIs, including many interfaces commonly used by C programmers.”*

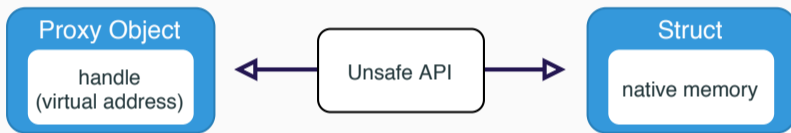
— Oracle, *Project Panama* <sup>3</sup>

---

<sup>2</sup><https://docs.oracle.com/javase/8/docs/technotes/guides/jni>

<sup>3</sup><https://openjdk.java.net/projects/panama>

In our previous project **neutrino**<sup>4</sup>, we explored creating so called "*Proxy Objects*" for accessing structs in native space.



Using `sun.misc.Unsafe`<sup>5</sup>, these proxy objects can write and read to their associated structs fields (in off-heap memory) directly.

<sup>4</sup><https://github.com/hhu-bsinfo/neutrino>

<sup>5</sup><http://www.docjar.com/docs/api/sun/misc/Unsafe.html>

Proxy objects create a mapping for struct fields using their names.

AddressHandle.java

Java

```
1 @LinkNative("ibv_ah")
2 public class AddressHandle extends Struct {
3     private final Context ctx = referenceField("context");
4     private final ProtectionDomain pd = referenceField("pd");
5     private final NativeInteger handle = integerField("handle");
6 }
```

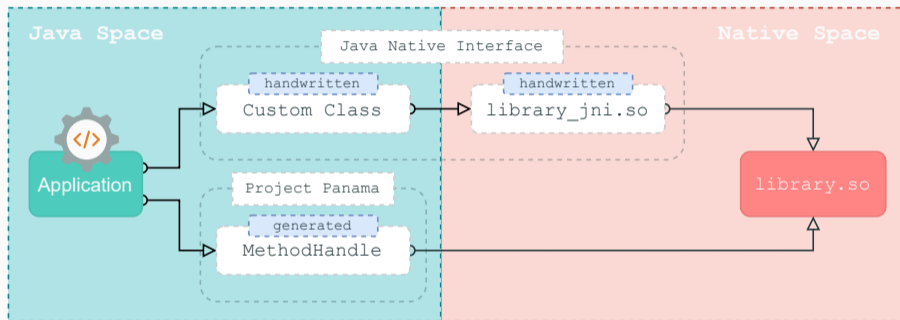
verbs.h

C

```
1 struct ibv_ah {
2     struct ibv_context *context;
3     struct ibv_pd *pd;
4     uint32_t handle;
5 }
```

A lookup table must be managed **by hand** in native space.

The Java Native Interface requires handwritten "glue code" <sup>6</sup>.



Project Panama automatically generates bindings through jextract <sup>7</sup>

<sup>6</sup>[https://en.wikipedia.org/wiki/Glue\\_code](https://en.wikipedia.org/wiki/Glue_code)

<sup>7</sup><https://github.com/openjdk/jextract>

We provide a Gradle plugin <sup>8</sup> to automate jextract's process

**build.gradle**

Groovy

```
1 plugins {
2     id "io.github.krakowski.jextract"
3     version "0.3.1"
4 }
5 jextract {
6     header("${project
7         .projectDir}/src/main/c/stdio.h")
8     {
9         libraries = [ 'stdc++' ]
10        targetPackage = 'org.unix'
11        className = 'Linux'
12        functions = [ 'printf' ]
13    }
14 }
```

**NativeHelloWorld.java**

Java

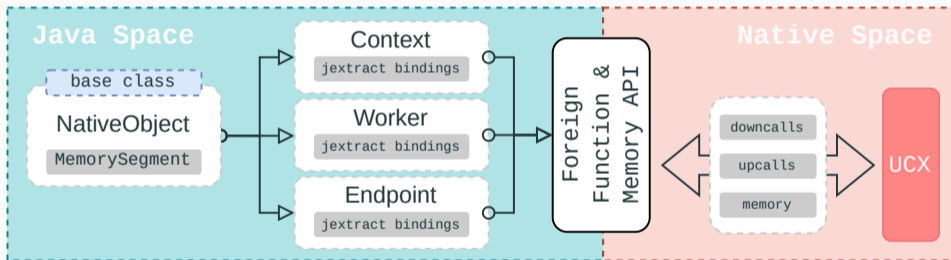
```
1 import static org.unix.Linux.*;
2
3 public final class NativeHelloWorld {
4
5     public static void main(String... args) {
6         try (var session =
7             MemorySession.openConfined()) {
8             var format =
9                 session.allocateUtf8String("Hello %s");
10            var value =
11                session.allocateUtf8String("World");
12            printf(format, value.address());
13        }
14    }
15 }
```

<sup>8</sup><https://github.com/krakowski/gradle-jextract>

# Framework Design

---

Using the **Foreign Function & Memory API**<sup>9</sup>, we implement an object-oriented framework called **Infinileap**<sup>10</sup> wrapping the native UCX library.



<sup>9</sup><https://openjdk.org/jeps/424>

<sup>10</sup><https://github.com/hhu-bsinfo/infinileap>



## NativeObject.java

Java

```
1 public class NativeObject {
2     private final MemorySegment segment;
3     protected NativeObject(MemorySegment segment) { ... }
4     protected NativeObject(MemoryAddress address, MemoryLayout layout) { ... }
5     protected NativeObject(MemoryAddress address, long byteSize) { ... }
6 }
```

- All relevant UCX structs are wrapped inside a class having the same name.
- User **can** provide `MemorySession`

## RequestParameters.java

Java

```
1 public class RequestParameters extends NativeObject {
2     public RequestParameters() {
3         this(MemorySession.openImplicit());
4     }
5
6     public RequestParameters(MemorySession session) {
7         super(ucp_request_param_t.allocate(session));
8     }
9 }
```

UCX makes heavy use of flags within its provided structs.

InfiniLeap wraps these flags and encapsulates them in **enum values implementing a common interface**, so that they can be passed as **varargs** to our API.

## LongFlag.java

Java

```
1 @FunctionalInterface
2 public interface LongFlag {
3     long getValue();
4 }
```

## RequestParameters.java

Java

```
1 public enum Field implements LongFlag {
2     CLIENT_ADDR(UCP_CONN_REQUEST_ATTR_FIELD_CLIENT_ADDR()),
3     CLIENT_ID(UCP_CONN_REQUEST_ATTR_FIELD_CLIENT_ID());
4
5     private final long value;
6     Field(int value) { this.value = value; }
7
8     @Override
9     public long getValue() { return value; }
10 }
```

Callbacks / Upcalls are provided as **functional interfaces**, which the user may implement. The framework calls `upcallStub` and stores the `MemorySegment` reference to prevent garbage collection.

SendCallback.java

Java

```
1 @FunctionalInterface
2 public interface SendCallback extends ucp_send_nbx_callback_t {
3     void onRequestSent(long request, Status status, MemoryAddress data);
4
5     @Override
6     default void apply(MemoryAddress request, byte status, MemoryAddress data) {
7         onRequestSent(request.toRawLongValue(), Status.of(status), data);
8     }
9
10    default MemorySegment upcallStub() {
11        return ucp_send_nbx_callback_t.allocate(this, MemorySession.openImplicit());
12    }
13 }
```

UCX enables atomic operations (e.g compare and swap) on remote memory. InfiniLeap abstracts this process by implementing so called "Native Primitives" which provide typed (`int`, `long`, ...) access to the underlying memory.

NativeLong.java

Java

```
1 public final class NativeLong extends NativePrimitive {
2     private static final int SIZE = Long.BYTES;
3     public NativeLong() { this(MemorySession.openImplicit()); }
4     public NativeLong(MemorySession session) { this(0, session); }
5     public NativeLong(long initialValue, MemorySession session) {
6         super(MemorySegment.allocateNative(SIZE, session), DataType.CONTIGUOUS_64_BIT);
7         set(initialValue);
8     }
9
10    private NativeLong(MemorySegment segment) {super(segment, DataType.CONTIGUOUS_64_BIT);}
11    public void set(long value) {segment().set(ValueLayout.JAVA_LONG, 0L, value);}
12    public long get() {return segment().get(ValueLayout.JAVA_LONG, 0L);}
}
```

Native Primitives can be compared to the JDK's `Atomic`{`Integer`, `Long`, ...} classes.

AtomicAddExample.java

Java

```
1 MemoryDescriptor descriptor = /* Received from other network participant */
2 RemoteKey remoteKey = endpoint.unpack(descriptor);
3
4 // Create a memory segment for atomic operations
5 var memorySegment = context.allocateMemory(Long.BYTES);
6 var nativeLong = NativeLong.map(memorySegment.segment());
7 nativeLong.set(32);
8
9 long request = endpoint.atomic(AtomicOperation.ADD, nativeLong,
    descriptor.remoteAddress(), remoteKey, new
    RequestParameters().setDataType(nativeLong.dataType()));
10
11 Requests.await(worker, request);
```

## Sender.java

Java

```
1 final var buffer = MemorySegment.allocateNative(64L, MemorySession.openImplicit());
2
3 long request = endpoint.sendTagged(buffer);
4 Requests.await(worker, request); // busy-spin until request finishes
5 Requests.release(request); // clean up request
```

## Receiver.java

Java

```
1 final var buffer = MemorySegment.allocateNative(64L, MemorySession.openImplicit());
2
3 long request = worker.receiveTagged(buffer);
4 Requests.await(worker, request); // busy-spin until request finishes
5 Requests.release(request); // clean up request
6
7 MemoryUtil.dump(buffer, "Buffer"); // Print content
```

Full examples are available in our GitHub repository <sup>11</sup>

<sup>11</sup><https://github.com/hhu-bsinfo/infinileap/tree/master/example>

# Evaluation

---

<b>CPU</b>	Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz (22 MB Cache)
<b>RAM</b>	4x Micron Technology 36ASF2G72PZ-2G6E1 16GB
<b>NIC</b>	Mellanox Technologies MT27800 Family [ConnectX-5] (100Gbit/s)

**Figure 1:** System specifications of the hardware used in all experiments.

All benchmarks are implemented using the *Java Microbenchmark Harness (JMH)*<sup>12</sup> Framework. The benchmark's source code is available on GitHub<sup>13</sup>

All measurements shown are average values with standard deviation error bars

---

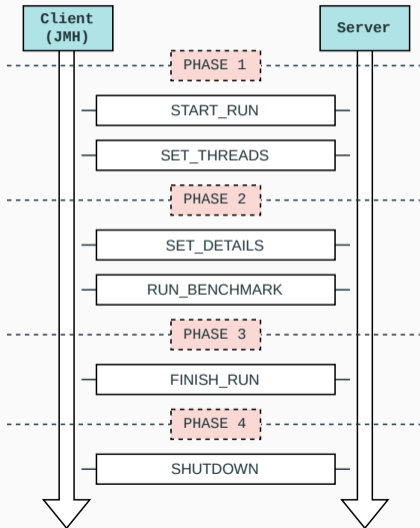
<sup>12</sup><https://github.com/openjdk/jmh>

<sup>13</sup><https://github.com/hhu-bsinfo/infinileap>



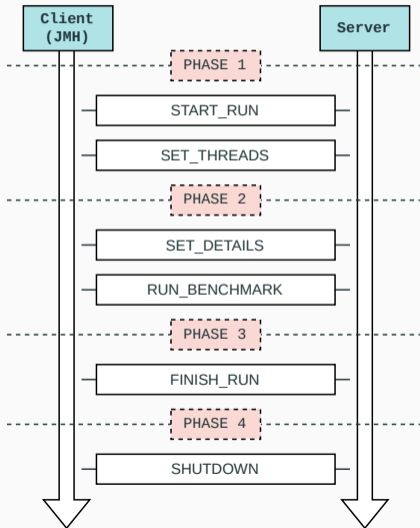
**PHASE 1** Initialize the server instance and set the number of threads.

1. Send a `START_RUN` command to tell the server to start the next or first run.
2. Instruct the server to start a specified number of worker threads by sending a `SET_THREADS` command.



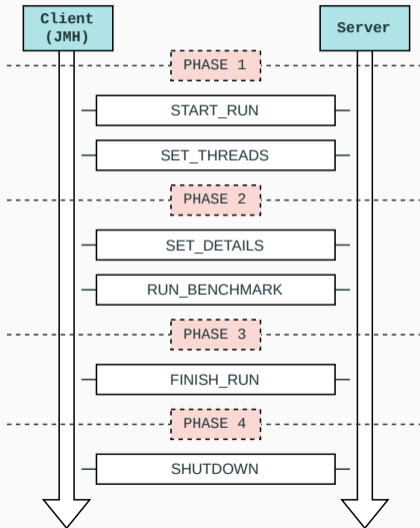
**PHASE 2** JMH invokes benchmark methods and makes measurements.

1. Inform the server about the benchmark details (operation count, buffer size, etc.) by sending a SET\_DETAILS command.
2. Execute the specified number of benchmark method invocations until a configured time has expired.
3. Synchronize with the server to let it receive new commands.
4. Repeat PHASE 2 until all configured buffer sizes have been measured; else go to PHASE 3



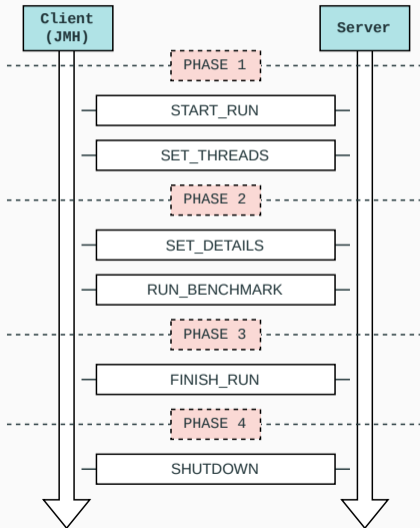
## PHASE 3 - Resources are released.

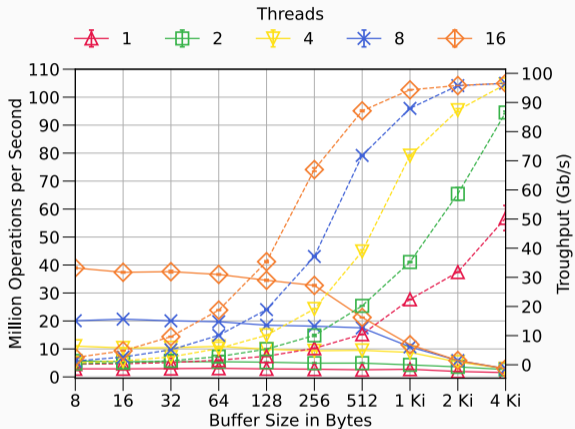
1. Send a `FINISH_RUN` command to tell the server to release its resources and terminate all worker threads.
2. Start a new run with a different number of threads by reentering PHASE 1 or move on to PHASE 4.



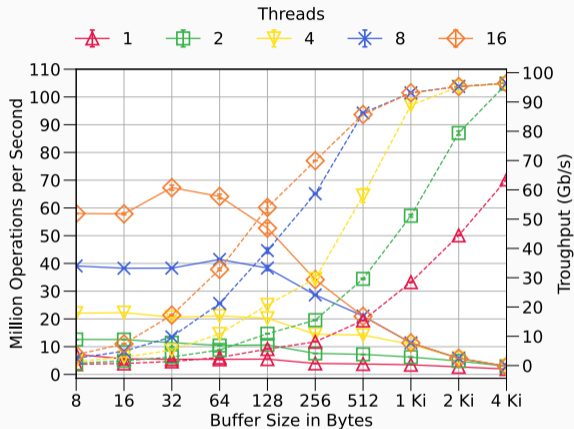
## PHASE 4 - Benchmark ends.

1. Send a SHUTDOWN command to tell the server it should terminate. Alternatively another benchmark run may be started by sending a START\_RUN command and starting again from the beginning.

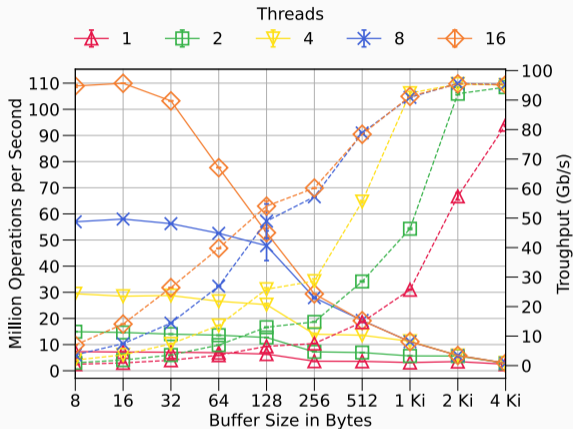




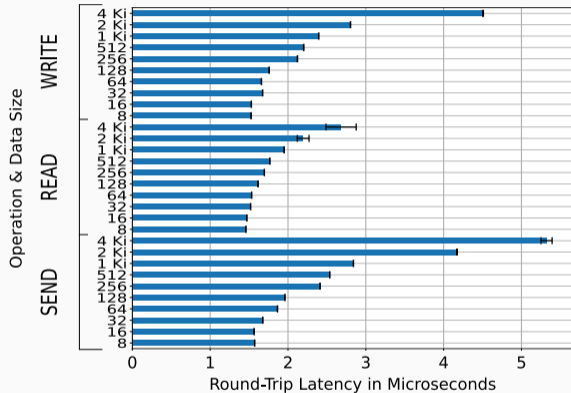
**Figure 2: Average read operation (solid line) and network (dashed line) throughput**



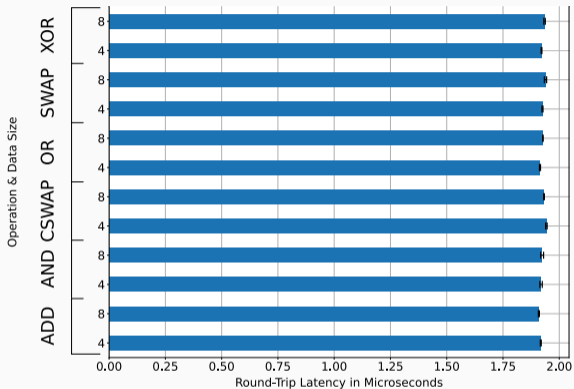
**Figure 3: Average write operation (solid line) and network (dashed line) throughput**



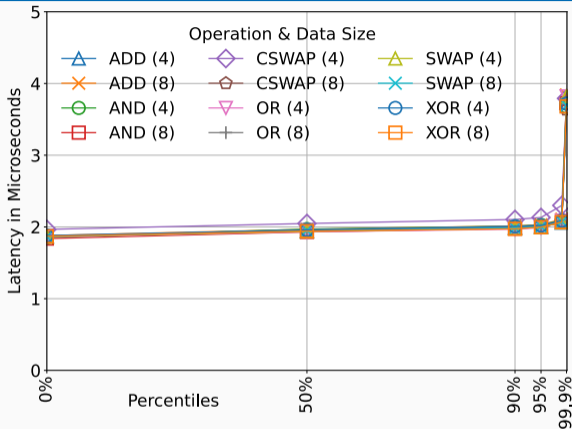
**Figure 4:** Average **send** operation (solid line) and network (dashed line) throughput



**Figure 5:** Average Round-trip latency for **RDMA write, RDMA read and send** operations



**Figure 6:** Average operation latency for all UCX-supported atomic operations and data sizes



**Figure 7:** Operation latency for all supported atomic operations by percentiles

## **Conclusion & Lookout**

---



- Project Panama's `jextract` tool enables easy integration of existing C libraries (direct call vs. glue code).
- Using the Foreign Function API in the "hot path" yields good results regarding performance.
- Thanks to the safety features of the Foreign Memory API, "off-by-one" errors are detected earlier in the development cycle.
- Higher-Level frameworks like Netty<sup>14</sup> can be accelerated using InfiniLeap.

---

<sup>14</sup><https://netty.io/>

# Questions