



MIC-MPI-lab: Simple MPI programs and Intel MPI Benchmarks in various modes on the MIC

In this lab you will gain first experiences with the various MPI execution scenarios supported by the Intel compiler for the Intel Xeon Phi architecture: Native mode, MPI ranks on hosts and/or coprocessors, MPI ranks on hosts which do offloading etc. We use a simple hello world program to understand the various execution modes. The second part of the exercises uses the Intel MPI Benchmark suite to do some performance measurements on SuperMIC.

Appropriate Environment

```
module load intel
export I_MPI_MIC=enable
export I_MPI_FABRICS=shm:dapl
export I_MPI_DAPL_PROVIDER_LIST=ofa-v2-mlx4_0-1u,ofa-v2-scif0,ofa-v2-
mcm-1
export MIC_LD_LIBRARY_PATH=$MIC_LD_LIBRARY_PATH:$I_MPI_ROOT/mic/lib/
```

Lab 1: MPI on the hosts

- Compile the simple hello world program `hellompi.c` using
`mpicc hellompi.c -o hellompi-host`
- Run 2 MPI ranks on a compute-host using
`mpirun -n 2 -host r23u18n627./hellompi-host`
- Modify the code that every MPI process outputs its rank, the total number of ranks and the hostname it is running on.

Lab 2: MPI on a single MIC in native mode

- Compile the code for the Intel MIC architecture using:
`mpicc -mmic hellompi.c -o hellompi-mic`

- Login to the MIC and run MPI natively, e.g.

```
export PATH=$PATH/apps/all/impi/2017.0.098-iccifort-  
2017.0.098-GCC-5.4.0-2.26/mic/bin/  
mpirun -n 10 ./hellmpi-mic
```

Lab 3: MPI on the MICs launched from the host

- Launch many MPI tasks of the program from the host on one single MIC using, e.g.

```
mpirun -n 5 -host r22u07n620-mic0 program
```

- Run one MPI task per MIC on the 2 MICs attached to one single host using

```
mpirun -n 2 -perhost 1 -host mic0,mic1 program
```

Lab 4: MPI on the host and its 2 attached MICs

- Run 2 MPI tasks on 1 host and 2 tasks per MIC on both MICs attached to the host using

```
mpirun -n 2 -host hostname hostprogram : -n 2 -host  
hostname-mic0 micprogram : -n 2 -host hostname-mic1  
micprogram
```

- If enough nodes are free, run MPI tasks on multiple hosts and MICs.

`qstat -f` provides you the names of the allocated hosts.

Lab 5: MPI on the host using Offloading

- Modify the code so that an offload region with another simple “hello world from MIC” is offloaded to the MIC.
- Run multiple MPI ranks just on 1 host and verify that the offload region runs on a MIC.
- Modify the file so that the code offloaded to a MIC tells the hostname of the coprocessor and the hostname and rank of the MPI process on the host which has done the offload.
- Offload MIC code from multiple MPI processes on the host to its 2 MICs, ensure that not only mic0 is used.

Lab 6: MPI Performance

In this Lab you use the Intel MPI Benchmark (IMB) Suite to compare the performance of various MPI communication patterns.

On the MIC the benchmark binary is available under `/apps/all/impi/2017.0.098-iccifort-2017.0.098-GCC-5.4.0-2.26/mic/bin/`.

On the host it is available under `$I_MPI_ROOT/bin64/IMB-MPI1`

To propagate the PATHS properly to the MIC coprocessors, you can use `export`

```
MIC_LD_LIBRARY_PATH=$MIC_LD_LIBRARY_PATH:/apps/all/impi/2017.0.098-iccifort-2017.0.098-GCC-5.4.0-2.26/mic/lib
```

```
mpirun -genv PATH /apps/all/impi/2017.0.098-iccifort-2017.0.098-GCC-5.4.0-2.26/mic/bin/ -genv LD_LIBRARY_PATH $MIC_LD_LIBRARY_PATH ...
```

- What is the maximum PingPong bandwidth between the host and 1 MIC? Use a command similar to

```
mpirun -n 1 -host hostname IMB-MPI1 PingPong : -n 1 -host hostname-mic0 IMB-MPI1 PingPong
```


Compare with the max. PCIe2 speed!
- How about the latency? Look at `t[usec]` for small `#bytes`.
- Also measure the PingPong Bandwidth between the 2 MICs attached to one host, between 2 hosts, and between 2 MICs attached to different hosts. If you are familiar with `gnuplot` or other plotting programs use them to compare the performance graphically.
- The recommended MPI Fabric is `I_MPI_FABRICS=shm:dapl`. to use DAPL via Infiniband. Change the MPI Fabric setting to use `tcp` for internode communication and compare the performance with MPI using DAPL.

Beyond PingPong the Intel MPI Benchmark suite also offers the following benchmarks:

```
# PingPong
# PingPing
# Sendrecv
# Exchange
# Allreduce
# Reduce
# Reduce_scatter
# Allgather
# Allgatherv
# Gather
# Gatherv
# Scatter
# Scatterv
# Alltoall
# Alltoallv
# Bcast
# Barrier
```

Look into the Intel IMB Docu <https://goo.gl/Yq8ESe> and the MPI Standard <http://www.mpi-forum.org/docs/docs.html> what these MPI functions do and compare the performance in native mode and using `host1↔host2`, `host1-mic0↔host1-mic1` etc. patterns.