# Working with Multiple MPIs
## Overcoming ABI Incompatibility

**Edgar León**
**Laurent Nguyen**
**Nathan Hanford**

Marc Joos
Eric Green
Adrien Cotte
Clément Fontenaille

21 May 2023

**IMAGINE TOMORROW**

ISC High Performance

MAY 21 – 25, 2023 | HAMBURG, GERMANY

**EOLEN** CONSEIL ET INGÉNIERIE
ALLIANCE SERVICES PLUS, UNE SOCIÉTÉ DU GROUPE EOLEN
AS+

CEA DE LA RECHERCHE À L'INDUSTRIE

**Lawrence Livermore National Laboratory**

---



Wrapper interface for MPI

| 14:00 | 14:35 | **Introduction and setup** |
|-------|-------|-----------------------------|
|       |       | **Changing MPIs dynamically with Wi4MPI** |
| 14:35 | 16:00 | • **Spack setup**<br>• **Preload and interface mode** |
| 16:00 | 16:30 | *Café* |
|       |       | **Applying Wi4MPI to HPC** |
| 16:30 | 18:00 | • **MPI for Python**<br>• **Mini-application exploration**<br>• **Podman** |

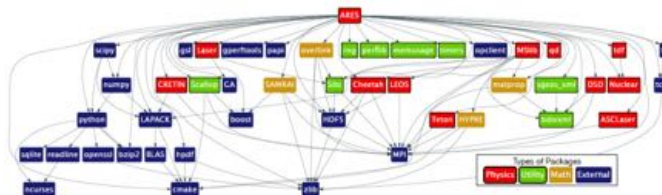https://github.com/LLNL/productivity-frameworks/tree/main/tutorials/isc23

## What if MPI applications could use any MPI library at runtime <u>without recompiling</u>?

- Working around limitations of an MPI library
  — Help diagnose the source of a problem
  — Choose the best MPI

- Enabling fast/portable containers
  — Containers provide flexibility and portability
  — Loss of portability to match the host MPI library

- Adding flexibility to high-level languages
  — High-level languages can depend on a specific MPI library

- Running on bleeding-edge/early-access systems
  — State-of-the-art systems may come with a single, vendor-optimized library

## An MPI library may use its own Application Binary Interface!

- MPI has a single API 😃

- MPI may have several ABIs 😫
  — OpenMPI, MPICH, MPC

- `typedef int MPI_Comm;`
  `#define MPI_COMM_WORLD ((MPI_Comm) 0x44000000)`

- `typedef struct ompi_communicator_t *MPI_Comm;`
  `#define MPI_COMM_WORLD ((MPI_Comm) &ompi_mpi_comm_world)`

- Need to recompile to use a different MPI library
  — May or may not be feasible

Gamblin et al. The Spack package manager: bringing order to HPC software chaos. SC'15.

## Wi4MPI: A general approach to ABI translation

- $T: f_\beta \to f_\rho$
  1. Translate input arguments from the β ABI to the ρ ABI
  2. Call $f_\rho$
  3. Translate output arguments and return value from the ρ ABI to the β ABI

- Automatically generate translation functions from two JSON files
  — The **functions**' signatures
  — The **mappers**' descriptor used to translate arguments

| | | |
|---|---|---|
| **Function** | *name* | Symbol's name |
| | *args* | Hash describing each argument |
| | *return* | Return mapping |
| **Func. Argument** | *name* | Argument's name |
| | *mapper* | Mapping descriptor to translate arg. |
| | *in* | Should argument be converted before $f$? |
| | *out* | Should argument be converted after $f$? |
| | *size* | Name of size argument for non-scalars |
| **Mapper** | *type* | Argument's type |
| | *b2r* | Func. to translate argument from β to ρ |
| | *r2b* | Func. to translate argument from ρ to β |
| | *alloc* | Should argument be allocated locally? |

## int MPI_Comm_rank(MPI_Comm comm, int* rank)

```
name: MPI_Comm_rank,

args: [
  {  name: comm,
   mapper: comm_converter,
       in: 1,
      out: 0
  },
  {  name: rank,
   mapper: int_ptr_mapper,
       in: 0,
      out: 1
  }
],

return: {
     name: rc
   mapper: error_converter,
       in: 0,
      out: 1
}...
```
<div align="right">Function JSON file</div>

```
comm_converter: {
  type: MPI_Comm,
  b2r: comm_conv_b2r,
  r2b: comm_conv_r2b,
  alloc: 0
}...
```
<div align="right">Mapper JSON file</div>

# Translating MPI dynamically

Build time

Run time

EXE | OpenMPI → Wi4MPI → Intel MPI / MPICH / MPC

β: OpenMPI ——→ ρ: Intel MPI
      Wi4MPI

MPI_File_open                    MPI_File_open        MPI_Allreduce

PMPI_File_open
Check: Not in wrapper

PMPI_Allreduce
Check: In wrapper

B_MPI_File_open
Translate arguments β → ρ

B_MPI_Allreduce
Translate arguments β → ρ

# Wi4MPI is easy to use

Build time

Run time

EXE | OpenMPI → Wi4MPI → Intel MPI / MPICH / MPC

(1)
```
# Specify Intel MPI as the target
export INTELMPI_ROOT=<path/to/intelmpi>
export WI4MPI_RUN_MPI_C_LIB=<path/to/intelmpi/libmpi.so>
export WI4MPI_RUN_MPIIO_C_LIB=<path/to/intelmpi/libmpi.so>
```

```
# Specify the Wi4MPI library to use
export LD_PRELOAD=<path/to/wi4mpi>/libwi4mpi_openmpi_intelmpi.so
```

(1)
```
# Load the appropriate Wi4MPI module
module load wi4mpi/3.4.0_openmpi_to_intelmpi
```

(2)
```
# Launch binary built with OpenMPI
srun —N20 —n80 <app_binary>
```

## Use cases demonstrating Wi4MPI

- Choosing the best MPI for my code
  - Arm architecture
  - x86 architectures
  - x86 + GPUs hybrid architectures

- Enabling applications on new hardware
  - New GPU systems
  - Bleeding-edge and early-access systems

- Running fast and portable containers
  - Use any MPI with host-tuned performance



On-the-Fly, Robust Translation of MPI Libraries
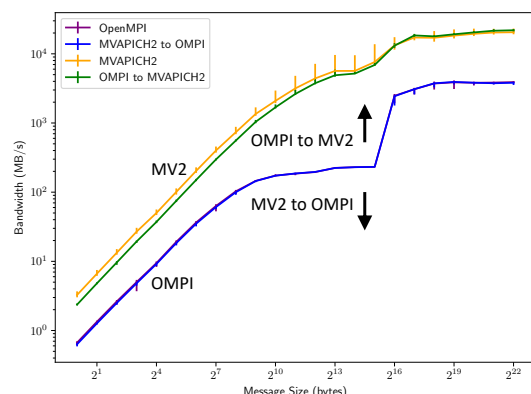
IEEE CLUSTER 2021
Portland, OR, USA • 7–10 September

---

## Wi4MPI allows me to choose the best MPI at runtime

Performance of translated and native MPI operations using OMB on x86 nodes



Better ↓  Startup time



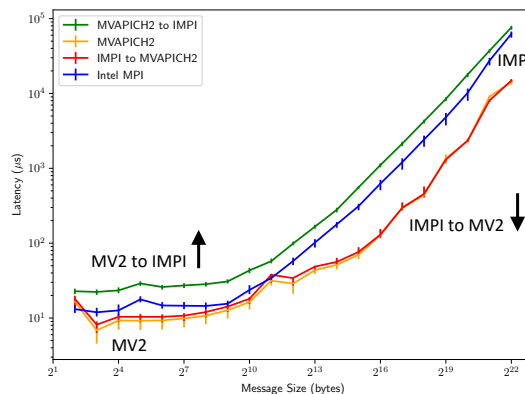Better ↑  Messaging rate across 2 nodes

## The performance gains of a better MPI often outweigh small-message overhead

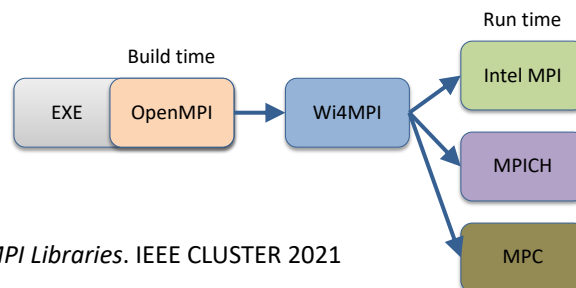Performance of translated and native MPI operations using OMB on x86 nodes



Bidirectional bandwidth across 2 nodes

Allreduce time of 144 tasks across 4 nodes

---

## Wi4MPI solves ABI incompatibility in the context of MPI

- Wi4MPI provides an effective framework for ABI translation

- Efficiently translating MPI has many benefits
  — Choose the best MPI for your code
  — Portable and fast Containers

- Wi4MPI has limitations
  — High overhead in some cases
  — Partial support for MPI_THREAD_MULTIPLE

- Wi4MPI is open source
  — León et al., *On-the-fly, Robust Translation of MPI Libraries*. IEEE CLUSTER 2021
  — https://github.com/cea-hpc/wi4mpi

## MPI Forum: Standardizing the ABI layer

- MPI Forum *likely* to define a C ABI in the future
  — Two ABIs cover over 90% of HPC platforms

- Plan is to have a single-feature ABI-only release for MPI 4.2
  — SC 2024

- MPICH has a prototype and would support the standard ABI
  — https://github.com/jeffhammond/mukautuva

- More information at the MPI ABI Working Group
  — https://github.com/mpiwg-abi

*Credit: Jeff Hammond on behalf of the MPI ABI Working Group*

## Hands-on: AWS Parallel Cluster

- `ssh` client needed

- 35 accounts
  — `user1`, …, `user35`

- `ssh user<k>@`

- IP address:

- Username: `user<k>`

- Password:

Wrapper interface for MPI

https://github.com/LLNL/productivity-frameworks/tree/main/tutorials/isc23