

# Using PMIx to Help Replace MPI\_Init

14<sup>th</sup> November 2018

Dan Holmes, Howard Pritchard, Nathan Hjelm

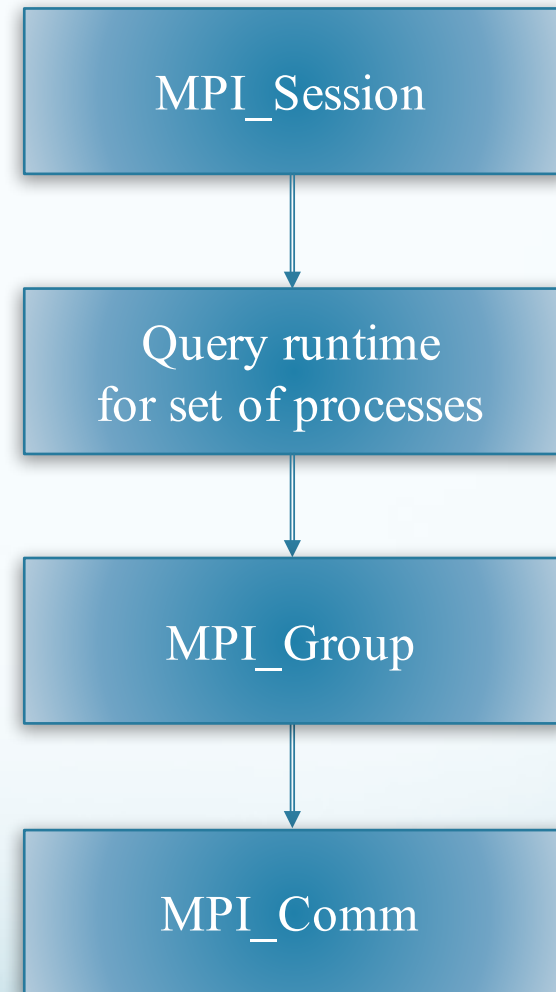
LA-UR-18-30830

# Problems with MPI\_Init

- All MPI processes must initialize MPI exactly once
- MPI cannot be initialized within an MPI process from different application components without coordination
- MPI cannot be re-initialized after MPI is finalized
- Error handling for MPI initialization cannot be specified

# Sessions – a new way to start MPI

- General scheme:
  - Query the underlying run-time system
    - Get a “set” of processes *Could be PMIx*
  - Determine the processes you want
    - Create an MPI\_Group
  - Create a communicator with just those processes
    - Create an MPI\_Comm



# MPI Sessions proposed API

- Create (or destroy) a session:
  - MPI\_SESSION\_INIT (and MPI\_SESSION\_FINALIZE)
- Get names of sets of processes:
  - MPI\_SESSION\_GET\_NUM\_PSETS,  
MPI\_SESSION\_GET\_NTH\_PSET
- Create an MPI\_GROUP from a process set name:
  - MPI\_GROUP\_CREATE\_FROM\_SESSION
- Create an MPI\_COMM from an MPI\_GROUP:
  - MPI\_COMM\_CREATE\_FROM\_GROUP

PMIx groups  
helps here

# MPI\_COMM\_CREATE\_FROM\_GROUP

```
MPI_Create_comm_from_group(IN MPI_Group group,
                            IN cc
                            *uri,
                            IN MF
                            info,
                            IN
                            MPI_Erhandler hndl,
                            OUT M
                            *comm);
```

The 'uri' is supplied by the application.

Implementation challenge: 'group' is a local object.

Need some way to synchronize with other "joiners" to the communicator. The 'uri' is different than a process set name.

# Using PMIx Groups

- PMIx Groups - a collection of processes desiring a unified identifier for purposes such as passing events or participating in PMIx fence operations
  - Invite/join/leave semantics
- Sessions prototype implementation currently uses `PMIX_Group_construct/PMIX_Group_destruct`
- Can be used to generate a “unique” 64-bit identifier for the group. Used by the sessions prototype to generate a communicator ID.
- Useful options for future work
  - Timeout for processes joining the group
  - Asynchronous notification when a process leaves the group

# Using PMIx\_Group\_Construct

```
PMIx_Group_Construct(const char id[],  
                    const pmix_proc_t procs[],  
                    const pmix_info_t info[],  
                    MPI_Comm *comm, MPI_Aint  
                    Size_t ninfo);
```

- ‘id’ maps to/from the ‘uri’ in MPI\_Comm\_create\_from\_group (plus additional Open MPI internal info)
- ‘procs’ array comes from information previously supplied by PMIx
  - “mpi://world” and “mpi://self” already available
  - `mpiexec -np 2 --pset user://ocean ocean.x : \`  
`-np 2 --pset user://atmosphere atmosphere.x`

*Work in progress*

# MPI Sessions Prototype Status

- All “MPI\_\*\_from\_group” functions have been implemented
  - Only pml/ob1 supported at this time
- Working on MPI\_Group creation (from PSET) now
  - Will start with mpi://world and mpi://self
  - User-defined process sets need **additional support from PMIx**
- Up next: break apart MPI initialization
  - Goal is to reduce startup time and memory footprint



# Summary

- PMIx Groups provides an OOB mechanism for MPI processes to bootstrap the formation of a communication context (MPI Communicator) from a group of MPI processes
- Functionality for future work
  - Handling (unexpected) process exit
  - User-defined process sets
  - Group expansion

# Funding Acknowledgments



**EPIGRAM HS**