

CHAPTER 12

Process Sets and Groups

1 PMIx supports two slightly related, but functionally different concepts known as *process sets* and
2 *process groups*. This chapter describes the two definitions and how they are utilized, along with
3 their corresponding APIs.

4 12.1 Process Sets

5 A PMIx *Process Set* is a user-provided label associated with a given set of application processes.
6 Definition of a PMIx process set typically occurs at time of application execution - e.g., on a
7 PRRTE command line:

```
8 $ prun -n 4 --pset ocean myoceanapp : -n 3 --pset ice myiceapp
```

9 In this example, the processes in the first application will be labeled with a **PMIX_PSET_NAME**
10 attribute of *ocean* while those in the second application will be labeled with an *ice* value. During
11 the execution, application processes could lookup the process set attribute for any other process
12 using **PMIx_Get** . Alternatively, other executing applications could utilize the
13 **PMIx_Query_info_nb** API to obtain the number of declared process sets in the system, a list
14 of their names, and other information about them. In other words, the *process set* identifier provides
15 a label by which an application can derive information about a process and its application - it does
16 *not*, however, confer any operational function.

17 Thus, *process sets* differ from *process groups* in several key ways:

- 18 • *Process sets* have no implied relationship between their members - i.e., a process in a process set
19 has no concept of a “pset rank” as it would in a *process group*
- 20 • Processes can only have one *process set* identifier, but can simultaneously belong to multiple
21 *process groups*
- 22 • *Process set* identifiers are considered job-level information set at launch. No PMIx API is
23 provided by which a user can change the *process set* value of a process on-the-fly. In contrast,
24 PMIx *process groups* can only be defined dynamically by the application.

- Process *groups* can be used in calls to PMIx operations. Members of process *groups* that are involved in an operation are translated by their PMIx server into their *native* identifier prior to the operation being passed to the host environment. For example, an application can define a process group to consist of ranks 0 and 1 from the host-assigned namespace of 210456, identified by the group id of *foo*. If the application subsequently calls the **PMIx_Fence** API with a process identifier of {foo, PMIX_RANK_WILDCARD}, the PMIx server will replace that identifier with an array consisting of {210456, 0} and {210456, 1} - the host-assigned identifiers of the participating processes - prior to passing the request up to the host environment

The two concepts do, however, overlap in one specific area. Process *groups* are included in the process *set* information returned by calls to **PMIx_Query_info_nb**. Thus, a *process group* can effectively be considered an extended version of a *process set* that adds dynamic definition and operational context to the *process set* concept.

Advice to PMIx library implementers

PMIx implementations are required to include all active *group* identifiers in the returned list of process *set* names provided in response to the appropriate **PMIx_Query_info_nb** call.

12.2 Process Groups

PMIx *Groups* are defined as a collection of processes desiring a common, unique identifier for purposes such as passing events or participating in PMIx fence operations. As with processes that assemble via **PMIx_Connect**, each member of the group is provided with both the job-level information of any other namespace represented in the group, and the contact information for all group members. However, *groups* differ from **PMIx_Connect** assemblages in the following key areas:

- Relation to the host environment
 - Calls to **PMIx_Connect** are relayed to the host environment. This means that the host RM should treat the failure of any process in the specified assemblage as a reportable event and take appropriate action. However, the environment is not required to define a new identifier for the connected assemblage or any of its member processes, nor does it define a new rank for each process within that assemblage. In addition, the PMIx server does not provide any tracking support for the assemblage. Thus, the caller is responsible for addressing members of the connected assemblage using their RM-provided identifiers.

- 1 – Calls to PMIx Group APIs are first processed within the local PMIx server. When constructed,
2 the server creates a tracker that associates the specified processes with the user-provided group
3 identifier, and assigns a new *group rank* based on their relative position in the array of
4 processes provided in the call to **PMIx_Group_construct** . Members of the group can
5 subsequently utilize the group identifier in PMIx function calls to address the group’s
6 members, using either **PMIX_RANK_WILDCARD** to refer to all of them or the group-level
7 rank of specific members. The PMIx server will translate the specified processes into their
8 RM-assigned identifiers prior to passing the request up to its host. Thus, the host environment
9 has no visibility into the group’s existence or membership.

Advice to users

10 User-provided group identifiers must be distinct from anything provided by the RM so as to
11 avoid collisions between group identifiers and RM-assigned namespaces. This can usually be
12 accomplished through the use of an application-specific prefix – e.g., “myapp-foo”

• Construction procedure

- 14 – **PMIx_Connect** calls require that every process call the API before completing – i.e., it is
15 modeled upon the bulk synchronous traditional MPI connect/accept methodology. Thus, a
16 given application thread can only be involved in one connect/accept operation at a time, and is
17 blocked in that operation until all specified processes participate. In addition, there is no
18 provision for replacing processes in the assemblage due to failure to participate, nor a
19 mechanism by which a process might decline participation.
- 20 – PMIx Groups are designed to be more flexible in their construction procedure by relaxing
21 these constraints. While a standard blocking form of constructing groups is provided, the event
22 notification system is utilized to provide a designated *group leader* with the ability to replace
23 participants that fail to participate within a given timeout period. This provides a mechanism
24 by which the application can, if desired, replace members on-the-fly or allow the group to
25 proceed with partial membership. In such cases, the final group membership is returned to all
26 participants upon completion of the operation.

27 Additionally, PMIx supports dynamic definition of group membership based on an invite/join
28 model. A process can asynchronously initiate construction of a group of any processes via the
29 **PMIx_Group_invite** function call. Invitations are delivered via a PMIx event (using the
30 **PMIX_GROUP_INVITED** event) to the invited processes which can then either accept or
31 decline the invitation using the **PMIx_Group_join** API. The initiating process tracks
32 responses by registering for the events generated by the call to **PMIx_Group_join** ,
33 timeouts, or process terminations, optionally replacing processes that decline the invitation,
34 fail to respond in time, or terminate without responding. Upon completion of the operation,
35 the final list of participants is communicated to each member of the new group.

• Destruct procedure

- 1 – Processes that assemble via `PMIx_Connect` must all depart the assemblage together – i.e.,
2 no member can depart the assemblage while leaving the remaining members in it. Even the
3 non-blocking form of `PMIx_Disconnect` retains this requirement in that members remain
4 a part of the assemblage until all members have called `PMIx_Disconnect_nb`
- 5 – Members of a PMIx Group may depart the group at any time via the `PMIx_Group_leave`
6 API. Other members are notified of the departure via the `PMIX_GROUP_LEFT` event to
7 distinguish such events from those reporting process termination. This leaves the remaining
8 members free to continue group operations. The `PMIx_Group_destruct` operation offers
9 a collective method akin to `PMIx_Disconnect` for deconstructing the entire group.

10 Note that applications supporting dynamic group behaviors such as asynchronous departure
11 take responsibility for ensuring global consistency in the group definition prior to executing
12 group collective operations - i.e., it is the application's responsibility to either ensure that
13 knowledge of the current group membership is globally consistent across the participants, or to
14 register for appropriate events to deal with the lack of consistency during the operation.

15 In other words, members of PMIx Groups are *loosely coupled* as opposed to *tightly connected*
16 when constructed via `PMIx_Connect` . The relevant APIs are explained below.

▼ Advice to users ▼

17 The reliance on PMIx events in the PMIx Group concept dictates that processes utilizing these APIs
18 must register for the corresponding events. Failure to do so will likely lead to operational failures.
19 Users are recommended to utilize the `PMIX_TIMEOUT` directive (or retain an internal timer) on
20 calls to PMIx Group APIs (especially the blocking form of those functions) as processes that have
21 not registered for required events will never respond.

22 12.2.1 `PMIx_Group_construct`

23 Summary

24 Construct a PMIx process group

1

Format

PMIx v4.0

C

2

`pmix_status_t`

3

`PMIx_Group_construct(const char grp[],`

4

`const pmix_proc_t procs[], size_t nprocs,`

5

`const pmix_info_t directives[], size_t ndirs,`

6

`pmix_info_t **results, size_t *nresults)`

C

7

IN grp

8

NULL-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the group identifier (string)

9

IN procs

10

Array of `pmix_proc_t` structures containing the PMIx identifiers of the member processes (array of handles)

11

12

IN nprocs

13

Number of elements in the *procs* array (`size_t`)

14

IN directives

15

Array of `pmix_info_t` structures (array of handles)

16

IN ndirs

17

Number of elements in the *directives* array (`size_t`)

18

INOUT results

19

Pointer to a location where the array of `pmix_info_t` describing the results of the operation is to be returned (pointer to handle)

20

21

INOUT nresults

22

Pointer to a `size_t` location where the number of elements in *results* is to be returned (memory reference)

23

24

Returns one of the following:

25

• `PMIX_SUCCESS`, indicating that the request has been successfully completed

26

• `PMIX_ERR_NOT_SUPPORTED` The PMIx library and/or the host RM does not support this operation

27

28

• a PMIx error constant indicating either an error in the input or that the request failed to be completed

29

30

Required Attributes

31

The following attributes are *required* to be supported by all PMIx libraries that support this operation:

32

33

PMIX_GROUP_LEADER "pmix.grp.ldr" (bool)

34

This process is the leader of the group

35

PMIX_GROUP_OPTIONAL "pmix.grp.opt" (bool)

1 Participation is optional - do not return an error if any of the specified processes terminate
2 without having joined. The default is false

3 **PMIX_GROUP_LOCAL_ONLY** "pmix.grp.lcl" (bool)

4 Group operation only involves local processes. PMIx implementations are *required* to
5 automatically scan an array of group members for local vs remote processes - if only local
6 processes are detected, the implementation need not execute a global collective for the
7 operation unless a context ID has been requested from the host environment. This can result
8 in significant time savings. This attribute can be used to optimize the operation by indicating
9 whether or not only local processes are represented, thus allowing the implementation to
10 bypass the scan. The default is false

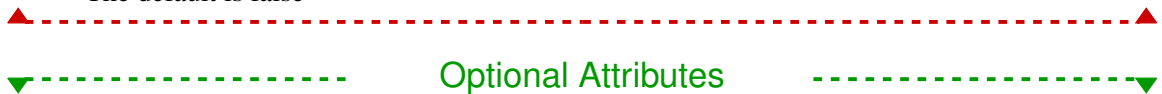
11 Host environments that support this operation are *required* to provide the following attributes:

12 **PMIX_GROUP_ASSIGN_CONTEXT_ID** "pmix.grp.actxid" (bool)

13 Notify remaining members when another member terminates without first leaving the group.
14 The default is false

15 **PMIX_GROUP_NOTIFY_TERMINATION** "pmix.grp.notterm" (bool)

16 Notify remaining members when another member terminates without first leaving the group.
17 The default is false

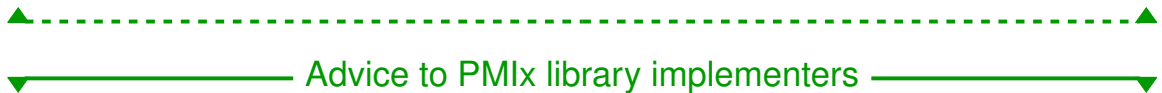


Optional Attributes

18 The following attributes are optional for host environments that support this operation:

19 **PMIX_TIMEOUT** "pmix.timeout" (int)

20 Time in seconds before the specified operation should time out (0 indicating infinite) in
21 error. The timeout parameter can help avoid "hangs" due to programming errors that prevent
22 the target process from ever exposing its data.



Advice to PMIx library implementers

23 We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host
24 environment due to race condition considerations between completion of the operation versus
25 internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT**
26 directly in the PMIx server library must take care to resolve the race condition and should avoid
27 passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not
28 created.

Description

Construct a new group composed of the specified processes and identified with the provided group identifier. The group identifier is a user-defined, **NULL**-terminated character array of length less than or equal to **PMIX_MAX_NSLEN**. Only characters accepted by standard string comparison functions (e.g., *strncmp*) are supported. Processes may engage in multiple simultaneous group construct operations so long as each is provided with a unique group ID. The *directives* array can be used to pass user-level directives regarding timeout constraints and other options available from the PMIx server.

If the **PMIX_GROUP_NOTIFY_TERMINATION** attribute is provided and has a value of **true**, then either the construct leader (if **PMIX_GROUP_LEADER** is provided) or all participants who register for the **PMIX_GROUP_MEMBER_FAILED** event will receive events whenever a process fails or terminates prior to calling **PMIx_Group_construct** – i.e. if a *group leader* is declared, *only* that process will receive the event. In the absence of a declared leader, *all* specified group members will receive the event.

The event will contain the identifier of the process that failed to join plus any other information that the host RM provided. This provides an opportunity for the leader or the collective members to react to the event – e.g., to decide to proceed with a smaller group or to abort the operation. The decision is communicated to the PMIx library in the results array at the end of the event handler. This allows PMIx to properly adjust accounting for procedure completion. When construct is complete, the participating PMIx servers will be alerted to any change in participants and each group member will receive an updated group membership (marked with the **PMIX_GROUP_MEMBERSHIP** attribute) as part of the *results* array returned by this API.

Failure of the declared leader at any time will cause a **PMIX_GROUP_LEADER_FAILED** event to be delivered to all participants so they can optionally declare a new leader. A new leader is identified by providing the **PMIX_GROUP_LEADER** attribute in the results array in the return of the event handler. Only one process is allowed to return that attribute, thereby declaring itself as the new leader. Results of the leader selection will be communicated to all participants via a **PMIX_GROUP_LEADER_SELECTED** event identifying the new leader. If no leader was selected, then the **pmix_info_t** provided to that event handler will include that information so the participants can take appropriate action.

Any participant that returns **PMIX_GROUP_CONSTRUCT_ABORT** from either the **PMIX_GROUP_MEMBER_FAILED** or the **PMIX_GROUP_LEADER_FAILED** event handler will cause the construct process to abort, returning from the call with a **PMIX_GROUP_CONSTRUCT_ABORT** status.

If the **PMIX_GROUP_NOTIFY_TERMINATION** attribute is not provided or has a value of **false**, then the **PMIx_Group_construct** operation will simply return an error whenever a proposed group member fails or terminates prior to calling **PMIx_Group_construct**.

Providing the **PMIX_GROUP_OPTIONAL** attribute with a value of **true** directs the PMIx library to consider participation by any specified group member as non-required - thus, the operation will return **PMIX_SUCCESS** if all members participate, or **PMIX_ERR_PARTIAL_SUCCESS** if

1 some members fail to participate. The *results* array will contain the final group membership in the
2 latter case. Note that this use-case can cause the operation to hang if the **PMIX_TIMEOUT**
3 attribute is not specified and one or more group members fail to call **PMIX_Group_construct**
4 while continuing to execute. Also, note that no leader or member failed events will be generated
5 during the operation.

6 Processes in a group under construction are not allowed to leave the group until group construction
7 is complete. Upon completion of the construct procedure, each group member will have access to
8 the job-level information of all namespaces represented in the group plus any information posted
9 via **PMIX_Put** (subject to the usual scoping directives) for every group member.

▼ **Advice to PMIx library implementers** ▼

10 At the conclusion of the construct operation, the PMIx library is *required* to ensure that job-related
11 information from each participating namespace plus any information posted by group members via
12 **PMIX_Put** (subject to scoping directives) is available to each member via calls to **PMIX_Get** .



▼ **Advice to PMIx server hosts** ▼

13 The collective nature of this API generally results in use of a fence-like operation by the backend
14 host environment. Host environments that utilize the array of process participants as a *signature* for
15 such operations may experience potential conflicts should both a **PMIX_Group_construct**
16 and a **PMIX_Fence** operation involving the same participants be simultaneously executed. As
17 PMIx allows for such use-cases, it is therefore the responsibility of the host environment to resolve
18 any potential conflicts.



19 **12.2.2 PMIX_Group_construct_nb**

20 **Summary**

21 Non-blocking form of **PMIX_Group_construct**

1

Format

PMIx v4.0

C

2

pmix_status_t

3

PMIx_Group_construct_nb(const char grp[],

4

const pmix_proc_t procs[], size_t nprocs,

5

const pmix_info_t directives[], size_t ndirs,

6

pmix_info_cbfunc_t cbfunc, void *cbdata)

C

7

IN grp

8

NULL-terminated character array of maximum size **PMIX_MAX_NSLEN** containing the group identifier (string)

9

IN procs

10

Array of **pmix_proc_t** structures containing the PMIx identifiers of the member processes (array of handles)

11

12

IN nprocs

13

Number of elements in the *procs* array (**size_t**)

14

IN directives

15

Array of **pmix_info_t** structures (array of handles)

16

IN ndirs

17

Number of elements in the *directives* array (**size_t**)

18

IN cbfunc

19

Callback function **pmix_info_cbfunc_t** (function reference)

20

IN cbdata

21

Data to be passed to the callback function (memory reference)

22

Returns one of the following:

23

24

• **PMIX_SUCCESS** indicating that the request has been accepted for processing and the provided callback function will be executed upon completion of the operation. Note that the library *must not* invoke the callback function prior to returning from the API.

25

26

27

• **PMIX_OPERATION_SUCCEEDED**, indicating that the request was immediately processed and returned *success* - the *cbfunc* will *not* be called

28

29

• **PMIX_ERR_NOT_SUPPORTED** The PMIx library does not support this operation - the *cbfunc* will *not* be called

30

31

• a non-zero PMIx error constant indicating a reason for the request to have been rejected - the *cbfunc* will *not* be called

32

33

If executed, the status returned in the provided callback function will be one of the following constants:

34

35

• **PMIX_SUCCESS** The operation succeeded and all specified members participated.

- 1 • **PMIX_ERR_PARTIAL_SUCCESS** The operation succeeded but not all specified members
2 participated - the final group membership is included in the callback function
- 3 • **PMIX_ERR_NOT_SUPPORTED** While the PMIx server supports this operation, the host RM
4 does not.
- 5 • a non-zero PMIx error constant indicating a reason for the request's failure

Required Attributes

6 PMIx libraries that choose not to support this operation *must* return
7 **PMIX_ERR_NOT_SUPPORTED** when the function is called.

8 The following attributes are *required* to be supported by all PMIx libraries that support this
9 operation:

10 **PMIX_GROUP_LEADER** "pmix.grp.ldr" (bool)
11 This process is the leader of the group

12 **PMIX_GROUP_OPTIONAL** "pmix.grp.opt" (bool)
13 Participation is optional - do not return an error if any of the specified processes terminate
14 without having joined. The default is false

15 **PMIX_GROUP_LOCAL_ONLY** "pmix.grp.lcl" (bool)
16 Group operation only involves local processes. PMIx implementations are *required* to
17 automatically scan an array of group members for local vs remote processes - if only local
18 processes are detected, the implementation need not execute a global collective for the
19 operation unless a context ID has been requested from the host environment. This can result
20 in significant time savings. This attribute can be used to optimize the operation by indicating
21 whether or not only local processes are represented, thus allowing the implementation to
22 bypass the scan. The default is false

23 Host environments that support this operation are *required* to provide the following attributes:

24 **PMIX_GROUP_ASSIGN_CONTEXT_ID** "pmix.grp.actxid" (bool)
25 Notify remaining members when another member terminates without first leaving the group.
26 The default is false

27 **PMIX_GROUP_NOTIFY_TERMINATION** "pmix.grp.notterm" (bool)
28 Notify remaining members when another member terminates without first leaving the group.
29 The default is false

Optional Attributes

The following attributes are optional for host environments that support this operation:

PMIX_TIMEOUT "pmix.timeout" (int)

Time in seconds before the specified operation should time out (*0* indicating infinite) in error. The timeout parameter can help avoid “hangs” due to programming errors that prevent the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT** directly in the PMIx server library must take care to resolve the race condition and should avoid passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not created.

Description

Non-blocking version of the **PMIx_Group_construct** operation. The callback function will be called once all group members have called either **PMIx_Group_construct** or **PMIx_Group_construct_nb**.

12.2.3 PMIx_Group_destruct

Summary

Destruct a PMIx process group

1

Format

PMIx v4.0

C

2

`pmix_status_t`

3

`PMIx_Group_destruct(const char grp[],`

4

`const pmix_info_t directives[], size_t ndirs)`

C

5

IN `grp`

6

NULL-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the identifier of the group to be destructed (string)

7

8

IN `directives`

9

Array of `pmix_info_t` structures (array of handles)

10

IN `ndirs`

11

Number of elements in the *directives* array (**size_t**)

12

Returns one of the following:

13

• `PMIX_SUCCESS`, indicating that the request has been successfully completed

14

• `PMIX_ERR_NOT_SUPPORTED` The PMIx library and/or the host RM does not support this operation

15

16

• a PMIx error constant indicating either an error in the input or that the request failed to be completed

17

Required Attributes

18

For implementations and host environments that support the operation, there are no identified required attributes for this API.

19

Optional Attributes

20

The following attributes are optional for host environments that support this operation:

21

`PMIX_TIMEOUT` "`pmix.timeout`" (**int**)

22

Time in seconds before the specified operation should time out (0 indicating infinite) in

23

error. The timeout parameter can help avoid "hangs" due to programming errors that prevent

24

the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the `PMIX_TIMEOUT` attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support `PMIX_TIMEOUT` directly in the PMIx server library must take care to resolve the race condition and should avoid passing `PMIX_TIMEOUT` to the host environment so that multiple competing timeouts are not created.

Description

Destruct a group identified by the provided group identifier. Processes may engage in multiple simultaneous group destruct operations so long as each involves a unique group ID. The *directives* array can be used to pass user-level directives regarding timeout constraints and other options available from the PMIx server.

The destruct API will return an error if any group process fails or terminates prior to calling `PMIx_Group_destruct` or its non-blocking version unless the `PMIX_GROUP_NOTIFY_TERMINATION` attribute was provided (with a value of `false`) at time of group construction. If notification was requested, then the `PMIX_GROUP_MEMBER_FAILED` event will be delivered for each process that fails to call destruct and the destruct tracker updated to account for the lack of participation. The `PMIx_Group_destruct` operation will subsequently return `PMIX_SUCCESS` when the remaining processes have all called destruct – i.e., the event will serve in place of return of an error.

Advice to PMIx server hosts

The collective nature of this API generally results in use of a fence-like operation by the backend host environment. Host environments that utilize the array of process participants as a *signature* for such operations may experience potential conflicts should both a `PMIx_Group_destruct` and a `PMIx_Fence` operation involving the same participants be simultaneously executed. As PMIx allows for such use-cases, it is therefore the responsibility of the host environment to resolve any potential conflicts.

12.2.4 `PMIx_Group_destruct_nb`

Summary

Non-blocking form of `PMIx_Group_destruct`

1

Format

PMIx v4.0

C

2

`pmix_status_t`

3

`PMIx_Group_destruct_nb(const char grp[],`

4

`const pmix_info_t directives[], size_t ndirs,`

5

`pmix_op_cbfunc_t cbfunc, void *cbdata)`

C

6

IN `grp`

7

NULL-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the identifier of the group to be destructed (string)

8

9

IN `directives`

Array of `pmix_info_t` structures (array of handles)

10

IN `ndirs`

Number of elements in the *directives* array (**size_t**)

11

IN `cbfunc`

Callback function `pmix_op_cbfunc_t` (function reference)

12

IN `cbdata`

Data to be passed to the callback function (memory reference)

13

14

15

16

17

Returns one of the following:

18

• **PMIX_SUCCESS**, indicating that the request is being processed - result will be returned in the provided *cbfunc*. Note that the library *must not* invoke the callback function prior to returning from the API.

19

20

21

• **PMIX_OPERATION_SUCCEEDED**, indicating that the request was immediately processed and returned *success* - the *cbfunc* will *not* be called

22

23

• **PMIX_ERR_NOT_SUPPORTED** The PMIx library does not support this operation - the *cbfunc* will *not* be called

24

25

• a PMIx error constant indicating either an error in the input or that the request was immediately processed and failed - the *cbfunc* will *not* be called

26

27

If executed, the status returned in the provided callback function will be one of the following constants:

28

29

• **PMIX_SUCCESS** The operation was successfully completed

30

• **PMIX_ERR_NOT_SUPPORTED** While the PMIx server supports this operation, the host RM does not.

31

32

• a non-zero PMIx error constant indicating a reason for the request's failure

Required Attributes

PMIx libraries that choose not to support this operation *must* return `PMIX_ERR_NOT_SUPPORTED` when the function is called. For implementations and host environments that support the operation, there are no identified required attributes for this API.

Optional Attributes

The following attributes are optional for host environments that support this operation:

`PMIX_TIMEOUT` "`pmix.timeout`" (`int`)

Time in seconds before the specified operation should time out (`0` indicating infinite) in error. The timeout parameter can help avoid “hangs” due to programming errors that prevent the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the `PMIX_TIMEOUT` attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support `PMIX_TIMEOUT` directly in the PMIx server library must take care to resolve the race condition and should avoid passing `PMIX_TIMEOUT` to the host environment so that multiple competing timeouts are not created.

Description

Non-blocking version of the `PMIx_Group_destruct` operation. The callback function will be called once all members of the group have executed either `PMIx_Group_destruct` or `PMIx_Group_destruct_nb`.

12.2.5 `PMIx_Group_invite`

Summary

Asynchronously construct a PMIx process group

1 **Format**

PMIx v4.0

C

```
2 pmix_status_t
3 PMIx_Group_invite(const char grp[],
4                   const pmix_proc_t procs[], size_t nprocs,
5                   const pmix_info_t directives[], size_t ndirs,
6                   pmix_info_t **results, size_t *nresult)
```

C

- 7 **IN grp**
8 NULL-terminated character array of maximum size **PMIX_MAX_NSLEN** containing the
9 group identifier (string)
- 10 **IN procs**
11 Array of **pmix_proc_t** structures containing the PMIx identifiers of the processes to be
12 invited (array of handles)
- 13 **IN nprocs**
14 Number of elements in the *procs* array (**size_t**)
- 15 **IN directives**
16 Array of **pmix_info_t** structures (array of handles)
- 17 **IN ndirs**
18 Number of elements in the *directives* array (**size_t**)
- 19 **INOUT results**
20 Pointer to a location where the array of **pmix_info_t** describing the results of the
21 operation is to be returned (pointer to handle)
- 22 **INOUT nresults**
23 Pointer to a **size_t** location where the number of elements in *results* is to be returned
24 (memory reference)

25 Returns one of the following:

- 26 • **PMIX_SUCCESS**, indicating that the request has been successfully completed
- 27 • **PMIX_ERR_NOT_SUPPORTED** The PMIx library and/or the host RM does not support this
28 operation
- 29 • a PMIx error constant indicating either an error in the input or that the request failed to be
30 completed

Required Attributes

31 The following attributes are *required* to be supported by all PMIx libraries that support this
32 operation:

- 33 **PMIX_GROUP_OPTIONAL** "pmix.grp.opt" (**bool**)
34 Participation is optional - do not return an error if any of the specified processes terminate
35 without having joined. The default is false

1 Host environments that support this operation are *required* to provide the following attributes:

2 **PMIX_GROUP_ASSIGN_CONTEXT_ID** "pmix.grp.actxid" (bool)
3 Notify remaining members when another member terminates without first leaving the group.
4 The default is false

5 **PMIX_GROUP_NOTIFY_TERMINATION** "pmix.grp.notterm" (bool)
6 Notify remaining members when another member terminates without first leaving the group.
7 The default is false



▼----- Optional Attributes -----▼

8 The following attributes are optional for host environments that support this operation:

9 **PMIX_TIMEOUT** "pmix.timeout" (int)
10 Time in seconds before the specified operation should time out (0 indicating infinite) in
11 error. The timeout parameter can help avoid “hangs” due to programming errors that prevent
12 the target process from ever exposing its data.



▼----- Advice to PMIx library implementers -----▼

13 We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host
14 environment due to race condition considerations between completion of the operation versus
15 internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT**
16 directly in the PMIx server library must take care to resolve the race condition and should avoid
17 passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not
18 created.



Description

Explicitly invite the specified processes to join a group. The process making the `PMIx_Group_invite` call is automatically declared to be the *group leader*. Each invited process will be notified of the invitation via the `PMIX_GROUP_INVITED` event - the processes being invited must therefore register for the `PMIX_GROUP_INVITED` event in order to be notified of the invitation. Note that the PMIx event notification system caches events - thus, no ordering of invite versus event registration is required.

The invitation event will include the identity of the inviting process plus the name of the group. When ready to respond, each invited process provides a response using either the blocking or non-blocking form of `PMIx_Group_join`. This will notify the inviting process that the invitation was either accepted (via the `PMIX_GROUP_INVITE_ACCEPTED` event) or declined (via the `PMIX_GROUP_INVITE_DECLINED` event). The `PMIX_GROUP_INVITE_ACCEPTED` event is captured by the PMIx client library of the inviting process – i.e., the application itself does not need to register for this event. The library will track the number of accepting processes and alert the inviting process (by returning from the blocking form of `PMIx_Group_invite` or calling the callback function of the non-blocking form) when group construction completes.

The inviting process should, however, register for the `PMIX_GROUP_INVITE_DECLINED` if the application allows invited processes to decline the invitation. This provides an opportunity for the application to either invite a replacement, declare “abort”, or choose to remove the declining process from the final group. The inviting process should also register to receive `PMIX_GROUP_INVITE_FAILED` events whenever a process fails or terminates prior to responding to the invitation. Actions taken by the inviting process in response to these events must be communicated at the end of the event handler by returning the corresponding result so that the PMIx library can adjust accordingly.

Upon completion of the operation, all members of the new group will receive access to the job-level information of each other’s namespaces plus any information posted via `PMIx_Put` by the other members.

The inviting process is automatically considered the leader of the asynchronous group construction procedure and will receive all failure or termination events for invited members prior to completion. The inviting process is required to provide a `PMIX_GROUP_CONSTRUCT_COMPLETE` event once the group has been fully assembled – this event is used by the PMIx library as a trigger to release participants from their call to `PMIx_Group_join` and provides information (e.g., the final group membership) to be returned in the *results* array.

Advice to users

Applications are not allowed to use the group in any operations until group construction is complete. This is required in order to ensure consistent knowledge of group membership across all participants.

1 Failure of the inviting process at any time will cause a `PMIX_GROUP_LEADER_FAILED` event to
2 be delivered to all participants so they can optionally declare a new leader. A new leader is
3 identified by providing the `PMIX_GROUP_LEADER` attribute in the results array in the return of
4 the event handler. Only one process is allowed to return that attribute, declaring itself as the new
5 leader. Results of the leader selection will be communicated to all participants via a
6 `PMIX_GROUP_LEADER_SELECTED` event identifying the new leader. If no leader was selected,
7 then the status code provided in the event handler will provide an error value so the participants can
8 take appropriate action.

9 12.2.6 `PMIx_Group_invite_nb`

10 Summary

11 Non-blocking form of `PMIx_Group_invite`

12 Format

PMIx v4.0

```
13 pmix_status_t  
14 PMIx_Group_invite_nb(const char grp[],  
15                     const pmix_proc_t procs[], size_t nprocs,  
16                     const pmix_info_t directives[], size_t ndirs,  
17                     pmix_info_cbfunc_t cbfunc, void *cbdata)
```

- 18 **IN** `grp`
19 NULL-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the
20 group identifier (string)
- 21 **IN** `procs`
22 Array of `pmix_proc_t` structures containing the PMIx identifiers of the processes to be
23 invited (array of handles)
- 24 **IN** `nprocs`
25 Number of elements in the `procs` array (`size_t`)
- 26 **IN** `directives`
27 Array of `pmix_info_t` structures (array of handles)
- 28 **IN** `ndirs`
29 Number of elements in the `directives` array (`size_t`)
- 30 **IN** `cbfunc`
31 Callback function `pmix_info_cbfunc_t` (function reference)
- 32 **IN** `cbdata`
33 Data to be passed to the callback function (memory reference)

- 1 Returns one of the following:
- 2 • **PMIX_SUCCESS**, indicating that the request is being processed - result will be returned in the
 - 3 provided *cbfunc*. Note that the library *must not* invoke the callback function prior to returning
 - 4 from the API.
 - 5 • **PMIX_OPERATION_SUCCEEDED**, indicating that the request was immediately processed and
 - 6 returned *success* - the *cbfunc* will *not* be called
 - 7 • **PMIX_ERR_NOT_SUPPORTED** The PMIx library does not support this operation - the *cbfunc*
 - 8 will *not* be called
 - 9 • a PMIx error constant indicating either an error in the input or that the request was immediately
 - 10 processed and failed - the *cbfunc* will *not* be called

11 If executed, the status returned in the provided callback function will be one of the following

12 constants:

- 13 • **PMIX_SUCCESS** The operation succeeded and all specified members participated.
- 14 • **PMIX_ERR_PARTIAL_SUCCESS** The operation succeeded but not all specified members
- 15 participated - the final group membership is included in the callback function
- 16 • **PMIX_ERR_NOT_SUPPORTED** While the PMIx server supports this operation, the host RM
- 17 does not.
- 18 • a non-zero PMIx error constant indicating a reason for the request's failure

▼----- Required Attributes -----▼

19 The following attributes are *required* to be supported by all PMIx libraries that support this

20 operation:

21 **PMIX_GROUP_OPTIONAL** "pmix.grp.opt" (bool)

22 Participation is optional - do not return an error if any of the specified processes terminate

23 without having joined. The default is false

24 Host environments that support this operation are *required* to provide the following attributes:

25 **PMIX_GROUP_ASSIGN_CONTEXT_ID** "pmix.grp.actxid" (bool)

26 Notify remaining members when another member terminates without first leaving the group.

27 The default is false

28 **PMIX_GROUP_NOTIFY_TERMINATION** "pmix.grp.notterm" (bool)

29 Notify remaining members when another member terminates without first leaving the group.

30 The default is false

▲-----

Optional Attributes

The following attributes are optional for host environments that support this operation:

PMIX_TIMEOUT "pmix.timeout" (int)

Time in seconds before the specified operation should time out (*0* indicating infinite) in error. The timeout parameter can help avoid “hangs” due to programming errors that prevent the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT** directly in the PMIx server library must take care to resolve the race condition and should avoid passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not created.

Description

Non-blocking version of the **PMIx_Group_invite** operation. The callback function will be called once all invited members of the group (or their substitutes) have executed either **PMIx_Group_join** or **PMIx_Group_join_nb**.

12.2.7 PMIx_Group_join

Summary

Accept an invitation to join a PMIx process group

1 **Format**

PMIx v4.0

C

```

2 pmix_status_t
3 PMIx_Group_join(const char grp[],
4                 const pmix_proc_t *leader,
5                 pmix_group_opt_t opt,
6                 const pmix_info_t directives[], size_t ndirs,
7                 pmix_info_t **results, size_t *nresult)

```

C

- 8 **IN grp**
9 NULL-terminated character array of maximum size **PMIX_MAX_NSLEN** containing the
10 group identifier (string)
- 11 **IN leader**
12 Process that generated the invitation (handle)
- 13 **IN opt**
14 Accept or decline flag (**pmix_group_opt_t**)
- 15 **IN directives**
16 Array of **pmix_info_t** structures (array of handles)
- 17 **IN ndirs**
18 Number of elements in the *directives* array (**size_t**)
- 19 **INOUT results**
20 Pointer to a location where the array of **pmix_info_t** describing the results of the
21 operation is to be returned (pointer to handle)
- 22 **INOUT nresults**
23 Pointer to a **size_t** location where the number of elements in *results* is to be returned
24 (memory reference)
- 25 Returns one of the following:
 - 26 • **PMIX_SUCCESS** , indicating that the request has been successfully completed
 - 27 • **PMIX_ERR_NOT_SUPPORTED** The PMIx library and/or the host RM does not support this
28 operation
 - 29 • a PMIx error constant indicating either an error in the input or that the request failed to be
30 completed

Required Attributes

31 There are no identified required attributes for implementers.

Optional Attributes

The following attributes are optional for host environments that support this operation:

PMIX_TIMEOUT "pmix.timeout" (int)

Time in seconds before the specified operation should time out (0 indicating infinite) in error. The timeout parameter can help avoid “hangs” due to programming errors that prevent the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT** directly in the PMIx server library must take care to resolve the race condition and should avoid passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not created.

Description

Respond to an invitation to join a group that is being asynchronously constructed. The process must have registered for the **PMIX_GROUP_INVITED** event in order to be notified of the invitation. When called, the event information will include the **pmix_proc_t** identifier of the process that generated the invitation along with the identifier of the group being constructed. When ready to respond, the process provides a response using either form of **PMIx_Group_join**.

Advice to users

Since the process is alerted to the invitation in a PMIx event handler, the process *must not* use the blocking form of this call unless it first “thread shifts” out of the handler and into its own thread context. Likewise, while it is safe to call the non-blocking form of the API from the event handler, the process *must not* block in the handler while waiting for the callback function to be called.

1 Calling this function causes the inviting process (aka the *group leader*) to be notified that the
2 process has either accepted or declined the request. The blocking form of the API will return once
3 the group has been completely constructed or the group's construction has failed (as described
4 below) – likewise, the callback function of the non-blocking form will be executed upon the same
5 conditions.

6 Failure of the leader during the call to `PMIx_Group_join` will cause a
7 `PMIX_GROUP_LEADER_FAILED` event to be delivered to all invited participants so they can
8 optionally declare a new leader. A new leader is identified by providing the
9 `PMIX_GROUP_LEADER` attribute in the results array in the return of the event handler. Only one
10 process is allowed to return that attribute, declaring itself as the new leader. Results of the leader
11 selection will be communicated to all participants via a `PMIX_GROUP_LEADER_SELECTED`
12 event identifying the new leader. If no leader was selected, then the status code provided in the
13 event handler will provide an error value so the participants can take appropriate action.

14 Any participant that returns `PMIX_GROUP_CONSTRUCT_ABORT` from the leader failed event
15 handler will cause all participants to receive an event notifying them of that status. Similarly, the
16 leader may elect to abort the procedure by either returning `PMIX_GROUP_CONSTRUCT_ABORT`
17 from the handler assigned to the `PMIX_GROUP_INVITE_ACCEPTED` or
18 `PMIX_GROUP_INVITE_DECLINED` codes, or by generating an event for the abort code. Abort
19 events will be sent to all invited participants.

20 12.2.8 `PMIx_Group_join_nb`

21 Summary

22 Non-blocking form of `PMIx_Group_join`

23 Format

PMIx v4.0

```
24 pmix_status_t  
25 PMIx_Group_join_nb(const char grp[],  
26                   const pmix_proc_t *leader,  
27                   pmix_group_opt_t opt,  
28                   const pmix_info_t directives[], size_t ndirs,  
29                   pmix_info_cbfunc_t cbfunc, void *cbdata)
```

30 IN `grp`

31 `NULL`-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the
32 group identifier (string)

1 **IN leader**
 2 Process that generated the invitation (handle)
 3 **IN opt**
 4 Accept or decline flag (`pmix_group_opt_t`)
 5 **IN directives**
 6 Array of `pmix_info_t` structures (array of handles)
 7 **IN ndirs**
 8 Number of elements in the *directives* array (`size_t`)
 9 **IN cbfunc**
 10 Callback function `pmix_info_cbfunc_t` (function reference)
 11 **IN cbdata**
 12 Data to be passed to the callback function (memory reference)

13 Returns one of the following:

- 14 • **PMIX_SUCCESS** , indicating that the request is being processed - result will be returned in the
 15 provided *cbfunc*. Note that the library *must not* invoke the callback function prior to returning
 16 from the API.
- 17 • **PMIX_OPERATION_SUCCEEDED** , indicating that the request was immediately processed and
 18 returned *success* - the *cbfunc* will *not* be called
- 19 • **PMIX_ERR_NOT_SUPPORTED** The PMIx library does not support this operation - the *cbfunc*
 20 will *not* be called
- 21 • a PMIx error constant indicating either an error in the input or that the request was immediately
 22 processed and failed - the *cbfunc* will *not* be called

23 If executed, the status returned in the provided callback function will be one of the following
 24 constants:

- 25 • **PMIX_SUCCESS** The operation succeeded and group membership is in the callback function
 26 parameters
- 27 • **PMIX_ERR_NOT_SUPPORTED** While the PMIx server supports this operation, the host RM
 28 does not.
- 29 • a non-zero PMIx error constant indicating a reason for the request's failure

▼----- Required Attributes -----▼

30 There are no identified required attributes for implementers.
 ▲-----

Optional Attributes

The following attributes are optional for host environments that support this operation:

PMIX_TIMEOUT "pmix.timeout" (int)

Time in seconds before the specified operation should time out (*0* indicating infinite) in error. The timeout parameter can help avoid “hangs” due to programming errors that prevent the target process from ever exposing its data.

Advice to PMIx library implementers

We recommend that implementation of the **PMIX_TIMEOUT** attribute be left to the host environment due to race condition considerations between completion of the operation versus internal timeout in the PMIx server library. Implementers that choose to support **PMIX_TIMEOUT** directly in the PMIx server library must take care to resolve the race condition and should avoid passing **PMIX_TIMEOUT** to the host environment so that multiple competing timeouts are not created.

Description

Non-blocking version of the **PMIx_Group_join** operation. The callback function will be called once all invited members of the group (or their substitutes) have executed either **PMIx_Group_join** or **PMIx_Group_join_nb**.

12.2.9 PMIx_Group_leave

Summary

Leave a PMIx process group

1

Format

PMIx v4.0

C

2

`pmix_status_t`

3

`PMIx_Group_leave(const char grp[],`

4

`const pmix_info_t directives[], size_t ndirs)`

C

5

IN `grp`

6

NULL-terminated character array of maximum size `PMIX_MAX_NSLEN` containing the group identifier (string)

7

IN `directives`

8

Array of `pmix_info_t` structures (array of handles)

9

IN `ndirs`

10

Number of elements in the *directives* array (**size_t**)

11

Returns one of the following:

12

- `PMIX_SUCCESS`, indicating that the request has been communicated to the local PMIx server
- `PMIX_ERR_NOT_SUPPORTED` The PMIx library and/or the host RM does not support this operation
- a PMIx error constant indicating either an error in the input or that the request is unsupported

13

14

15

16

Required Attributes

17

There are no identified required attributes for implementers.

18

Description

19

Leave a PMIx Group. Calls to `PMIx_Group_leave` (or its non-blocking form) will cause a `PMIX_GROUP_LEFT` event to be generated notifying all members of the group of the caller's departure. The function will return (or the non-blocking function will execute the specified callback function) once the event has been locally generated and is not indicative of remote receipt. All PMIx-based collectives such as `PMIx_Fence` in action across the group will automatically be adjusted if the collective was called with the `PMIX_GROUP_FT_COLLECTIVE` attribute (default is false) – otherwise, the standard error return behavior for that collective will be executed.

20

21

22

23

24

25

Advice to users

26

The `PMIx_Group_leave` API is intended solely for asynchronous departures of individual processes from a group as it is not a scalable operation – i.e., when a process determines it should no longer be a part of a defined group, but the remainder of the group retains a valid reason to continue in existence. Developers are advised to use `PMIx_Group_destruct` (or its non-blocking form) for all other scenarios as it represents a more scalable operation.

27

28

29

30

1 12.2.10 PMIx_Group_leave_nb

2 Summary

3 Non-blocking form of [PMIx_Group_leave](#)

4 Format

PMIx v4.0

C

```
5 pmix_status_t
6 PMIx_Group_leave_nb(const char grp[],
7                     const pmix_info_t directives[], size_t ndirs,
8                     pmix_op_cbfunc_t cbfunc, void *cbdata)
```

C

- 9 **IN** **grp**
10 **NULL**-terminated character array of maximum size [PMIX_MAX_NSLEN](#) containing the
11 group identifier (string)
- 12 **IN** **directives**
13 Array of [pmix_info_t](#) structures (array of handles)
- 14 **IN** **ndirs**
15 Number of elements in the *directives* array (**size_t**)
- 16 **IN** **cbfunc**
17 Callback function [pmix_op_cbfunc_t](#) (function reference)
- 18 **IN** **cbdata**
19 Data to be passed to the callback function (memory reference)

20 Returns one of the following:

- 21 • [PMIX_SUCCESS](#) , indicating that the request is being processed - result will be returned in the
22 provided *cbfunc*. Note that the library *must not* invoke the callback function prior to returning
23 from the API.
- 24 • [PMIX_OPERATION_SUCCEEDED](#) , indicating that the request was immediately processed and
25 returned *success* - the *cbfunc* will *not* be called
- 26 • [PMIX_ERR_NOT_SUPPORTED](#) The PMIx library does not support this operation - the *cbfunc*
27 will *not* be called
- 28 • a PMIx error constant indicating either an error in the input or that the request was immediately
29 processed and failed - the *cbfunc* will *not* be called

30 If executed, the status returned in the provided callback function will be one of the following
31 constants:

- 32 • [PMIX_SUCCESS](#) The operation succeeded - i.e., the [PMIX_GROUP_LEFT](#) event was
33 generated

1 • **PMIX_ERR_NOT_SUPPORTED** While the PMIx library supports this operation, the host RM
2 does not.

3 • a non-zero PMIx error constant indicating a reason for the request's failure

▼----- Required Attributes -----▼

4 There are no identified required attributes for implementers.

▲-----▲

5 **Description**

6 Non-blocking version of the **PMIx_Group_leave** operation. The callback function will be
7 called once the event has been locally generated and is not indicative of remote receipt.