#### CALIPER



# Enabling introspection of MPI libraries through the MPI\_T interface

Authors: Srinivasan Ramesh

### Introduction

MPI libraries today involve multiple components interacting in complex ways to affect performance. Together with the heterogeneous nature of current and future architectures, this means default MPI library settings may not always be optimal for performance and scalability — significant performance gains may be achieved by tailoring MPI library behaviour to suit application characteristics. To understand MPI library internals, there needs to be a way to introspect the MPI\_T libary at runtime. In order to modify MPI library behaviour dynamically at runtime, the library must expose a means to do so.

The MPI\_T interface, introduced in the MPI 3.0 standard provides external tools an opportunity to introspect and potentially modify MPI library behaviour at runtime by means of two semantics:

- Performance Variables (PVARs): Performance variables represent MPI internal information in the form of counters, state, watermarks, etc. The MPI specification details the various classes of PVARs supported, allowed datatypes and access semantics each of class.
- Control Variables (CVARs): Control variables are the means by which an external tool can modify MPI library behaviour and fine-tune application performance. They are essentially knobs that may represent the value of a particular setting inside the MPI library.

Caliper is an application introspection tool that relies on source code annotations to collect information and perform profiling related tasks. Caliper services are the basic building blocks that can be combined freely to realize advanced profiling / tracing capabilities. The MPI service utilizes the PMPI interface to profile MPI library calls. This document describes the design of the MPIT service that performs MPI library introspection through the MPI\_T interface. The motivation behind this document is to describe the rationale that went into the design of the service. Caliper is in active development — this document deliberately avoids description of code or filenames

used to implement the MPIT service. However, any design modifications to the service shall be reflected here.

We shall touch upon some Caliper concepts when required. For detailed description of Caliper design, kindly refer to Caliper documentation.

## Collecting PVARs through MPI\_T

#### 0.1 Creating a performance session

In order to use MPI\_T, a tool must first create a *performance session* and associate *handles* for the performance variables it wishes to read. Performance sessions allow the MPI library to distinguish between multiple tools / software modules that may be simultaneously querying the MPI\_T interface.

The MPIT Caliper service creates an MPI\_T performance session during the Caliper service registration phase.

#### 0.2 PVAR handle allocation

Before a tool can read the value of a PVAR, it must first allocate a *handle* for the PVAR. The MPI\_T interface specifies a function that allows a tool to know the number of PVARs exported by an MPI implementation at any given point in time. A couple of points need to be kept in mind when allocating PVAR handles:

- Number of PVARs can change: The number of PVARs exported by the library can change at any point in time. Typically, MPI libraries export additional PVARs after MPI\_Init. Caliper allocates handles for PVARs from multiple places: During the registration for the MPIT service, inside the wrapper for MPI\_Init, and from inside the wrapper for certain MPI calls (description follows).
- PVARs can be bound to MPI objects: The MPI\_pvar\_get\_info function returns the *bind* type for the PVAR. The idea here is that PVARs

can be associated with a specific object such as a communicator or message. As a result, there can be multiple handles allocated for a PVAR at any given index. These handles must be allocated appropriately depending on the bind type.

- MPI\_T\_BIND\_NO\_OBJECT: These PVARs are not bound to any MPI object Caliper allocates handles for such PVARs during Caliper registration and inside the wrapper for MPI\_Init.
- MPI\_T\_BIND\_MPI\_COMM: These PVARs are bound to MPI communicators, and is a special case. Handles for pvars bound to MPI\_COMM\_WORLD and MPI\_COMM\_SELF are allocated during Caliper registration phase. Additionally, handles are created each time MPI\_Comm\_create is invoked, by intercepting the call through the PMPI wrapper.
- MPI\_T\_BIND\_WIN: These PVARs are bound to MPI windows. As a result, handles for such pvars are allocated inside the Caliper PMPI wrapper for MPI\_Win\_create.
- MPI\_T\_BIND\_MPI\_ERR\_HANDLER: These PVARs are bound to MPI error handlers. Handles for such pvars are allocated inside the PMPI wrapper for MPI\_Errhandler\_create.
- MPI\_T\_BIND\_MPI\_FILE: These PVARs are bound to file objects.
  Handles are allocated inside the PMPI wrapper for MPI\_File\_open.
- MPI\_T\_BIND\_MPI\_GROUP: These PVARs are bound to MPI group objects. Handles are allocated inside the PMPI wrapper for MPI\_COMM\_GROUP.
- MPI\_T\_BIND\_MPI\_OP: These PVARs are bound to MPI reduction operators. Handles are allocated inside the PMPI wrapper for MPI\_Op\_create.
- MPI\_T\_BIND\_MPI\_INFO: These PVARs are bound to MPI Info objects. Handles are allocated inside the PMPI wrapper for MPI\_Info\_create.
- MPI\_T\_BIND\_MPI\_MESSAGE, MPI\_T\_BIND\_MPI\_REQUEST: Not supported inside Caliper currently. Open question how do we allocate handles for these?

#### 0.3 PVAR classes and notion of aggregability

Depending on what they represent, PVARs are categorized into counters, state variables, watermarks, etc., and are handled differently. For this pur-

pose, we define the notion of aggregatability as follows: Any PVAR on which it is *meaningful* to apply one or more of (SUM, MAX, MIN, AVG, COUNT) operators is defined as aggregatable.

Along with other information, a call to MPIT\_pvar\_get\_info returns the *CLASS* to which the PVAR belongs. The various classes, along with how Caliper handles them are:

- MPI\_T\_PVAR\_CLASS\_TIMER, MPI\_T\_PVAR\_CLASS\_AGGREGATE, MPI\_T\_PVAR\_CLASS\_COUNTERS: These are free-counting, monotonically increasing values. As such, they are not aggregatable, but by storing the "last" value for these counters and timers, the difference between the current and last value is a derived metric that is aggregatable by use of SUM, MAX, MIN, AVG operators. Storing this difference is more useful than just the raw counter values, as one would typically by interested in the change caused to any of these PVARs rather than the raw value itself.
- MPI\_T\_PVAR\_CLASS\_STATE: Represents MPI state at any instant in time. Non-aggregatable value.
- MPI\_T\_PVAR\_CLASS\_SIZE: Represents size of an MPI resource. Non-aggregatable value.
- MPI\_T\_PVAR\_CLASS\_LEVEL, MPI\_T\_PVAR\_CLASS\_PERCENTAGE: Represents the instantaneous level or percentage utilization of an MPI resource. It is meaningful to apply the AVG, MIN, MAX operators, and hence these classes are aggregatable.
- MPI\_T\_PVAR\_CLASS\_HIGHWATERMARK, MPI\_T\_PVAR\_CLASS\_LOWWATERMARK: As such, they are non-aggregatable. However, one can define aggregatable derived metrics out of these PVARs. Specifically, Caliper defines two derived metrics: A boolean that tells us if the watermark has gone up from the last time it was read, and a double value specifying the *change* in the value between successive reads. Both of these derived metrics are aggregatable quantities as one can apply the *COUNT* and / or *SUM* operator to them.
- MPI\_T\_PVAR\_CLASS\_GENERIC: Represents PVARs that do not fall into any of the above classes. These PVARs would need to handled on a case-by-case basis, and thus for now, we define these as non-aggregatable values.

#### 0.4 Creating Caliper attributes for PVARs

The basic data unit in Caliper is an attribute. An attribute is a key value pair that has certain properties. For each PVAR exposed by the MPI library, Caliper defines an attribute with the same name as the PVAR. Each PVAR attribute has the following properties:

- CALI\_ATTR\_AS\_VALUE We do not want "stacking" for PVAR values. They should be treated much the same way as PAPI counters.
- CALI\_ATTR\_SCOPE\_PROCESS PVARs are defined on a per-rank basis
- CALI\_ATTR\_SKIP\_EVENTS We do not want callbacks to be triggered everytime the attribute for a PVAR is updated
- Metadata (class.aggregatable) Boolean value specifying if the PVAR is aggregatable or not. Aggregatability is determined in accordance with the rules above.

Apart from creating a Caliper attribute for each PVAR exported, there are two additional attributes created for each watermark PVAR exported — one that represents the number of times the watermark changes, and another that represents the cumulative change in the watermark PVAR.

#### 0.5 Querying and storing PVARs

In the current design, all PVARs exported by the MPI library are queried when a snapshot is triggered. By integrating the *MPIT* service along with the *MPI* service, this would be useful in determining how various MPI function calls contribute to changes in PVAR values. Moreover, one can gather meaningful information by aggregating using MPI function names or annotated code regions as keys. Currently, we note about 10-15% overheads in collecting PVARs during every snapshot event. This may not be a very scalable option, as this overhead would increase with a rise in number of PVARs exported.

Depending on the class of the PVAR, we either store the raw value read from the interface in the snapshot, or a derived metric.

MPI\_T\_PVAR\_CLASS\_TIMER, MPI\_T\_PVAR\_CLASS\_AGGREGATE,
 MPI\_T\_PVAR\_CLASS\_COUNTERS: We store the difference between the current value and the "last value" for such PVARs in the snapshot.
 Storing and aggregating this derived value is more meaningful — it

helps us answer questions such as: How do different MPI functions contribute to this PVAR? Which MPI function is responsible for the highest value?

- MPI\_T\_PVAR\_CLASS\_STATE, MPI\_T\_PVAR\_CLASS\_SIZE: These PVARs are stored as is in the snapshot. Perhaps it may be more meaningful to view changes *over time*, such as in a trace.
- MPI\_T\_PVAR\_CLASS\_HIGHWATERMARK, MPI\_T\_PVAR\_CLASS\_LOWWATERMARK: Along with storing the raw value for watermark PVARs, we store the derived metrics that represent the number of times the watermark changed, along with how much the watermark changed in the snapshot. By aggregating across MPI functions for example, we can answer questions such as: Which function most frequently pushed up / down a watermark? Which function was responsible for the highest cumulative change in a given watermark?
- MPI\_T\_PVAR\_CLASS\_LEVEL, MPI\_T\_PVAR\_CLASS\_PERCENTAGE: We store these PVARs as is in the snapshot. It maybe meaningful to view the average, maximum or minimum value for these PVARs, aggregated across MPI functions.