

# Scalable MPI Completion

9<sup>TH</sup> JLESC Workshop

04/15/2019

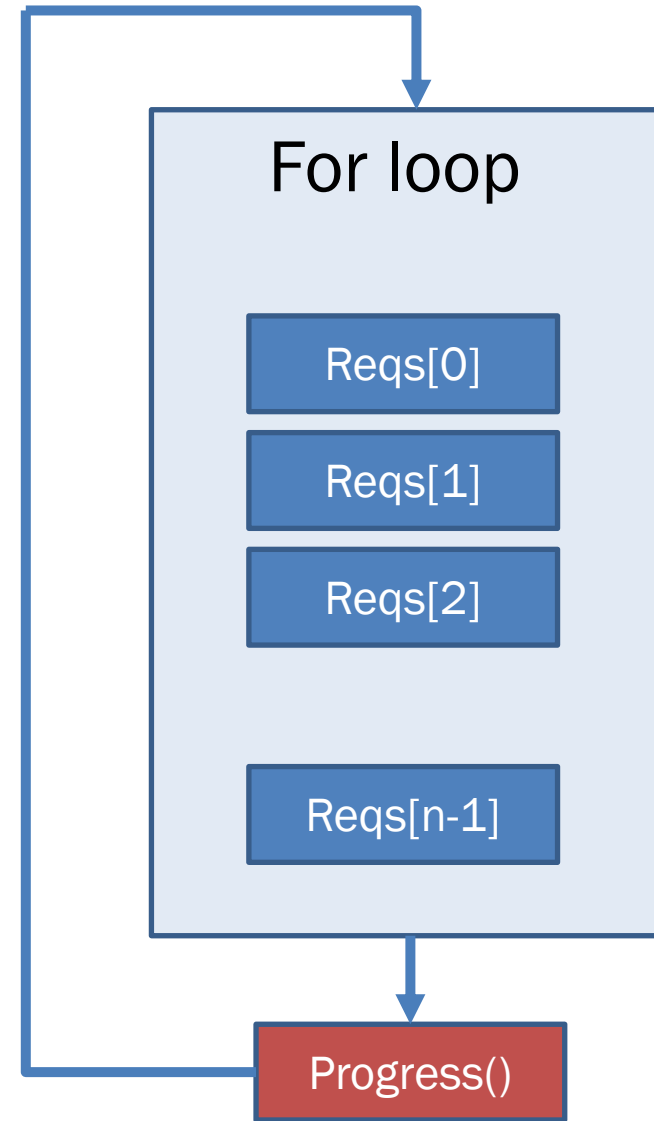
Thananon “Arm” Patinyasakdikul



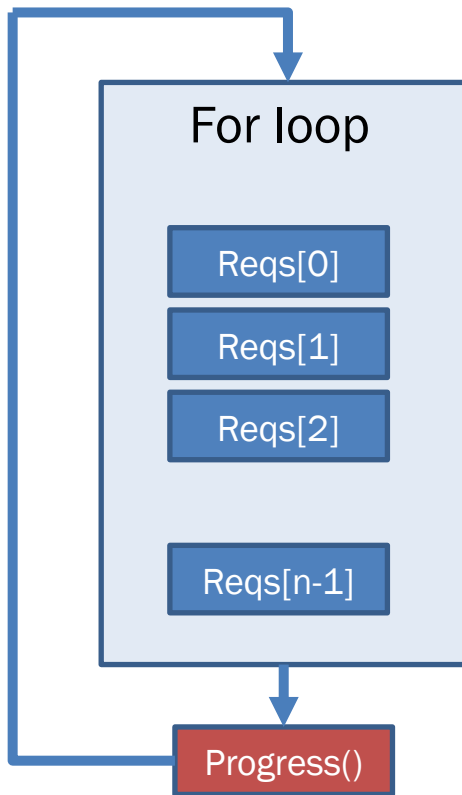
# MPI Completion - Design

MPI\_Waitall, MPI\_Waitsome, MPI\_Waitany  
MPI\_Testall, MPI\_Testsome, MPI\_Testany

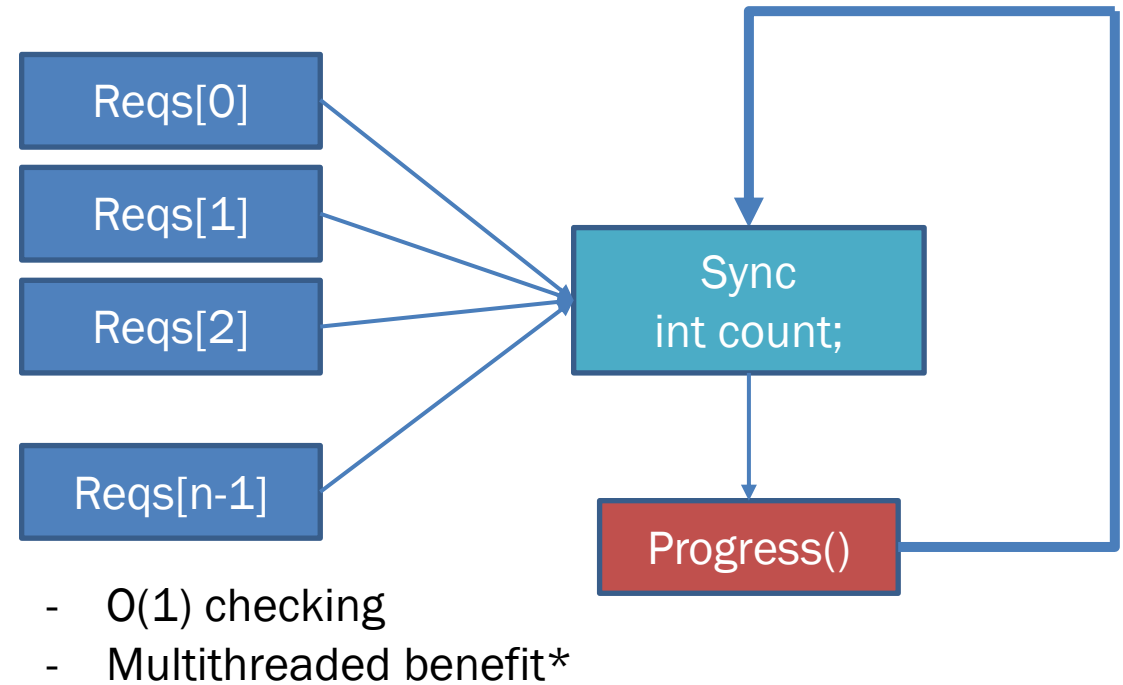
```
MPI_Isend(..., reqs[0]);  
MPI_Isend(..., reqs[1]);  
MPI_Isend(..., reqs[2]);  
..  
..  
MPI_Isend(..., reqs[n-1]);  
  
MPI_Waitall(n, &reqs, &statuses);
```



# MPI Completion – Sync object



Previous OMPI implementation



Current OMPI implementation

# MPI Completion – Bad case

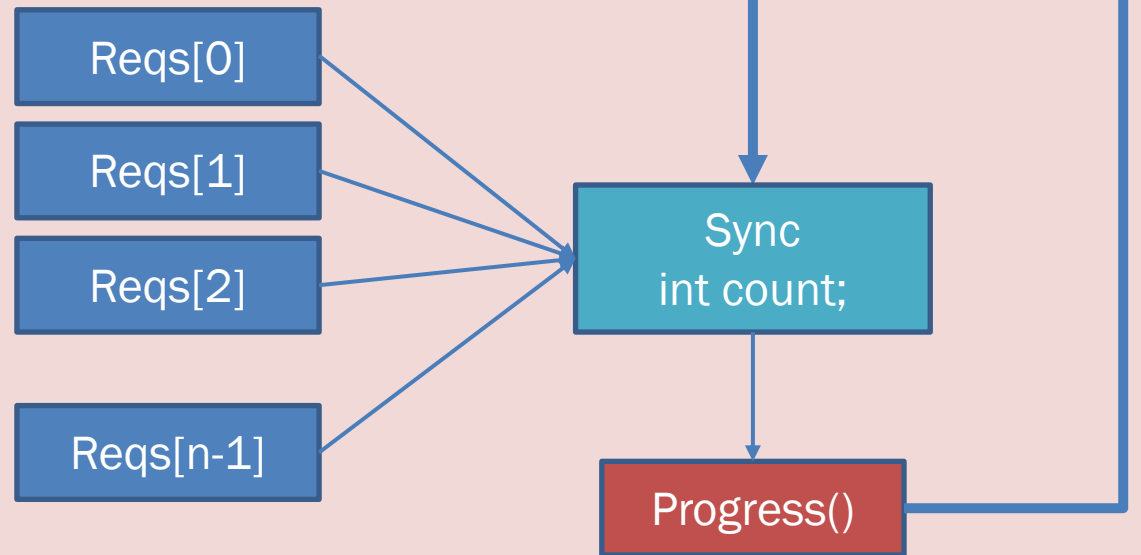
```
MPI_Isend(..., reqs[0]);  
MPI_Isend(..., reqs[1]);  
MPI_Isend(..., reqs[2]);  
..  
..  
MPI_Isend(..., reqs[n-1]);
```

```
MPI_Waitsome(n, &reqs, &statuses);
```

```
...
```

```
MPI_Waitsome(n, &reqs, &statuses);
```

## MPI space



DRAWBACK: If the request is not completed, we *detach* it from sync.

# MPI Completion – Worst case

```
MPI_Init_thread(MULTIPLE);
```

```
MPI_Isend(.., reqs[0]);
```

```
MPI_Isend(.., reqs[1]);
```

```
MPI_Isend(..,reqs[2]);
```

```
..
```

```
..
```

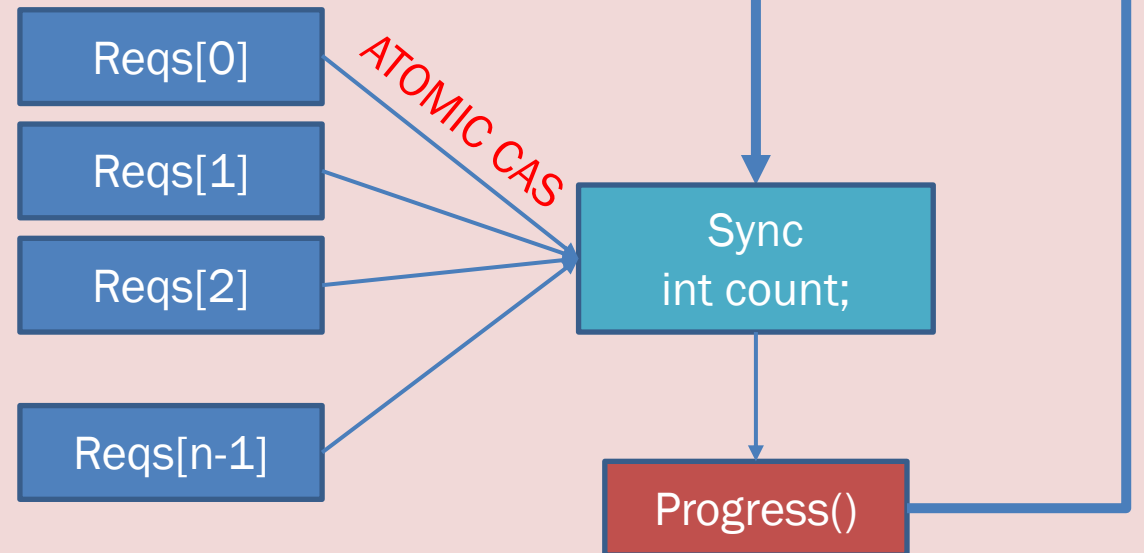
```
MPI_Isend(..,reqs[n-1]);
```

```
MPI_Waitsome(n, &reqs, &statuses);
```

```
...
```

```
MPI_Waitsome(n, &reqs, &statuses);
```

MPI space



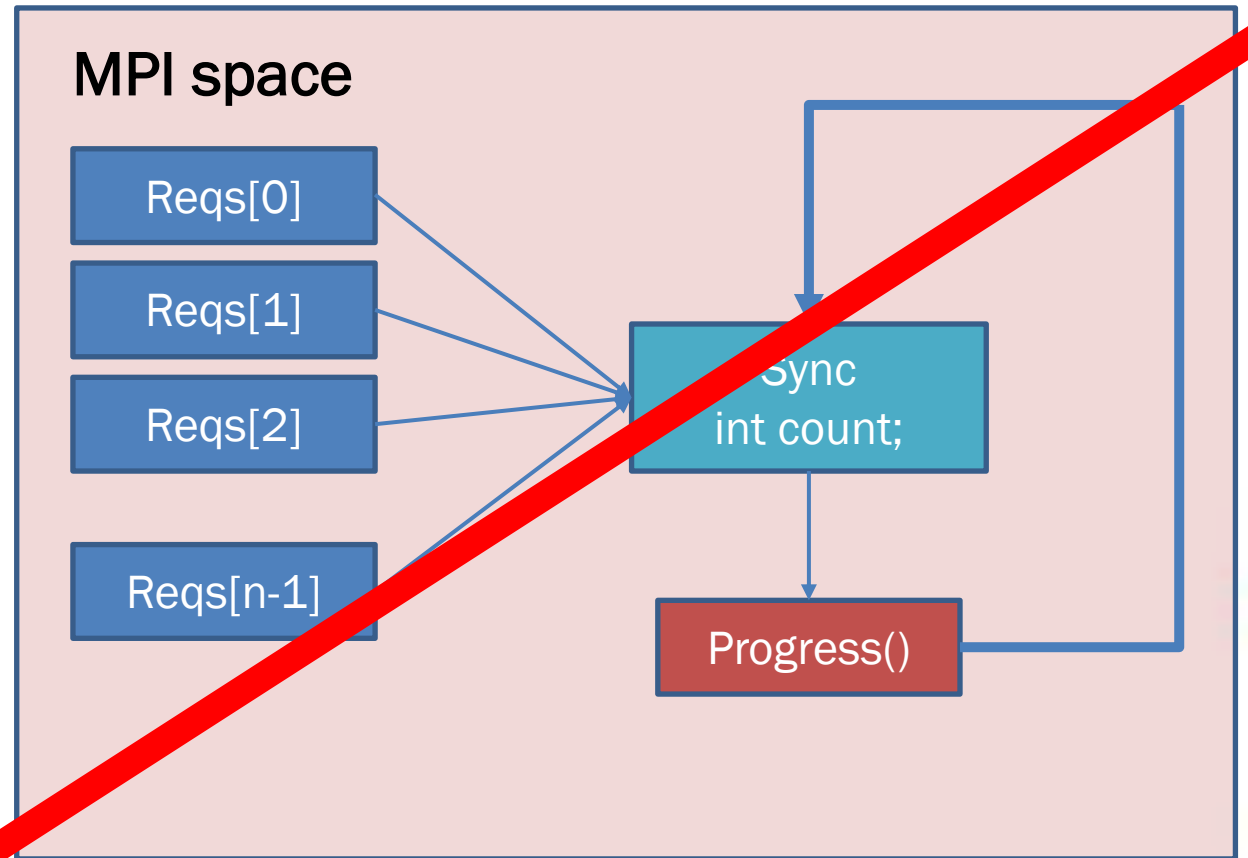
DRAWBACK: If the request is not completed, we *detach* it from sync.

# MPI Completion – Not for MPI\_Test\*

\*too costly to create an object for single use.

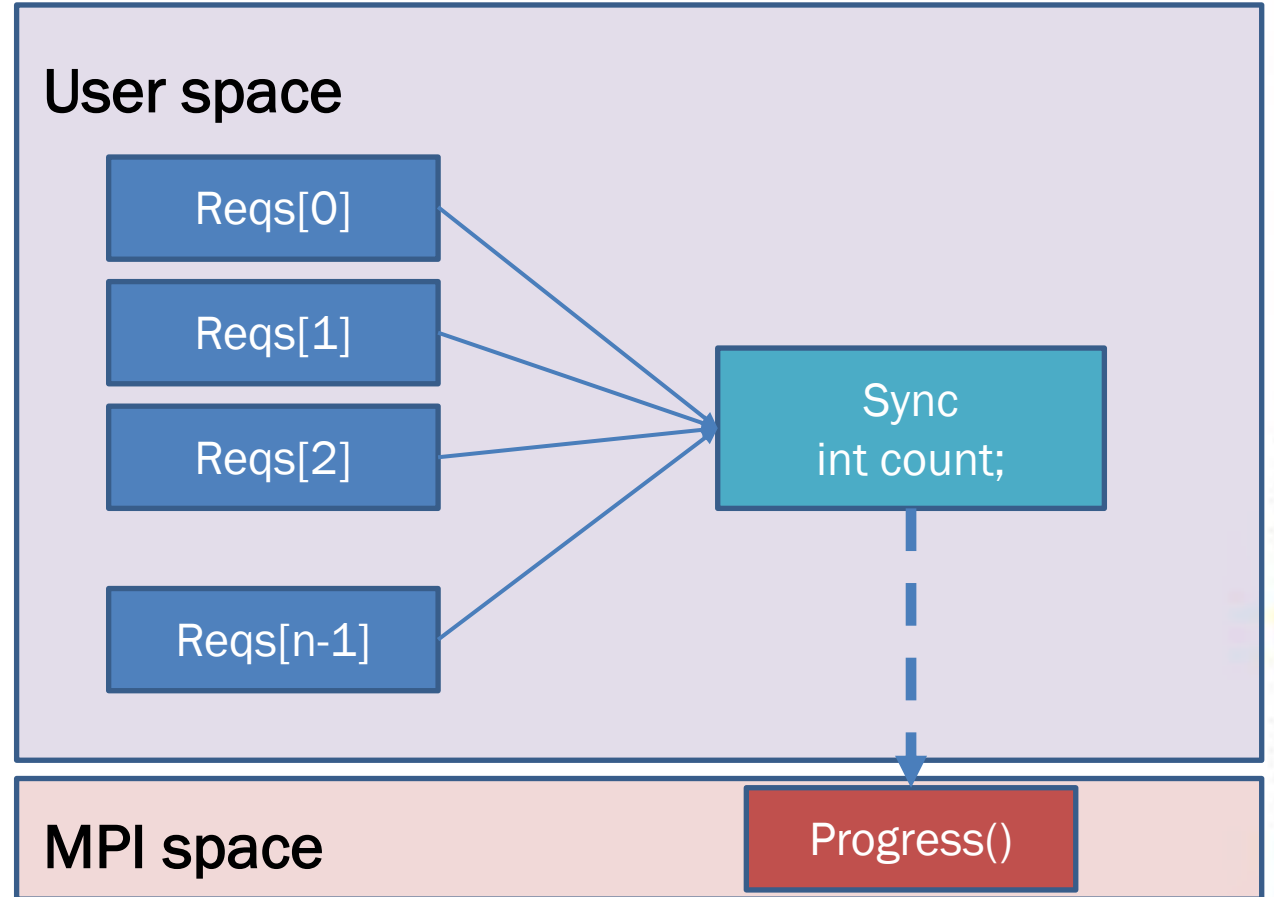
```
MPI_Isend(..., reqs[0]);  
MPI_Isend(..., reqs[1]);  
MPI_Isend(..., reqs[2]);  
..  
..  
MPI_Isend(..., reqs[n-1]);
```

```
MPI_Testsome(n, &reqs, &statuses);  
...  
MPI_Testsome(n, &reqs, &statuses);
```



# MPIX\_Sync- Design

- Give users control of the sync object.
  - Reusable (No need for array rebuilding)
  - Relieves MPI from detaching.
  - Support user-level callback function\*.



# MPIX\_Sync- API

```
int MPIX_Sync_init(MPIX_Sync *sync);
```

```
int MPIX_Sync_attach(  
    MPIX_Sync sync,  
    MPI_Request request,  
    void *completion_data);
```

```
void* MPIX_Sync_query(  
    MPIX_Sync sync,  
    MPI_Status *status);
```

```
int MPIX_Sync_query_bulk(  
    int incount;  
    MPIX_Sync sync,  
    int *outcount;  
    void **completion_data  
    MPI_Status *status);
```

Initialize the object.

Attach the request with sync object and *completion data*.

Return the first *completion\_data* associated with sync object.

Return n *completion\_data* associated with sync object. (testsome API)



# MPIX\_Sync- API

```
int MPIX_Sync_waitall(MPIX_Sync sync);  
*not sure if needed.
```

Wait for all attached requests to complete.

```
int MPIX_Sync_testall(MPIX_Sync sync);  
*not sure if needed.
```

Check if all attached requests are completed.

```
int MPIX_Sync_size(MPIX_Sync sync);
```

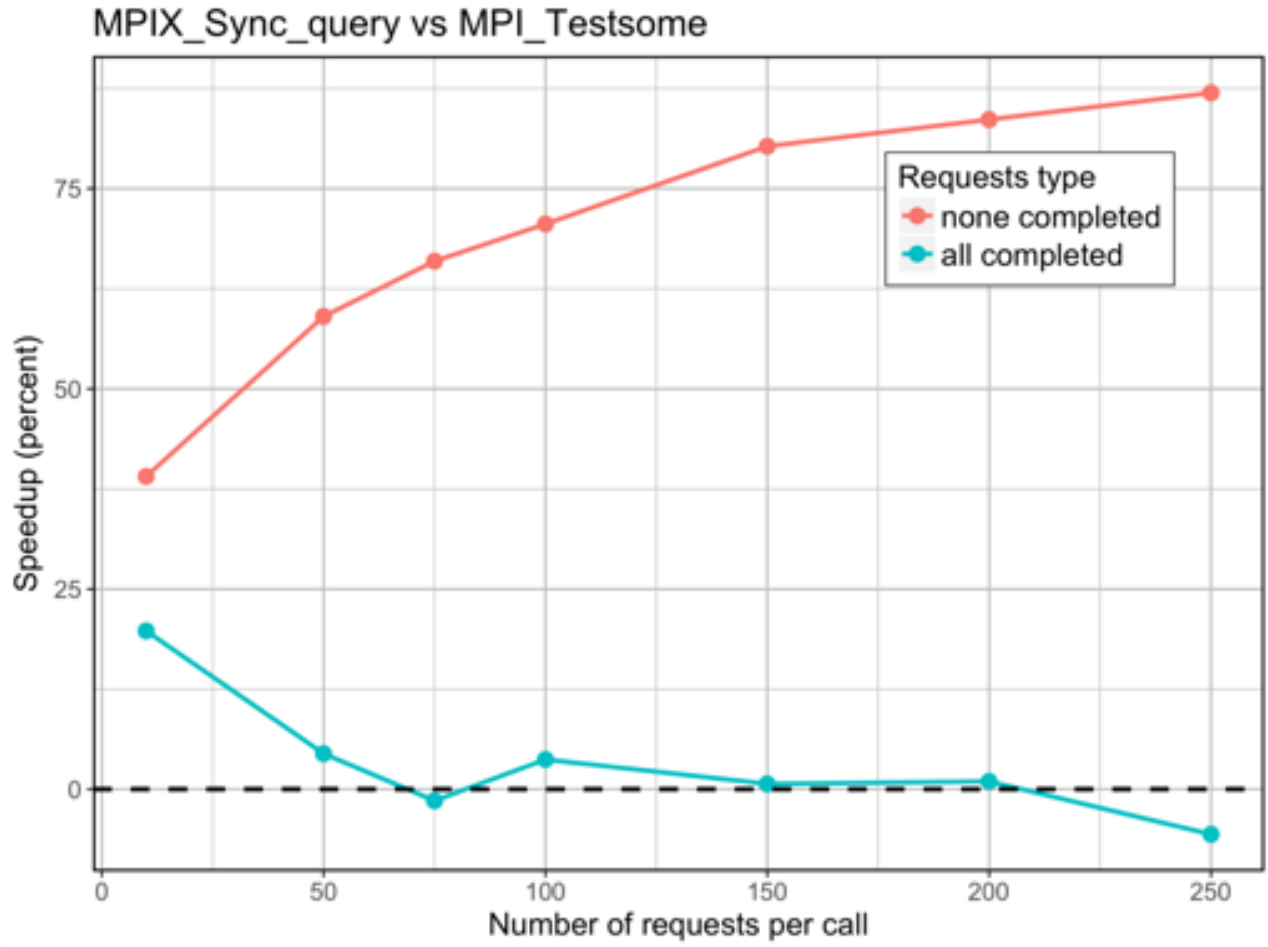
Get the number of incomplete requests.

```
int MPIX_Sync_probe(MPIX_Sync sync);
```

Get the number of *unconsumed* completed requests.

```
int MPIX_Sync_free(MPIX_Sync *sync);
```

# MPIX\_Sync\_query – speedup



- Faster turnaround time if nothing completed  $\sim O(1)$ .
- Slightly slower if every request completed. (Tuning possible)
- MPIX\_Sync\_query\_bulk gives slightly better performance.
- API provides more flexibilities.

# MPIX\_Sync\_query – Application

\*prelim results from PaRSEC (Runtime task scheduler)  
- ~10-30 persistent requests + MPI\_Testsome()

Pingpong (communication only)

message size (bytes)	speedup
400	13.73 %
4000	7.15 %
40000	3.1 %
400000	2.55 %

DPLASMA (matrix factorization)

kernel	speedup
dpotrf	~0%*
dgeqrf	~0%*

\*within error margin

# MPIX\_Sync\_query – User level callback

```
int MPIX_Sync_init(MPIX_Sync *sync);
```

```
int MPIX_Sync_attach(  
    MPIX_Sync sync,  
    MPI_Request request,  
    void *completion_data);
```

```
void* MPIX_Sync_query(  
    MPIX_Sync sync,  
    MPI_Status *status);
```

```
int MPIX_Sync_query_bulk(  
    int incount;  
    MPIX_Sync sync,  
    int *outcount;  
    void **completion_data  
    MPI_Status *status);
```

## Discussion:

- Users can already call the function after query.  
(Passing a structure as *completion\_data*)  
\*more control to the user
- MPI can call callback function as soon as ..
  - A request is completed
  - Certain number of requests are completed.  
\*might affect communication performance



# MPIX\_Sync- API

```
int MPIX_Sync_init(MPIX_Sync *sync);
```

```
int MPIX_Sync_attach(  
    MPIX_Sync sync,  
    MPI_Request request,  
    void *completion_data);
```

```
void* MPIX_Sync_query(  
    MPIX_Sync sync,  
    MPI_Status *status);
```

```
int MPIX_Sync_query_bulk(  
    int incount;  
    MPIX_Sync sync,  
    int *outcount;  
    void **completion_data  
    MPI_Status *status);
```

- *completion\_data* can be  
MPIX\_SYNC\_NO\_COMPLETION\_DATA.

- Return MPIX\_SYNC\_EMPTY if there is no  
completion.

- *incount* can be any number.  
(unlike MPI\_Testsome where it has to be *nrequests*)

- Return MPIX\_SYNC\_MORE if there is more  
completion than asked for.

# MPIX\_Sync- API

```
int MPIX_Sync_init(MPIX_Sync *sync);
```

```
int MPIX_Sync_attach(  
    MPIX_Sync sync,  
    MPI_Request request,  
    void *completion_data);
```

```
void* MPIX_Sync_query(  
    MPIX_Sync sync,  
    MPI_Status *status);
```

```
int MPIX_Sync_query_bulk(  
    int incount;  
    MPIX_Sync sync,  
    int *outcount;  
    void **completion_data  
    MPI_Status *status);
```

- User relinquish the request in attach.

- We free the request here.

Users are responsible to keep track of the requests before passing it to attach if they want to *cancel* them.