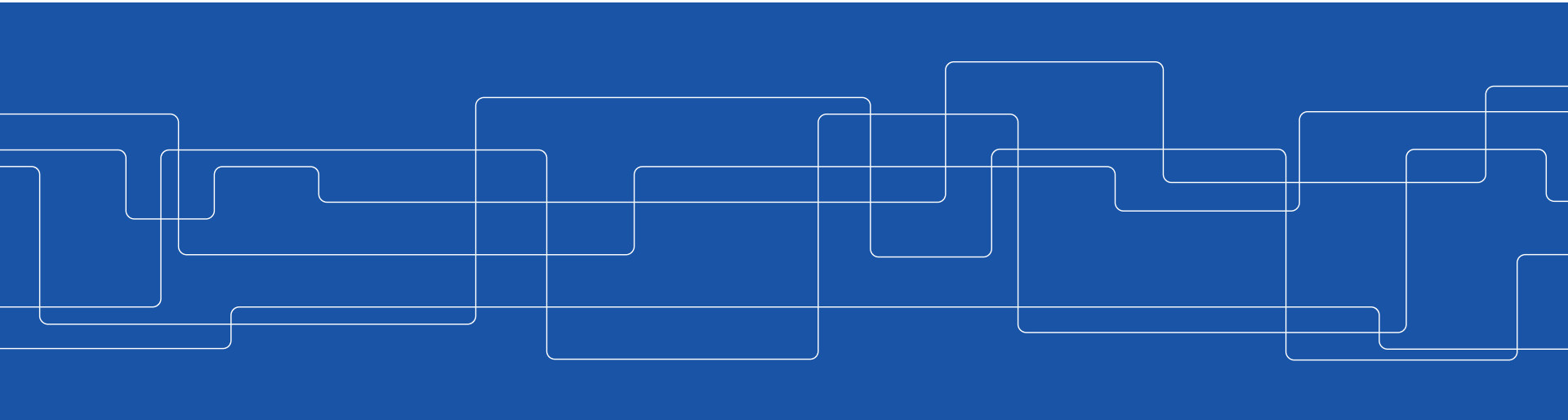# Introduction to MPI I/O

Stefano Markidis

*KTH Royal Institute of Technology*

# What does Parallel I/O Mean?

- **At the program level:**
  - Concurrent reads or writes from multiple processes to a common file
- **At the system level:**
  - A parallel file system and hardware that support such concurrent access
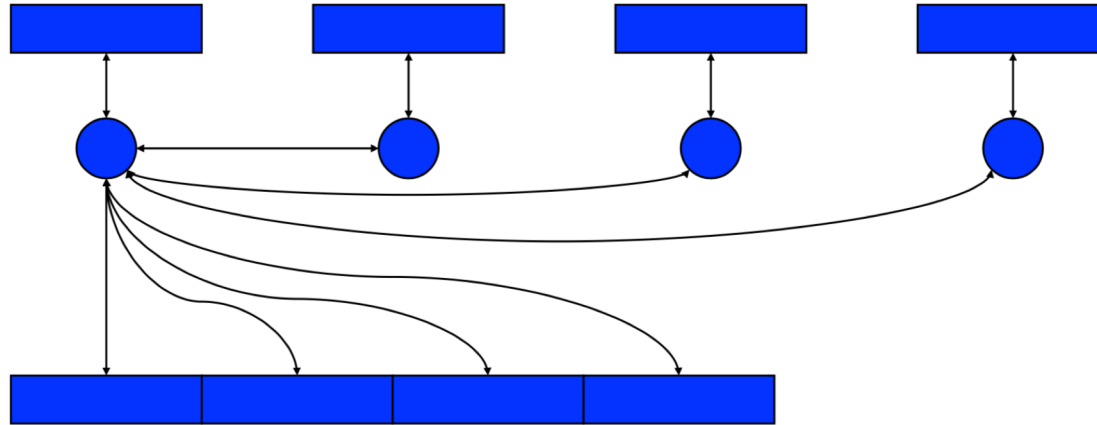
# Parallel I/O in MPI

- Why do I/O in MPI? Why not just POSIX?
    - Parallel performance
    - Single file (instead of one file / process)
- MPI has replacement functions for POSIX I/O
    - Provides migration path
- Multiple styles of I/O can all be expressed in MPI
    - Including some that cannot be expressed without MPI

# Why MPI is a Good Setting for Parallel I/O

- Writing is like sending and reading is like receiving
- Any parallel I/O system will need:
  - user-defined datatypes to describe both memory and file layout
  - non-blocking operations
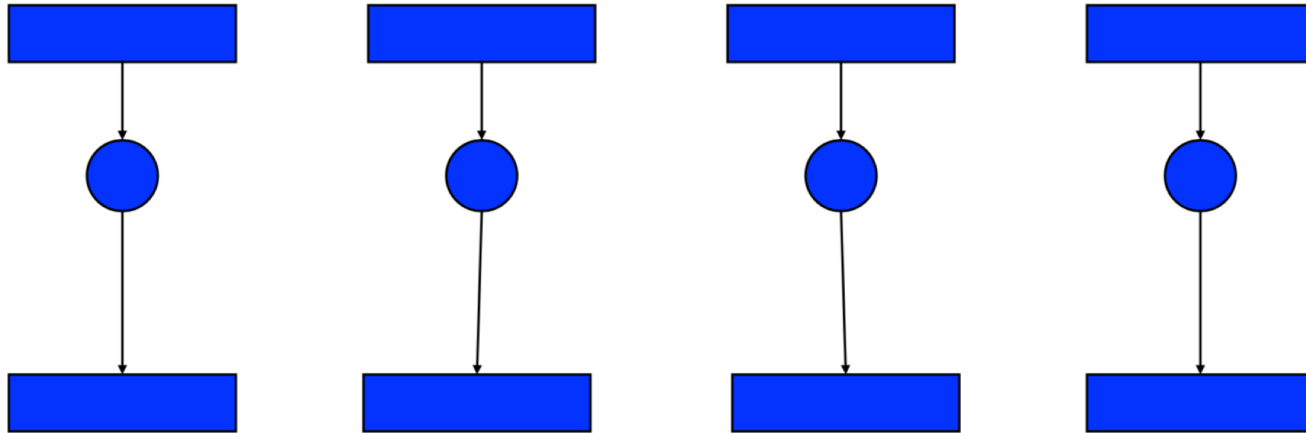  - collective operations

# Non Parallel I/O with MPI I/O



- Non-parallel
- Performance worse than sequential
- Legacy from before application was parallelized

# Example of non-parallel I/O

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include "mpi.h"
4   #define MASTER 0
5   #define BUFSIZE 1024
6
7   int main(int argc, char *argv[])
8   {
9     int rank;
10    int size;
11    int dest;              /* destination rank for message */
12    int source;           /* source rank of a message */
13    int tag = 0;          /* scope for adding extra information to a message *
14    MPI_Status status;    /* struct used by MPI_Recv */
15    int buf;  /* buffer is a single int */
16    FILE *fh;
17
18    MPI_Init(&argc, &argv);
19    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
20    MPI_Comm_size(MPI_COMM_WORLD, &size);
21
22    buf = rank * 100 + 1;  /* all ranks set buf to a unique value */
23    if (rank != MASTER) {
24      dest = MASTER;
25      MPI_Send(&buf, 1, MPI_INT, dest, tag, MPI_COMM_WORLD);
26    }
27    else {
28      fh = fopen("collated.txt", "w");
29      source=0;
30      fprintf(fh,"rank %d: %d\n", source, buf);  /* MASTER'S buf value */
31      for (source=1; source<size; source++) {
32        MPI_Recv(&buf, 1, MPI_INT, source, tag, MPI_COMM_WORLD, &status);
33        fprintf(fh,"rank %d: %d\n", source, buf);
34      }
35      fclose(fh);
36    }
37
38    MPI_Finalize();
39    return EXIT_SUCCESS;
40  }
41
```

**6**

# Independent Parallel I/O



- Pro: parallelism
- Con: lots of small files to manage
- Legacy from before MPI I/O
  - MPI or not

# Example of Independent Parallel I/O (No MPI)

```c
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"

#define MASTER 0
#define BUFSIZE 1024

int main(int argc, char *argv[])
{
  int rank;
  int size;
  int buf;  /* buffer is a single int */
  char outfile[BUFSIZE];
  FILE *fh;

  MPI_Init(&argc, &argv);
  MPI_Comm_rank(MPI_COMM_WORLD, &rank);
  MPI_Comm_size(MPI_COMM_WORLD, &size);

  buf = rank * 100 + 1;  /* all ranks set buf to a unique value */

  sprintf(outfile,"individual-%d.txt", rank);
  fh = fopen(outfile, "w");
  fprintf(fh,"rank %d: %d\n", rank, buf);
  fclose(fh);

  MPI_Finalize();

  return EXIT_SUCCESS;
}
```

# Independent Parallel I/O with MPI

- Just like POSIX I/O, you need to
  - Open the file
  - Read or Write data to the file
  - Close the file
- In MPI, these steps are almost the same:
  - Open the file: MPI_File_open
  - Write to the file: MPI_File_write / MPI_File_read
  - Close the file: MPI_File_close

# Writing to a File with MPI I/O

- MPI_File_open: use MPI_MODE_WRONLY or MPI_MODE_RDWR as the flags. If the file doesn't exist previously, the flag MPI_MODE_CREATE must also be passed to MPI_File_open
    - We can pass multiple flags by using bitwise-or '|' in C, or addition '+" in Fortran

- Use MPI_File_write to write to file.
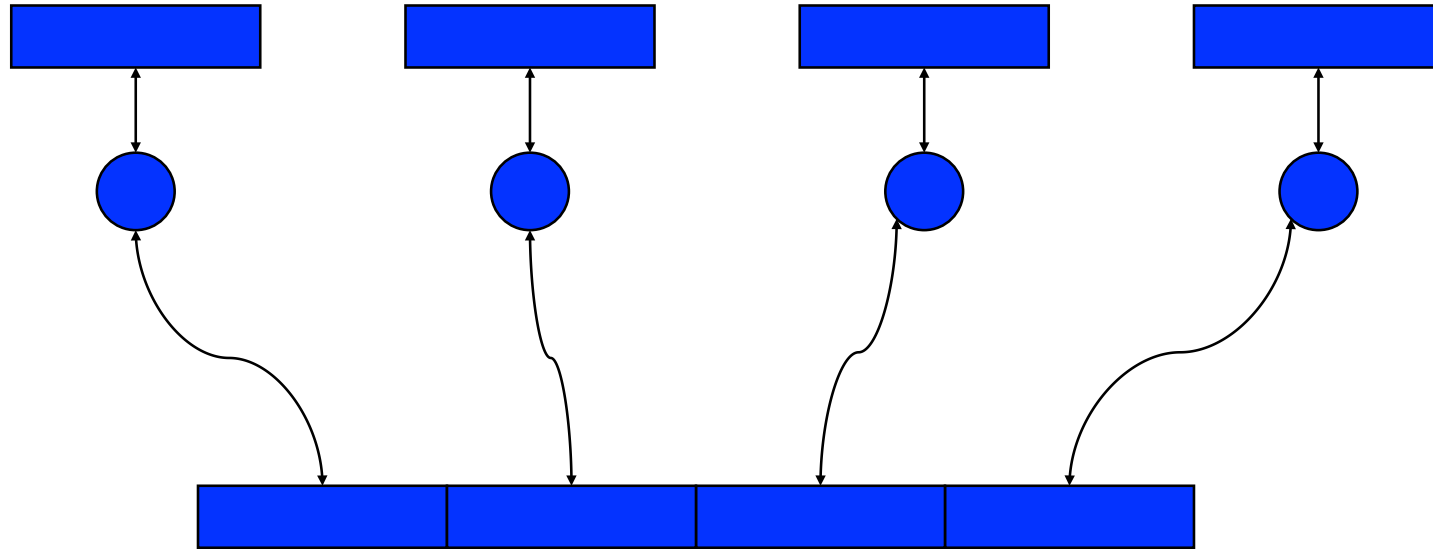
# Example of Independent I/O with MPI

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3   #include "mpi.h"
4
5   #define BUFSIZE 1024
6
7   int main(int argc, char *argv[])
8   {
9     int rank;
10    int size;
11    int buf;   /* buffer is a single int */
12    char outfile[BUFSIZE];
13
14    MPI_File file_handle = NULL;    /* parallel file handle */
15
16    MPI_Init(&argc, &argv);
17    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
18    MPI_Comm_size(MPI_COMM_WORLD, &size);
19
20    buf = rank * 100 + 1;   /* all ranks set buf to a unique value */
21
22    sprintf(outfile,"MPI-individual-%d.txt", rank);
23    /* all processes open the specified file for writing only */
24    MPI_File_open(MPI_COMM_SELF,outfile,MPI_MODE_CREATE | MPI_MODE_WRONLY, MPI_INFO_NULL, &file_handle);
25
26    /* now we are clear to write to the file */
27    MPI_File_write(file_handle, &buf, 1, MPI_INT, MPI_STATUS_IGNORE);
28
29    /* close the file when we're done */
30    MPI_File_close(&file_handle);
31
32
33    MPI_Finalize();
34
35    return EXIT_SUCCESS;
36  }
```

- File Open is collective over the communicator Modes similar to Unix open
  MPI_Info provides additional hints for performance

- File Write is independent

- Many important variations covered in later slides

- File close is collective; similar in style to MPI_Comm_free

Why MPI I/O for Independent I/O?
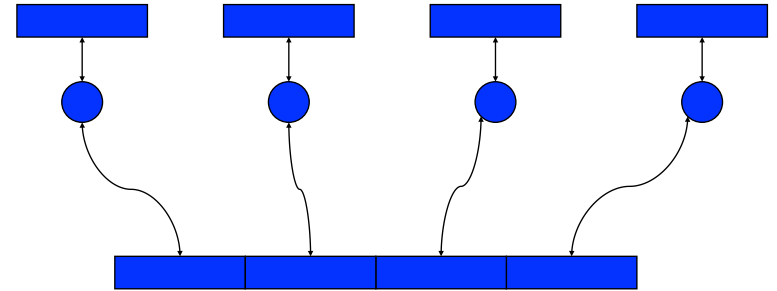
# Cooperative Parallel I/O – Single File



- Parallelism
  - Can only be expressed in MPI

# File Views

- Use MPI File View to tell each process which part of the file is allowed to write to or read from
  - Described in MPI with an offset and an MPI_Datatype

- Specified by a triplet (*displacement*, *etype*, and *filetype*) passed to **MPI_File_set_view**
  - *displacement* = number of bytes to be skipped from the start of the file
  - *etype* = basic unit of data access (can be any basic or derived datatype)
  - *filetype* = etype if contiguous access or MPI derived data type for non-contiguous access

13

# Example of Writing with File View

```c
#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include "mpi.h"

#define MASTER 0
#define BUFSIZE 10        /* size of storage array we'll use on each process */
#define ALPHSIZE 26       /* how many chars in alphabet */
int main(int argc, char* argv[])
{
  int rank;               /* rank of process */
  int size;               /* number of processes started */
  int ii;                 /* simple counter */
  char buf[BUFSIZE];      /* values each process will set and write to file */
  char alphabet[ALPHSIZE] = "abcdefghijklmnopqrstuvwxyz";
  MPI_File file_handle = NULL;    /* parallel file handle */
  /* initialise processes */
  MPI_Init( &argc, &argv );
  MPI_Comm_size( MPI_COMM_WORLD, &size );
  MPI_Comm_rank( MPI_COMM_WORLD, &rank );
  /* arrange for a (looping) sequence of characters */
  for (ii=0; ii<BUFSIZE; ii++) {
    buf[ii] = alphabet[(ALPHSIZE + rank) % ALPHSIZE];
  }

  /* all processes open the specified file for writing only */
  MPI_File_open(MPI_COMM_WORLD, "view.txt",MPI_MODE_CREATE | MPI_MODE_WRONLY, MPI_INFO_NULL, &file_handle);
  /* establish a different 'view' of the file for each process */
  MPI_File_set_view(file_handle, (rank * BUFSIZE * sizeof(char)),
      MPI_CHAR, MPI_CHAR, "native", MPI_INFO_NULL);
  /* now we are clear to write to the file */
  MPI_File_write(file_handle, buf, BUFSIZE, MPI_CHAR, MPI_STATUS_IGNORE);
  /* close the file when we're done */
  MPI_File_close(&file_handle);

  MPI_Finalize();

  return EXIT_SUCCESS;
}
```

Blocking Write
(What does it mean?)
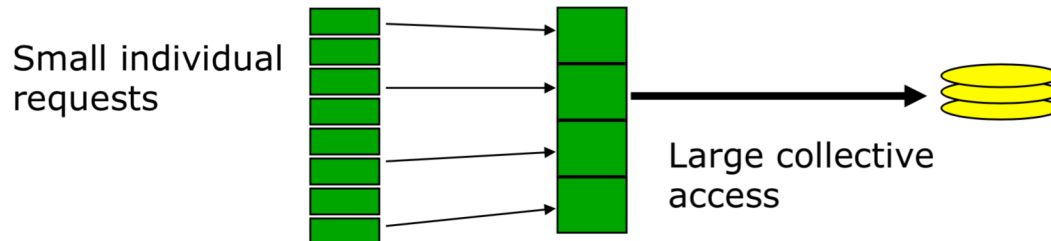
# Non-blocking I/O Operations

```c
13  #include <stdio.h>
14  #include <string.h>
15  #include <stdlib.h>
16  #include "mpi.h"
17  #define MASTER 0
18  #define BUFSIZE 10          /* size of storage array we'll use on each process */
19  #define ALPHSIZE 26         /* how many chars in alphabet */
20  int main(int argc, char* argv[])
21  {
22    int rank;                 /* rank of process */
23    int size;                 /* number of processes started */
24    int ii;                   /* simple counter */
25    char buf[BUFSIZE];        /* values each process will set and write to file */
26    char alphabet[ALPHSIZE] = "abcdefghijklmnopqrstuvwxyz";
27    MPI_File file_handle;     /* parallel file handle */
28    MPI_Request request;      /* request for action used by async functions */
29    /* initialise processes */
30    MPI_Init( &argc, &argv );
31    MPI_Comm_size( MPI_COMM_WORLD, &size );
32    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
33    /* arrange for a (looping) sequence of characters */
34    for (ii=0; ii<BUFSIZE; ii++) {
35      buf[ii] = alphabet[(ALPHSIZE + rank) % ALPHSIZE];
36    }
37    /* all processes open the specified file for writing only */
38    MPI_File_open(MPI_COMM_WORLD, "iwrite.txt",
39      MPI_MODE_CREATE | MPI_MODE_WRONLY,
40      MPI_INFO_NULL, &file_handle);
41    /* establish a different 'view' of the file for each process */
42    MPI_File_set_view(file_handle, (rank * BUFSIZE * sizeof(char)),
43          MPI_CHAR, MPI_CHAR, "native", MPI_INFO_NULL);
44
45    /* now we are clear to write to the file */
46    MPI_File_iwrite(file_handle, buf, BUFSIZE, MPI_CHAR, &request);
47    /* >>> we could get on with something else here <<< */
48
49    /* we must wait for the async request to be completed */
50    MPI_Wait(&request, MPI_STATUS_IGNORE);
51    /* close the file when we're done */
52    MPI_File_close(&file_handle);
53    MPI_Finalize();
54    return EXIT_SUCCESS;
55  }
```

← What is this equivalent to?

# Collective I/O and MPI

- A critical optimization in parallel I/O
- All processes (in the communicator) must call the collective I/O function
  - Allows communication of "big picture" to file system
- Framework for I/O optimizations at the MPI-IO layer
  - Basic idea: build large blocks, so that reads/writes in I/O
- **Requests from different processes may be merged together**



Small individual requests

Large collective access

# Collective I/O Functions

- **MPI_File_write_all**
  - **_all** indicates that all processes in the group specified by the communicator passed to **MPI_File_open** will call this function
  - Each process specifies only its own **access information** - the argument list is the same as for the non-collective functions

How do we specify *access information*?

# Example Collective MPI I/O

```c
13  #include <stdio.h>
14  #include <string.h>
15  #include <stdlib.h>
16  #include "mpi.h"
17
18  #define MASTER 0
19  #define BUFSIZE 10          /* size of storage array we'll use on each process */
20  #define ALPHSIZE 26         /* how many chars in alphabet */
21
22  int main(int argc, char* argv[])
23  {
24    int rank;                 /* rank of process */
25    int size;                 /* number of processes started */
26    int ii;                   /* simple counter */
27    char buf[BUFSIZE];        /* values each process will set and write to file */
28    char alphabet[ALPHSIZE] = "abcdefghijklmnopqrstuvwxyz";
29    MPI_File file_handle = NULL;    /* parallel file handle */
30    /* initialise processes */
31    MPI_Init( &argc, &argv );
32    MPI_Comm_size( MPI_COMM_WORLD, &size );
33    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
34    /* arrange for a (looping) sequence of characters */
35    for (ii=0; ii<BUFSIZE; ii++) {
36      buf[ii] = alphabet[(ALPHSIZE + rank) % ALPHSIZE];
37    }
38    /* all processes open the specified file for writing only */
39    MPI_File_open(MPI_COMM_WORLD, "view.txt",MPI_MODE_CREATE | MPI_MODE_WRONLY, MPI_INFO_NULL, &file_handle);
40    /* establish a different 'view' of the file for each process */
41    MPI_File_set_view(file_handle, (rank * BUFSIZE * sizeof(char)),
42          MPI_CHAR, MPI_CHAR, "native", MPI_INFO_NULL);
43
44    /* use collective I/O*/
45    MPI_File_write_all(file_handle, buf, BUFSIZE, MPI_CHAR, MPI_STATUS_IGNORE);
46
47    /* close the file when we're done */
48    MPI_File_close(&file_handle);
49
50    MPI_Finalize();
51
52    return EXIT_SUCCESS;
53  }
```

# **Summary**

- MPI I/O is a convenient way to express both independent (many files) and cooperative (shared file) parallel I/O

- MPI cooperative parallel I/O might use MPI File Views and supports different types of I/O as different communication:
    - Blocking
    - Non-blocking
    - Collective