Fig. 00.

# Parallel I/O for SwissTx
## Emin Gabrielyan
## Prof. Roger D. Hersch
## Peripheral Systems Laboratory
## Ecole Polytechnique Fédérale de Lausanne
## Switzerland

- Introduction

- SFIO (Striped File I/O)
  software architecture

- SFIO interface

- Performance measurements

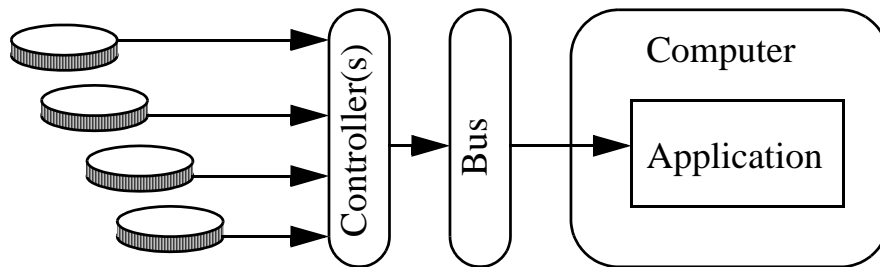- Integration into MPI-II I/O

- Conclusion

Fig. 01.

# Parallel I/O for Swiss-Tx (Introduction)

- Requirements for Parallel I/O

- Striped Files Parallel Access Solution

- Interface

- Network Communication and
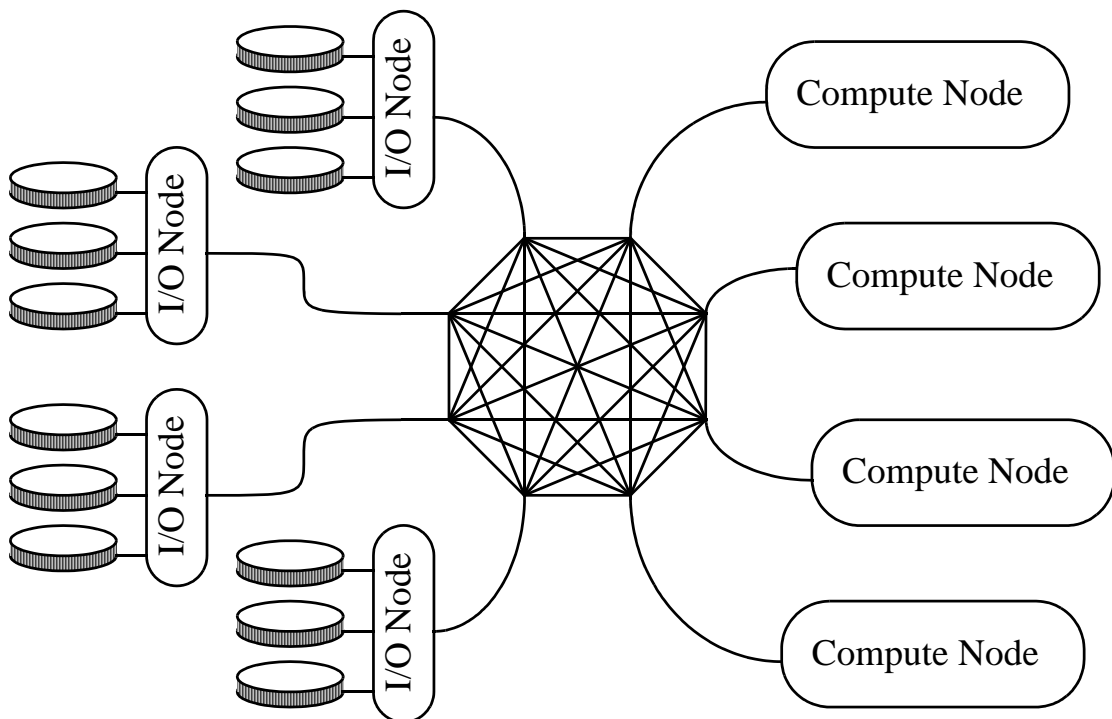  Disk Access optimisation

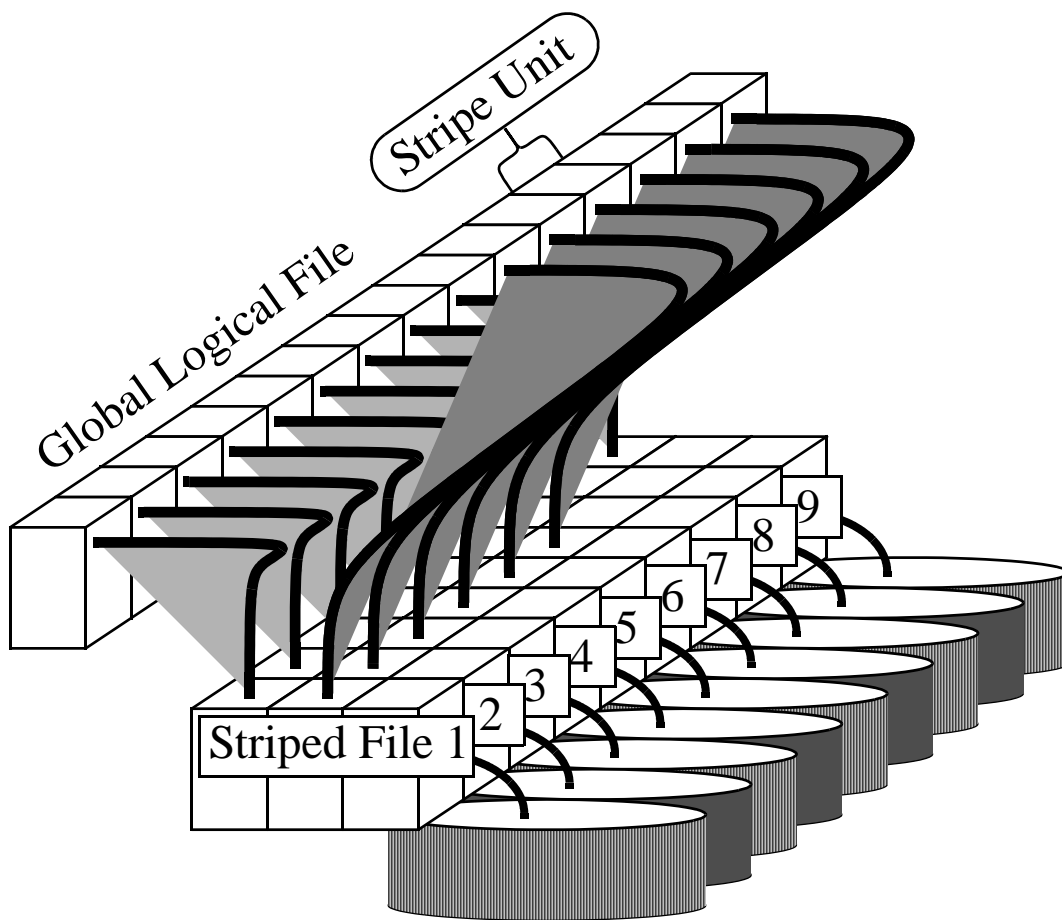# Requirements to Parallel I/O System

● Scalable Throughput



● Concurent Access



number of I/O devices. The upper limit of throughput can be the controller(s)/bus bandwidth or the TOTAL network bandwidth.

Fig. 03. The basic idea of our implementation is cyclical distribution of logical file

# Striped Files Parallel Access Solution

Stripe Unit

Global Logical File

Striped File 1

2 3 4 5 6 7 8 9

accross the set of striped files. Access to single part of logical file require parallel access to set of striped files.

Fig. 04. The native interface is SFIO, but we work to provide also MPI-II I/O inter-

# Interface

## ● SFIO

MFILE* **mopen**(char *s, int unitsz); // "t0-p1,/tmp/a;t0-p2,/tmp/a"

void **mclose**(MFILE *f);

void **mchsize**(MFILE *f, long size);

void **mdelete**(char *s);

void **mcreate**(char *s);

void **mread**(MFILE *f, long offset, char *buf, unsigned count);

void **mwrite**(MFILE *f, long offset, char *buf, unsigned count);

void **mreadb**(MFILE *f, unsigned bcount,
    long Offset[], char *Buf[], unsigned Count[]);

void **mwriteb**(MFILE *f, unsigned bcount,
    long Offset[], char *Buf[], unsigned Count[]);

## ● MPI-II I/O

| | |
|---|---|
| MPI_File_open | MPI_File_write_all |
| MPI_File_set_view | MPI_File_read_all |
| MPI_File_write | MPI_File_write_at_all |
| MPI_File_read | MPI_File_read_at_all |
| MPI_File_write_at | MPI_File_close |
| MPI_File_read_at | MPI_File_delete |

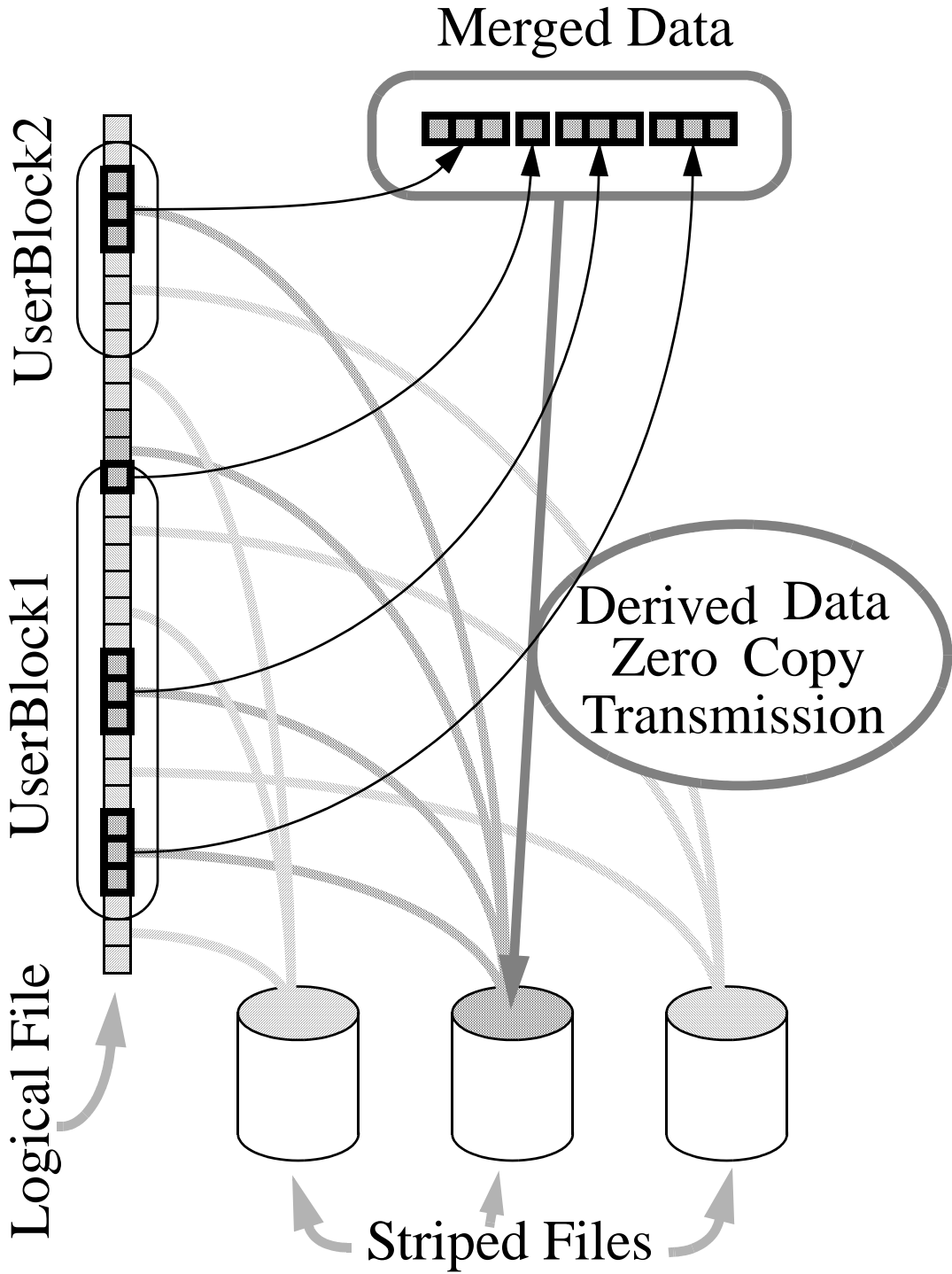face.

# Network and Disk Optimisation

⬤ Network optimisation for noncollective access

⬤ Disk optimisation for noncollective access

⬤ Collective access optimisation

# Network Transmission Optimisation

Merged Data

UserBlock2

UserBlock1

Derived Data
Zero Copy
Transmission

Logical File

Striped Files

# Disk Access Optimisation

Fig. 07. Data fragmentation is resolved for network communication

User Block 2

User Block 2

The images of user blocks 1 and 2 of the global logical file in the local file of specifyed I/O node

Fig. 08. Data fragmentation still exist at disk access

Compute
Node side

Compute Node

Temporary Derived Datatype

| Location 1 | Length 1 |
| Location 2 | Length 2 |
| Location 3 | Length 3 |
| Location 4 | Length 4 |
| Location 5 | Length 5 |
| Location 6 | Length 6 |
| Location 7 | Length 7 |

Singe large MPI transmission

I/O Node

I/O Node
side

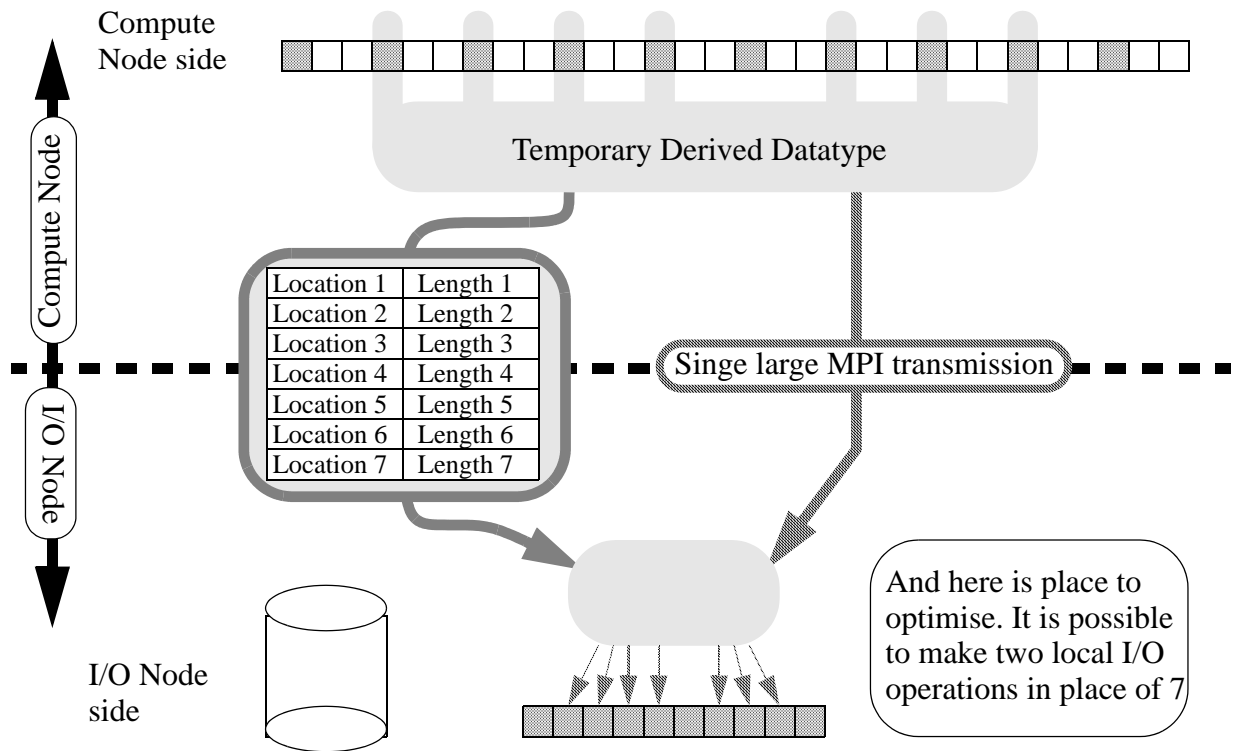And here is place to optimise. It is possible to make two local I/O operations in place of 7

Fig. 09. Data fragmentation is resolved on disk level also

# Disk Access Optimisation

Compute
Node side

Temporary Derived Datatype

| Location 1 | Length 1 |
| Location 2 | Length 2 |
| Location 3 | Length 3 |
| Location 4 | Length 4 |
| Location 5 | Length 5 |
| Location 6 | Length 6 |
| Location 7 | Length 7 |

compute node
side optimization

| Location 1 | Length 1 |
| Location 2 | Length 2 |

MPI transmission

Compute Node

I/O Node

I/O Node side

Two local I/O
oprations in place of 7

Fig. 10. Disk Access Optiisation of Collective Operations

# Disk Access Optimisation of Collective Operations

Logical File

Local Disk

and then writing or reading

waiting for all data, sorting and gathering

Proc1

Proc2

Proc3

Proc4

# Disk Access Optimisation
# of Collective Operations

**Compute Node 1**

Data related to
Specifyed striped file

**Compute Node 2**

Data related to
Specifyed striped file

**Compute Node 3**

Data related to
Specifyed striped file

The data to be accessed by
all of compute nodes in the
specified striped file

operations. Here each compute node tryes to access some part of logical file, but in
the picture is shown the translation of this part to the view of some specifyed striped
file. For each of compute node the striped file is same.

Fig. 12.

# SFIO Software Architecture

● Encapsulation over MPI

● Functionality

● Data organisation and optimisation
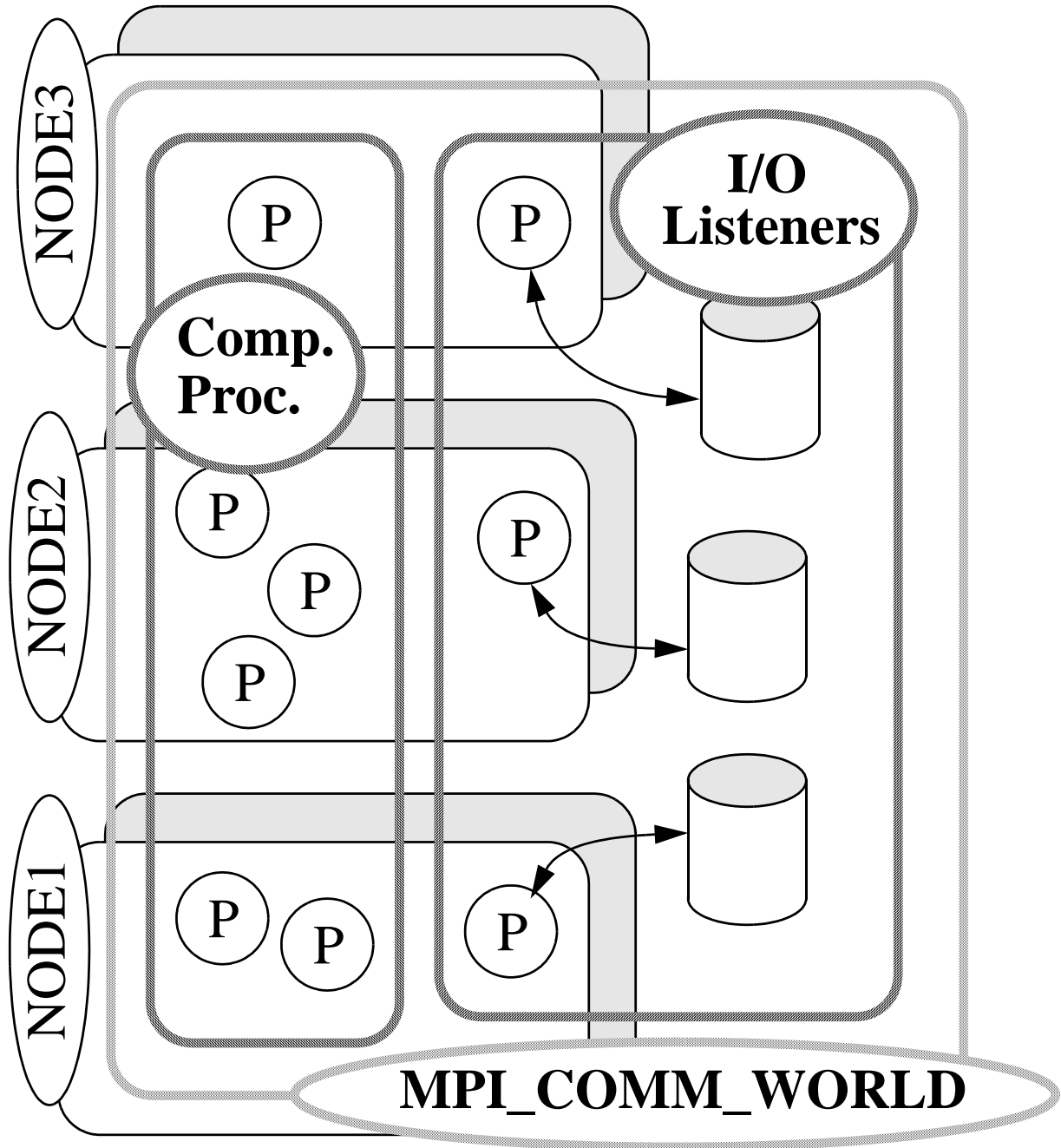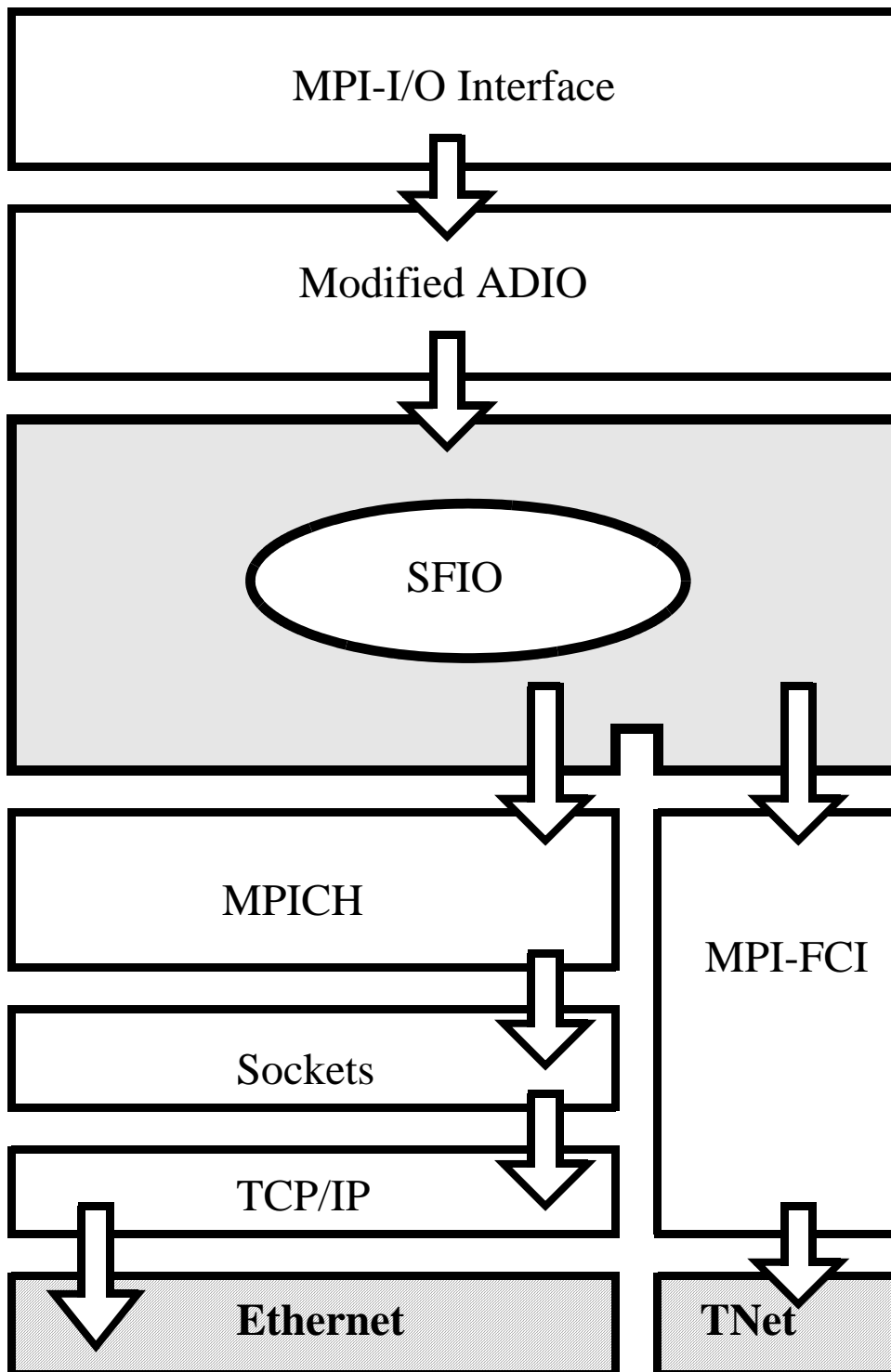
Fig. 13. Encapsulation over MPI

Fig. 14. Functionality. SFIO is implemented and tested with Digital Unix



(MPICH, FCI) and Intel Windos NT (MPICH).

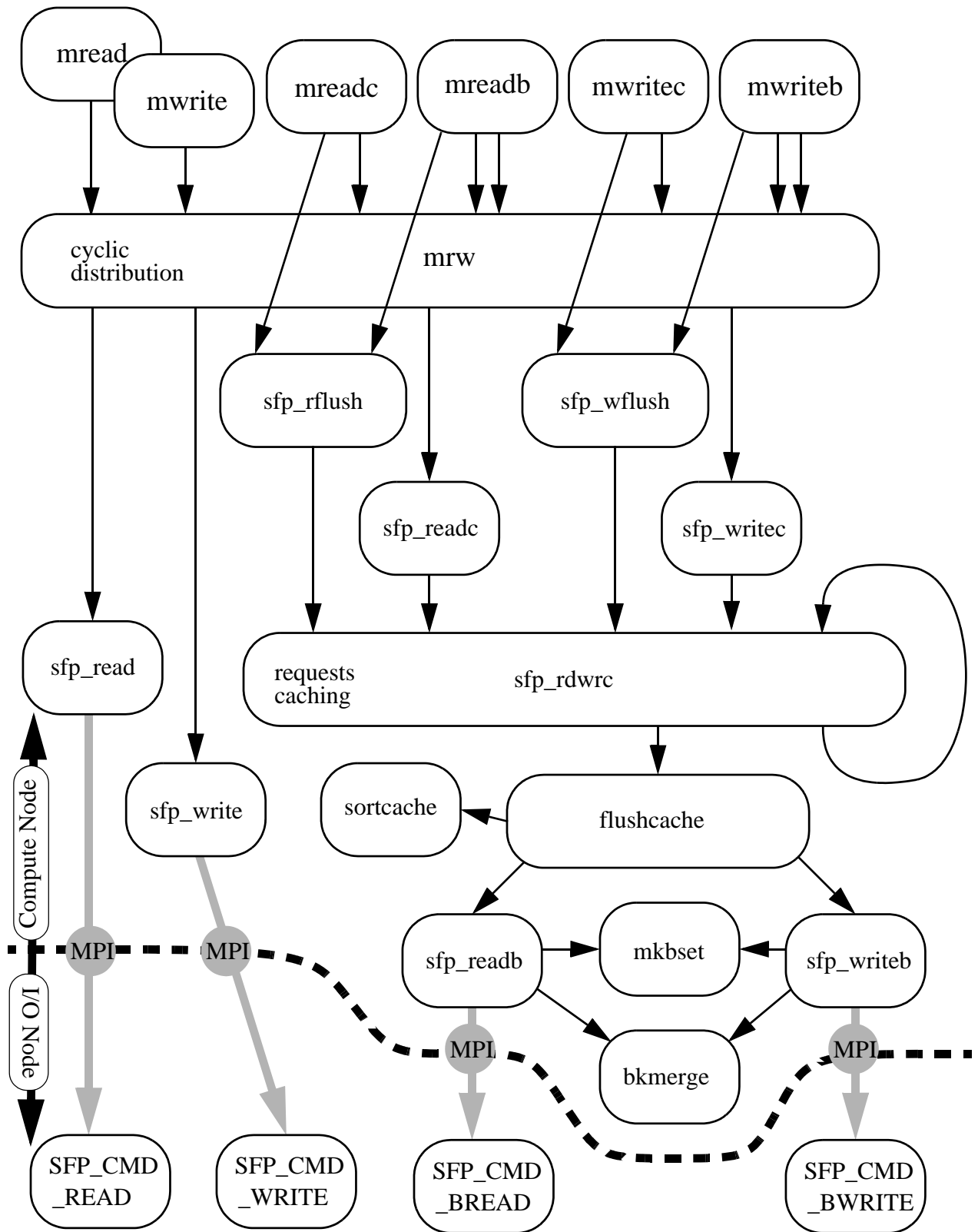Fig. 15. Functionality, data organisation and optimisation

Fig. 16.

SFIO Interface

● Single Block Access Interface

● Multiple Block Access Interface

Fig. 17. If there is one compute node

# Single Block Access Interface

```
#include "/usr/p5/gabriely/mpi/lib/mio.h"

int _main(int argc, char *argv[])
{
    char buf[]="ABCDEFGHIJKLMNOPQRSTUVWXYZ";
    MFILE *f;

    f=mopen( "t0-p1.epfl.ch,/tmp/aa;t0-p2.epfl.ch,/tmp/aa" , 10 );

    mwrite( f , 0 , buf , 26 );

    mclose(f);
    return 0;
}
```

Offset in Logical File

Names of Striped Files

Stripe Unit Size

Fig. 18.  If there are more than one compute nodes

# Single Block Access Interface
# more than one compute processes

```
#include "/usr/p5/gabriely/mpi/lib/mio.h"
int _main(int argc, char *argv[])
{
    char buf[]="abcdefghijklmnopqrstuvwxyz";
    MFILE *f;
    int rank;
    MPI_Comm_rank(_MPI_COMM_WORLD,&rank);

    f=mopen( "t0-p1.epfl.ch,/tmp/aa;t0-p2.epfl.ch,/tmp/aa" , 10 );

    mwrite(f, rank*26 , buf , 26 );

    mclose(f);
    printf("rank=%d\n",rank);
    return 0;
}
```

Offset in Logical File

Names of Striped Files

Buffer

Stripe Unit Size

Buffer Size

Fig. 19.

# Multiple Block Access Interface

```
#include "/usr/p5/gabriely/mpi/lib/mio.h"
int _main(int argc, char *argv[])
{
    char buf1[]="abcdefghijklmnopqrstuvwxyz";
    char buf2[]="ABCDEFGHIJKLMNOPQRSTUVWXYZ";
    char buf3[]="0123456789";
    char* Buf[]={buf1,buf2,buf3};
    long Offset[]={0,26,52};
    unsigned Count[]={26,26,10};
    MFILE *f;

    f=mopen("t0-p1.epfl.ch,/tmp/aa;t0-p2.epfl.ch,/tmp/aa",10);

    mwriteb(f,  3  ,Offset,Buf,Count);

    mclose(f);
    return 0;
}
```

Number of Blocks

Fig. 20.

## Performance Measurements

Data Access Performance on Windows NT cluster (4 I/O nodes). Effect of Optimisation

Data Access Performance on Swiss-Tx SFIO over MPICH with 100Mbps ethernet HUB (7 I/O nodes). Effect of Optimisation

Data Access throughput dependance from stripe unit size (susz) and user block size (bksz) on Windows NT cluster (4 I/O nodes)

Data Access throughput dependance from stripe unit size (susz) and user block size (bksz) on Swiss-Tx SFIO over MPICH with 100Mbps ethernet HUB (7 I/O nodes).

Performance measurements on Swiss-Tx SFIO over MPICH with GigaEthernet Switch (4 I/O nodes) for different number of concurently reading/writing compute processes.

Performance measurements on Swiss-Tx SFIO over MPI-FCI with TNET Crossbar Switch (4 I/O nodes). Effect of Optimisation

Fig. 21. Multiblock interface (WinNT)

**Fig. 22. Multiblock user interface (Sw-tx, stripe unit = 1005 bytes)**

Fig. 23.

Fig. 24. Swiss-Tx 100Mbps MPICH 10MB data, average of 11 measur.

speed (mbyte/s)

Stripe unit size (bytes)

User block size (bytes)

Fig. 25. mwrite 4 I/O nodes Swiss-tx Gigabit ethernet Crossbar switch

Fig. 26. TNET 3 I/O nodes 1 compute node 660Mbyte

Fig. 27.

# Integration into MPI-II I/O

**Portable MPI Model Implementation**
Bill Gropp, Rusty Lusk,
Debbie Swider, Rajeev Thakur
**Argonne National Laboratory**

SFIO

MPICH

MPI-FCI

Sockets

TCP/IP

**Ethernet**

**TNet**

Fig. 28.

Synchronism

noncollective    nonblocking

collective    blocking

Positioning

Coordination

explicit offsets

individual file pointers

shared file pointer

Coordination

Coordination

Synchronism

Coordination

Synchronism

Positioning

Synchronism

Positioning

Synchronism

Fig. 29.



**Coordination**

**Coordination**

**Synchronism**

| READ AT_ALL BEGIN | END |
| WRITE AT_ALL BEGIN | END |

| READ ALL BDEGIN | END |
| WRITE ALL BEGIN | END |

| READ ORDERED BEGIN | END |
| WRITE ORDERED BEGIN | END |

| READ AT_ALL |
| WRITE AT_ALL |

| READ ALL |
| WRITE ALL |

| READ ORDERED |
| WRITE ORDERED |

**Synchronism**

**Positioning**

| IREAD AT |
| IWRITE AT |

| IREAD |
| IWRITE |

| IREAD SHARED |
| IWRITE SHARED |

**Synchronism**

| READ AT |
| WRITE AT |

| READ |
| WRITE |

| READ SHARED |
| WRITE SHARED |

**Synchronism**

**Positioning**

Fig. 30.

mpich/include/mpio.h

typedef struct ADIOI_FileD ***MPI_File**;

mpich/romio/adio/include/adio.h

typedef struct ADIOI_FileD ***ADIO_File**;

mpich/romio/adio/include/adio.h

```
struct ADIOI_FileD {
    int cookie;
    int fd_sys;

    int fd_sys2;

    MFILE *fd_sys3;

    ADIO_Offset fp_ind;
    ADIO_Offset fp_sys_posn;
    ADIOI_Fns *fns;
    MPI_Comm comm;
    char *filename;
    int file_system;
    int access_mode;
    ADIO_Offset disp;
    MPI_Datatype etype;
    MPI_Datatype filetype;
    int etype_size;
    MPI_Info info;
    int async_count;
    int perm;
    int atomicity;
    int iomode;
};
```

Alternative file descriptor

Striped file descriptor

● We modify MPI_File_open operation

● We modify MPI_File_close operation

Fig. 31.

mpich/romio/adio/ad_nfs/ad_nfs_write.c

void **ADIOI_NFS_WriteStrided**(ADIO_File fd,
void *buf, int count, MPI_Datatype datatype,
int file_ptr_type, ADIO_Offset offset,
ADIO_Status *status, int *error_code)

Memory

FileView

File

noncontiguous in memory,
contiguous in file

Memory

FileView

File

noncontiguous in memory
as well as in file

Memory

FileView
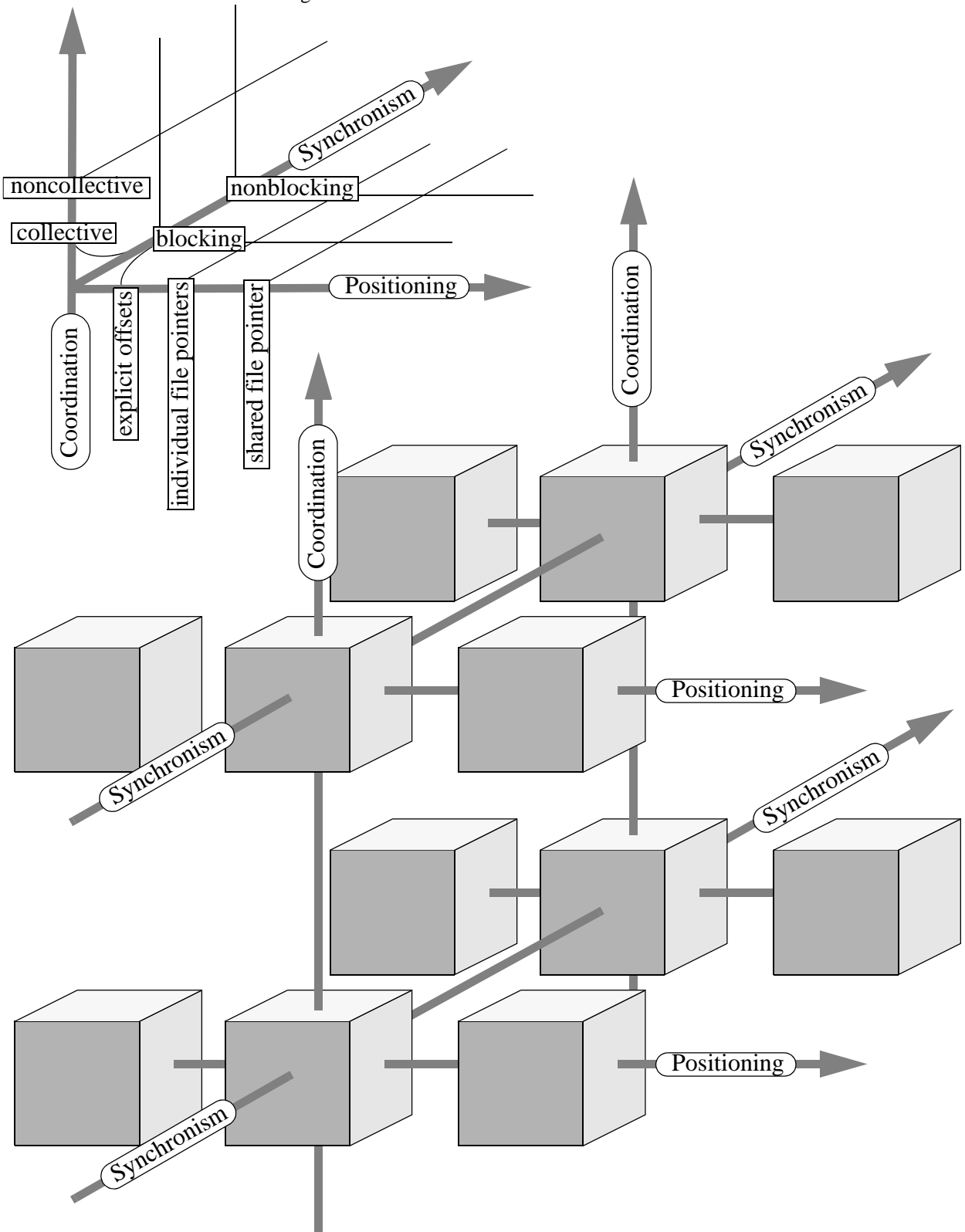
File

contiguous in memory,
noncontiguous in file

for (j=0; j<count; j++)
 for (i=0; i<flat_buf->count; i++)
 {
  userbuf_off = j*buftype_extent + flat_buf->indices[i];
  req_off = off;
  req_len = flat_buf->blocklens[i];

  lseek(fd->fd_sys2,req_off,SEEK_SET);
  write(fd->fd_sys2,(char*)buf+userbuf_off,req_len);
  **mwrite(fd->fd_sys3,req_off,(char*)buf+userbuf_off,req_len)**

  ADIOI_BUFFERED_WRITE_WITHOUT_READ
  off += flat_buf->blocklens[i];
 }

if(fwr_size)
 {
  req_off = off;
  req_len = fwr_size;
  userbuf_off = i;

  lseek(fd->fd_sys2,req_off,SEEK_SET)
  write(fd->fd_sys2,
   (char*)buf+userbuf_off,req_len);
  **mwrite(fd->fd_sys3,req_off,
   (char*)buf+userbuf_off,req_len);**

  ADIOI_BUFFERED_WRITE
 }

if(size)
 {
  req_off = off;
  req_len = size;
  userbuf_off = i;

  lseek(fd->fd_sys2,req_off,SEEK_SET);
  write(fd->fd_sys2,(char*)buf+userbuf_off,req_len);
  **mwrite(fd->fd_sys3,req_off,(char*)buf+userbuf_off,req_len);**
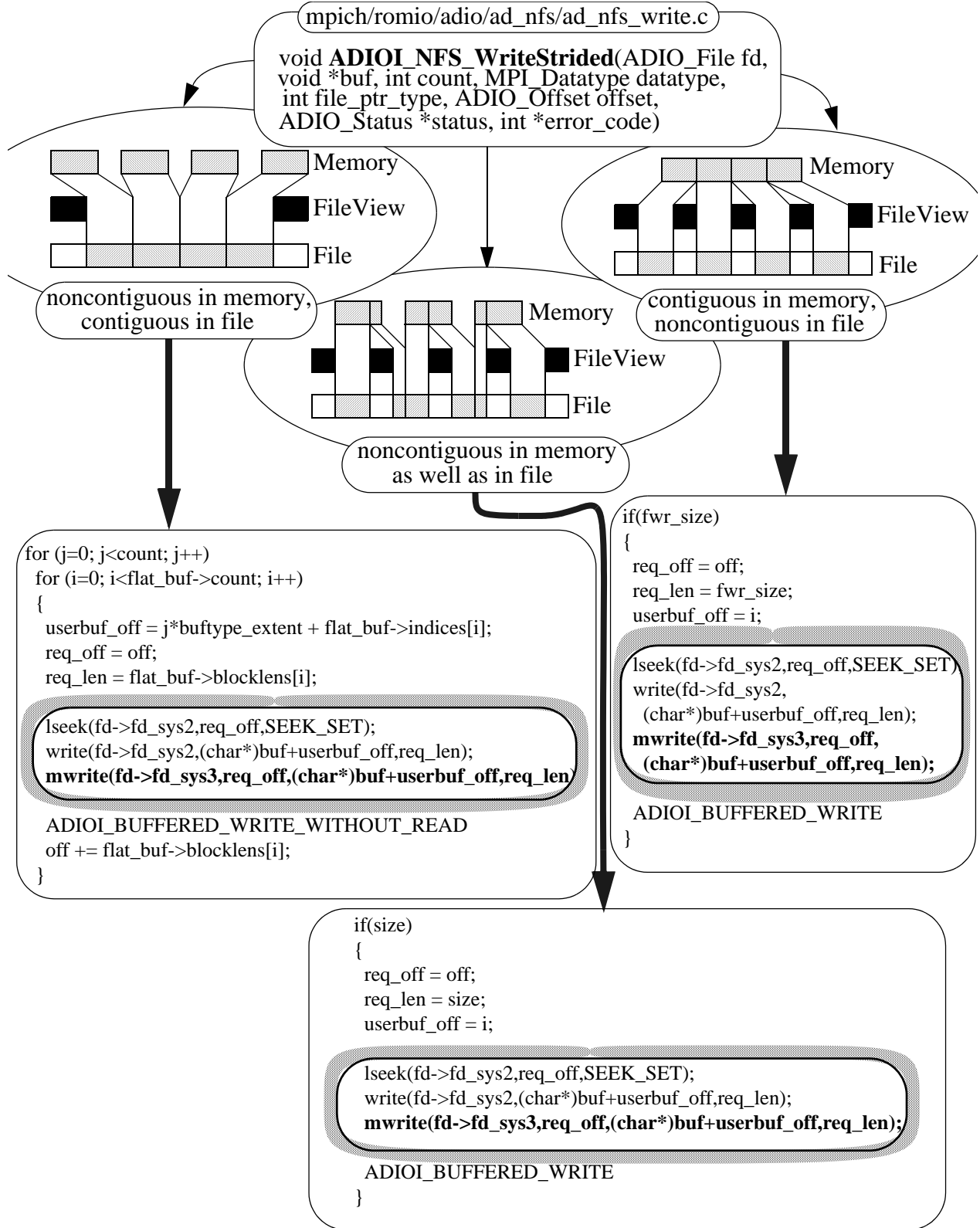
  ADIOI_BUFFERED_WRITE
 }

Fig. 32.

## Conclusion

- This is a cheap way of obtaining high performance I/O

Fig. 33.

Data in memory of 9 processes.
Each process keep only 1/9-th of matrix

Three files of three matrixes. Files are
accessible by all of processes, but for
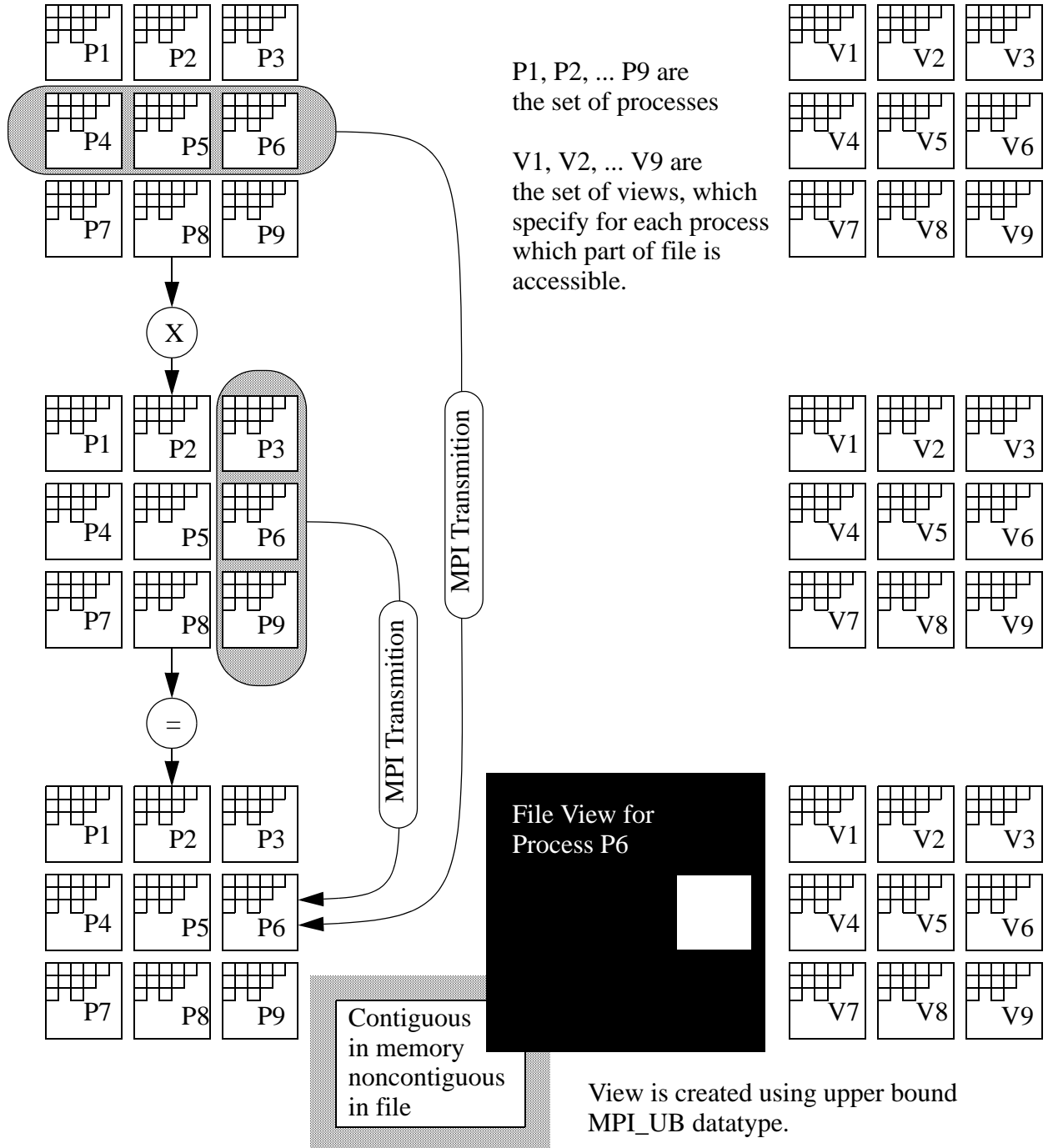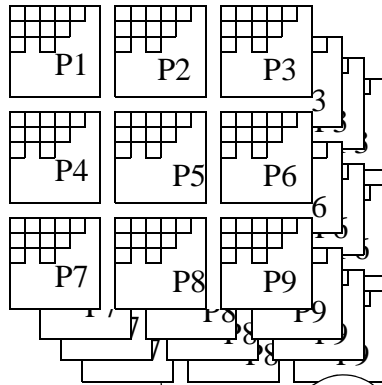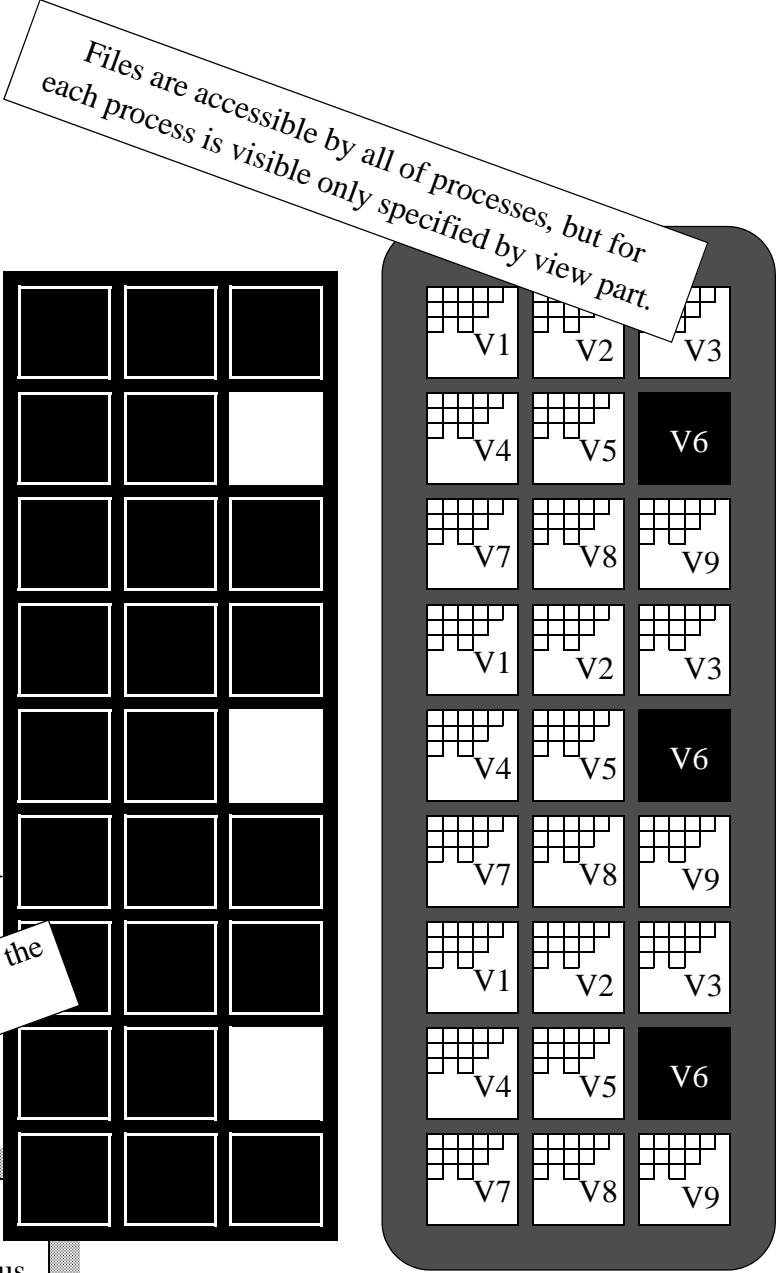each process is visible only spec. part.

P1, P2, ... P9 are
the set of processes

V1, V2, ... V9 are
the set of views, which
specify for each process
which part of file is
accessible.

P1 P2 P3
P4 P5 P6
P7 P8 P9

X

P1 P2 P3
P4 P5 P6
P7 P8 P9

=

P1 P2 P3
P4 P5 P6
P7 P8 P9

MPI Transmition

MPI Transmition

Contiguous
in memory
noncontiguous
in file

File View for
Process P6

View is created using upper bound
MPI_UB datatype.

V1 V2 V3
V4 V5 V6
V7 V8 V9

V1 V2 V3
V4 V5 V6
V7 V8 V9

V1 V2 V3
V4 V5 V6
V7 V8 V9

Fig. 34.

Data in memory of 9 processes.

Files are accessible by all of processes, but for each process is visible only specified by view part.

P1 P2 P3

P4 P5 P6

P7 P8 P9
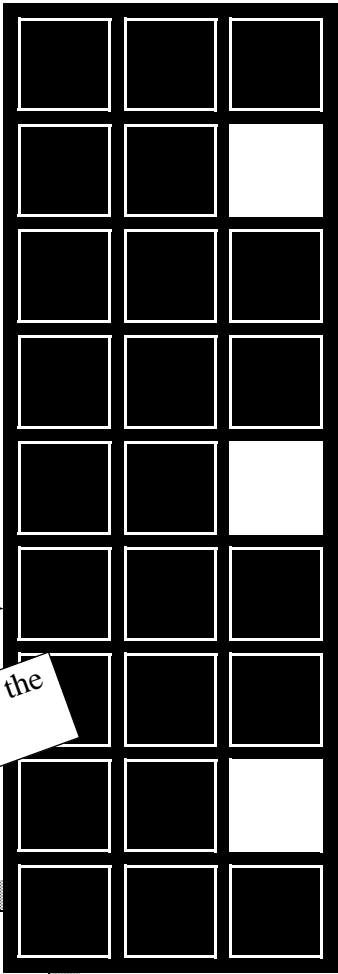
write uses nonzero offset, the position in the file relative to the current view

The view for the process P6

Contiguous in memory noncontiguous in file

P1, P2, ... P9 are the set of processes

V1, V2, ... V9 are the set of views

V1 V2 V3

V4 V5 V6

V7 V8 V9

V1 V2 V3

V4 V5 V6

V7 V8 V9

V1 V2 V3

V4 V5 V6

V7 V8 V9

View is created using upper bound and lower bound MPI_UB, MPI_LB datatype.