

Parallel I/O for SwissTx

**SFIO parallel I/O system implementation for
Swiss-Tx. MPI-II I/O interface support for
SFIO.**

LSP DI EPFL

1.0 Introduction

1.1 Requirements for parallel I/O

1.1.1 Scalable Throughput

1.1.2 Multiple Access

We try to design the system adapted to multiple access.

1.2 Striped Files Parallel Access Solution

Cyclical distribution of logical file block by block over the set of striped files located on different disks and different computers.

1.3 Interface

1.3.1 SFIO

1.3.2 MPI-II I/O

1.4 Network communication and Disk I/O access optimisation

This is data flow optimisation problems over the communication network and between memory and physical disk

1.4.1 Data fragmentation problem for noncollective operations

For noncollective operations, data fragmentation problem is resolved and maximally optimized as for the communication over the network as well as for writing and reading optimization.

1.4.2 Data fragmentation problem for collective operations

The optimization is realized for noncollective operations do his job here and resolves half of optimisation problems, but there are space of optimization specialy for the case of collective communications.

Fig. 01.

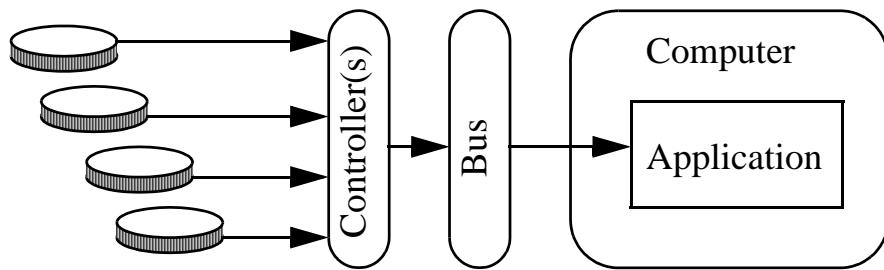
Parallel I/O for Swiss-Tx

- Requirements for Parallel I/O
- Striped Files Parallel Access Solution
- Interface
- Network Communication and Disk Access optimisation

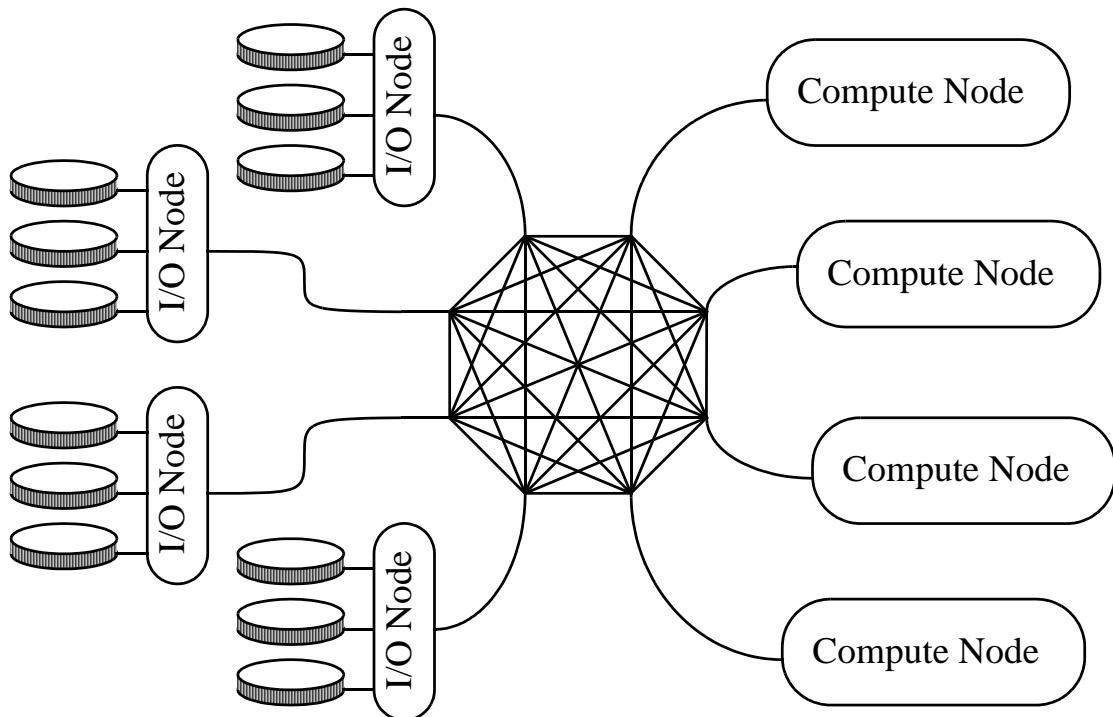
Fig. 02. We are trying to be close to the linear dependance of throughput from the

Requirements to Parallel I/O System

- Scalable Throughput



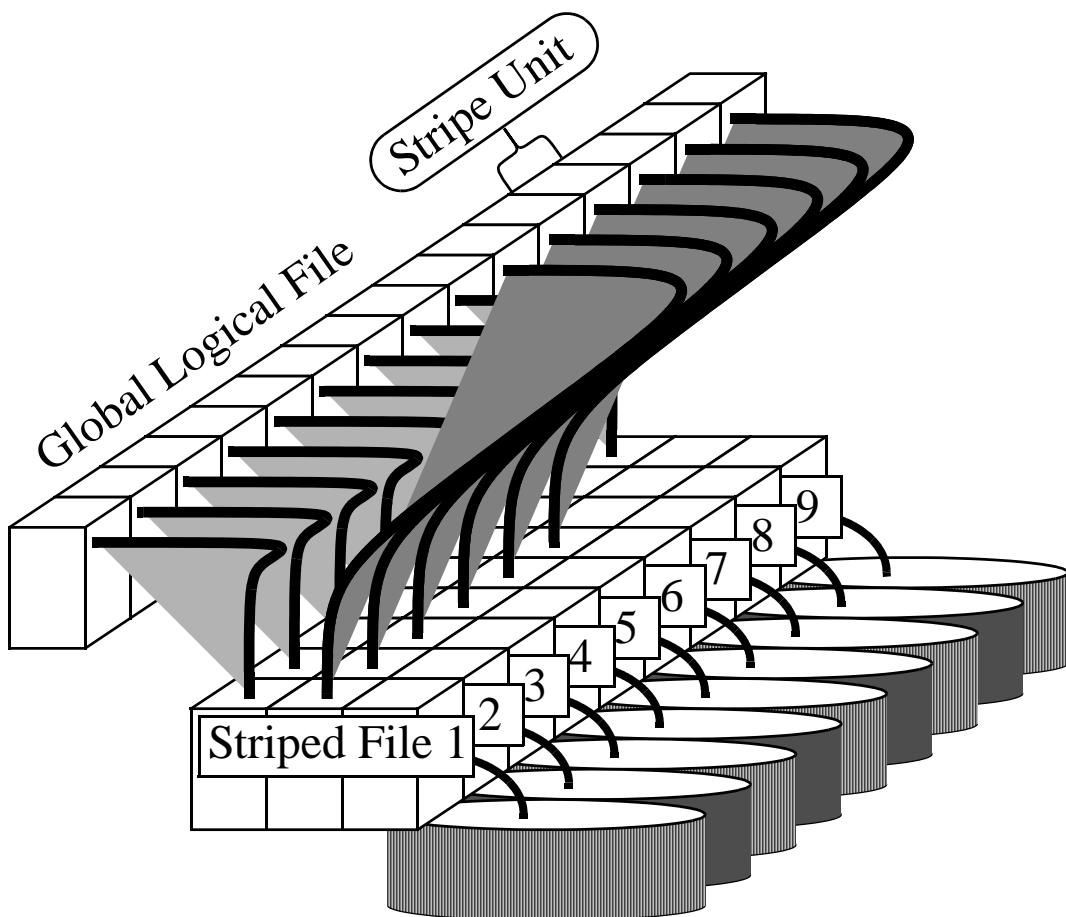
- Concurrent Access



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



across 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

| | |
|--------------------------------|------------------------------------|
| <code>MPI_File_open</code> | <code>MPI_File_write_all</code> |
| <code>MPI_File_set_view</code> | <code>MPI_File_read_all</code> |
| <code>MPI_File_write</code> | <code>MPI_File_write_at_all</code> |
| <code>MPI_File_read</code> | <code>MPI_File_read_at_all</code> |
| <code>MPI_File_write_at</code> | <code>MPI_File_close</code> |
| <code>MPI_File_read_at</code> | <code>MPI_File_delete</code> |

face.

Fog. 05.

Network and Disk Optimisation

- Network optimisation for noncollective access
- Disk optimisation for noncollective access
- Collective access optimisation

Pg. 06. Fragmented data optimisation for communication

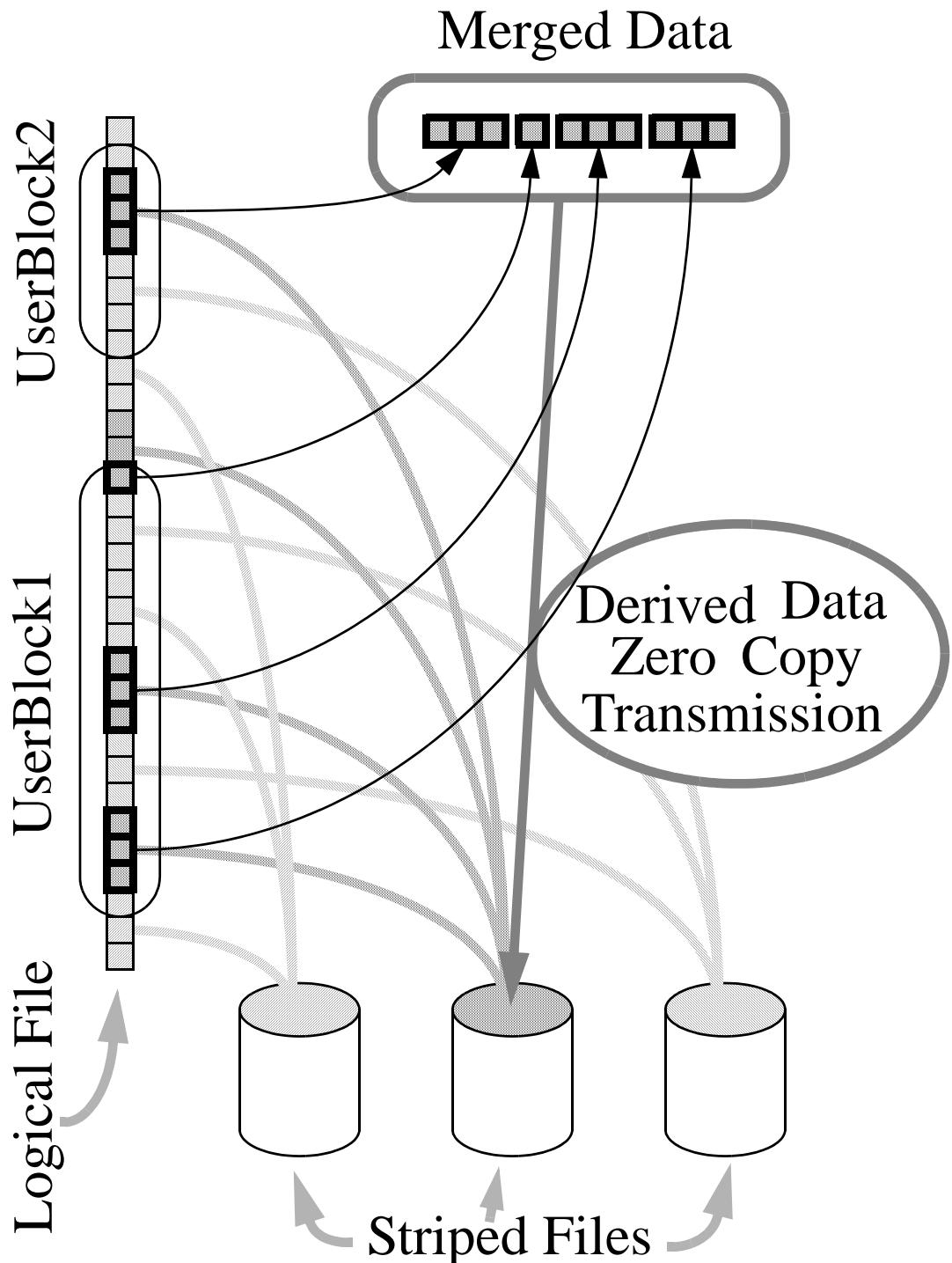


Fig. 07. Data fragmentation is resolved for network communication

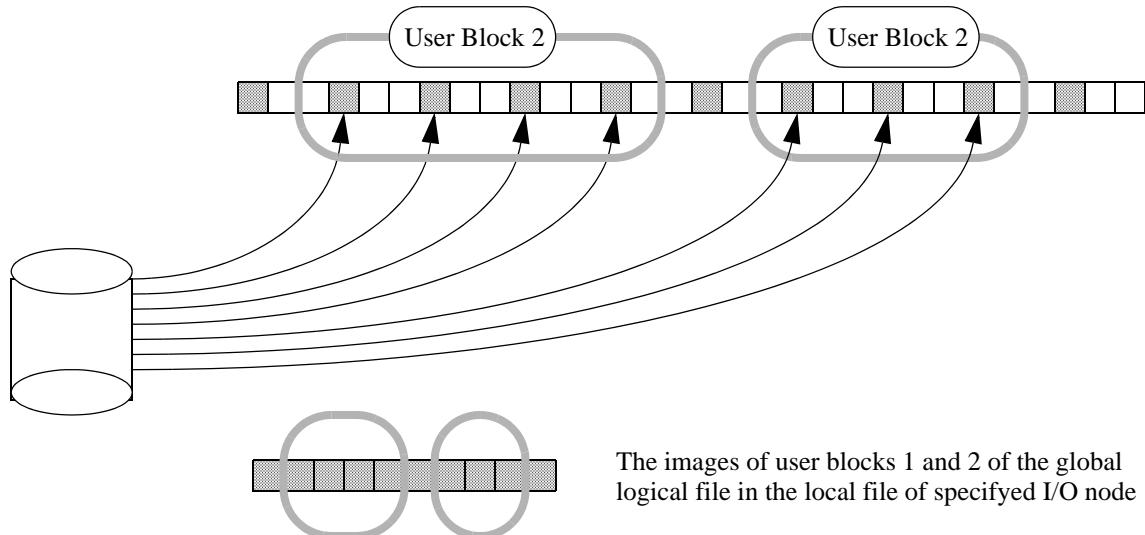


Fig. 08. Data fragmentation still exist at disk access

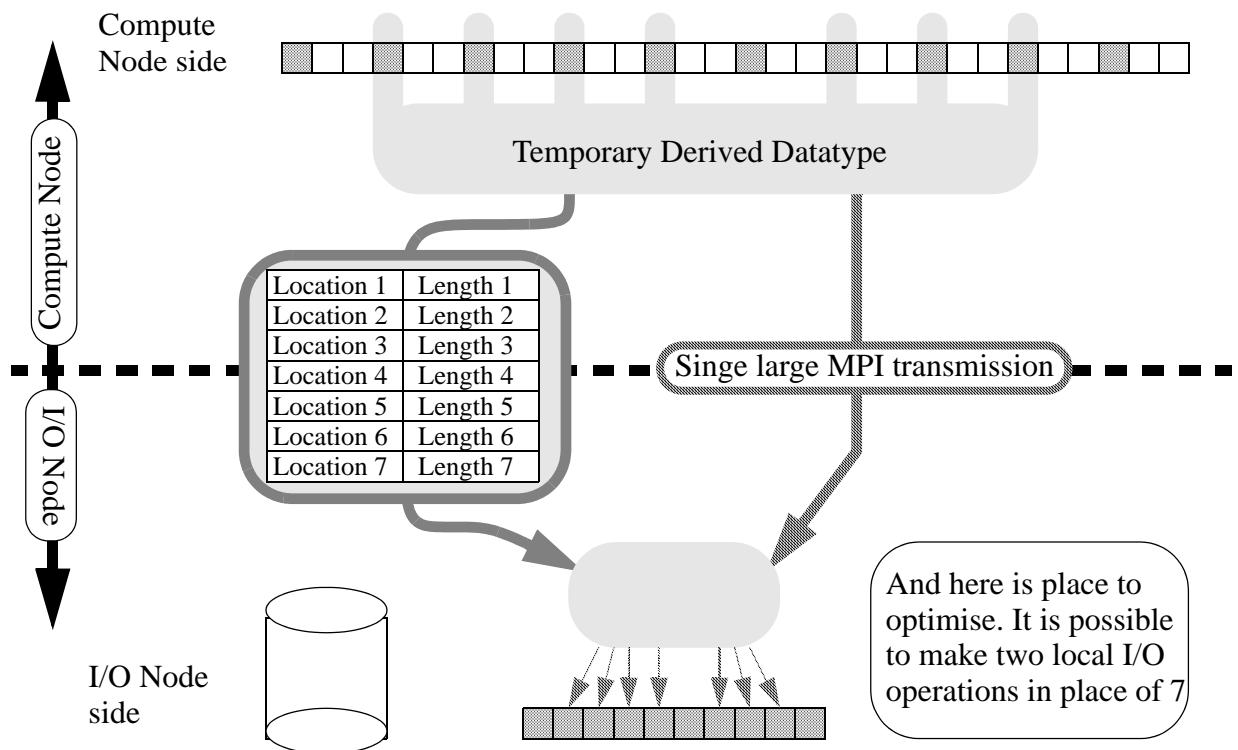


Fig. 09. Data fragmentation is resolved on disk level also

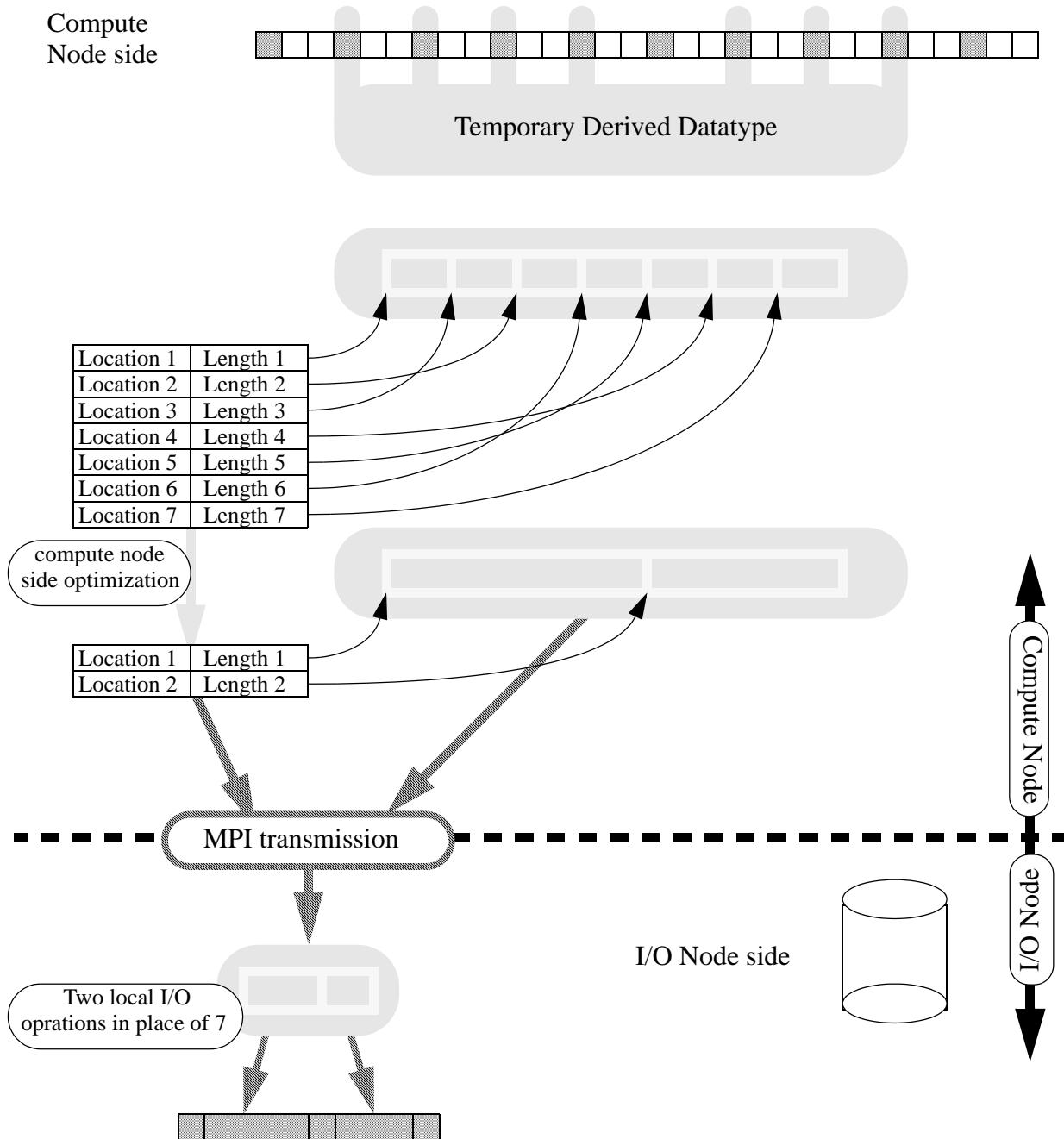


Fig. 10.

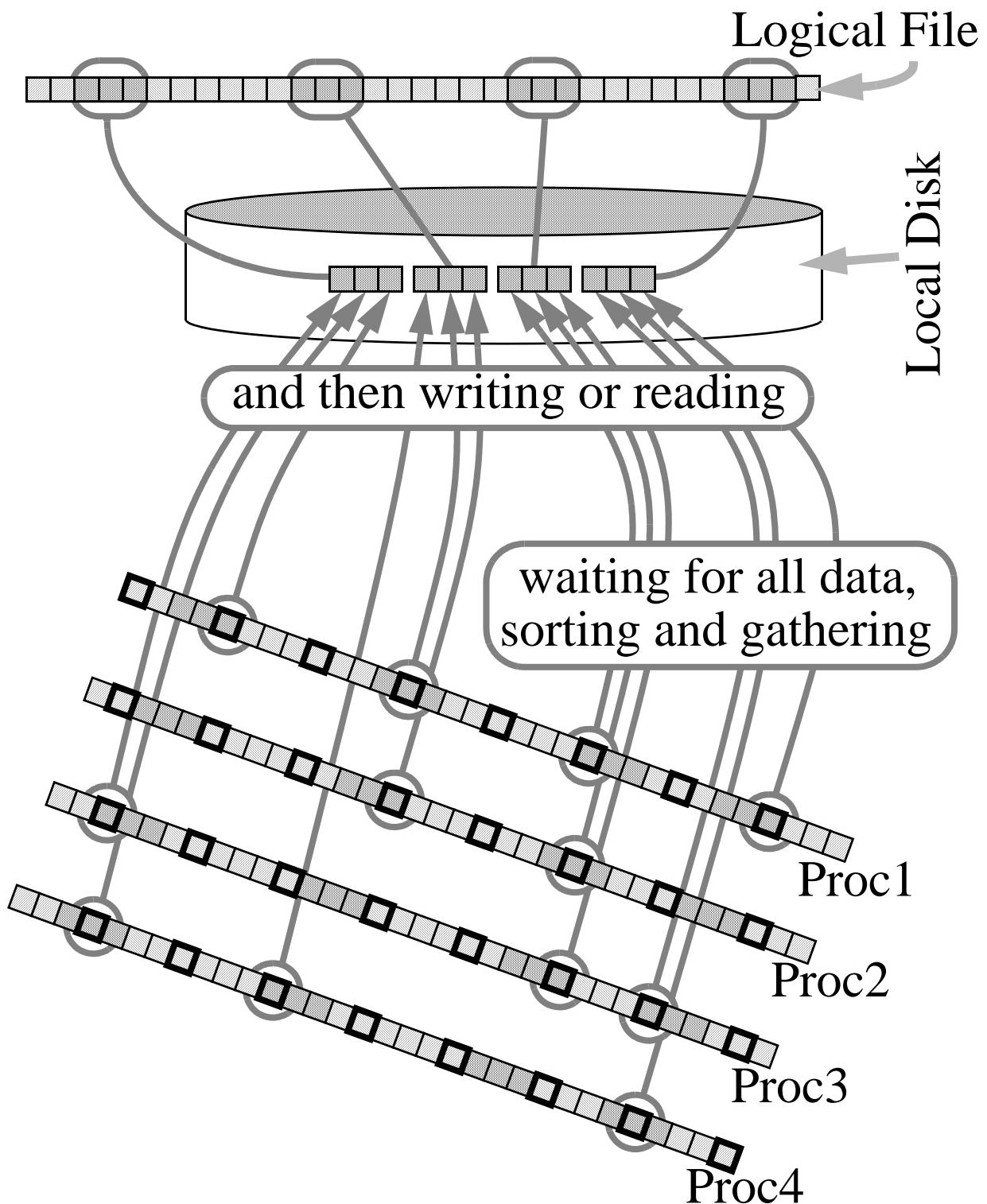
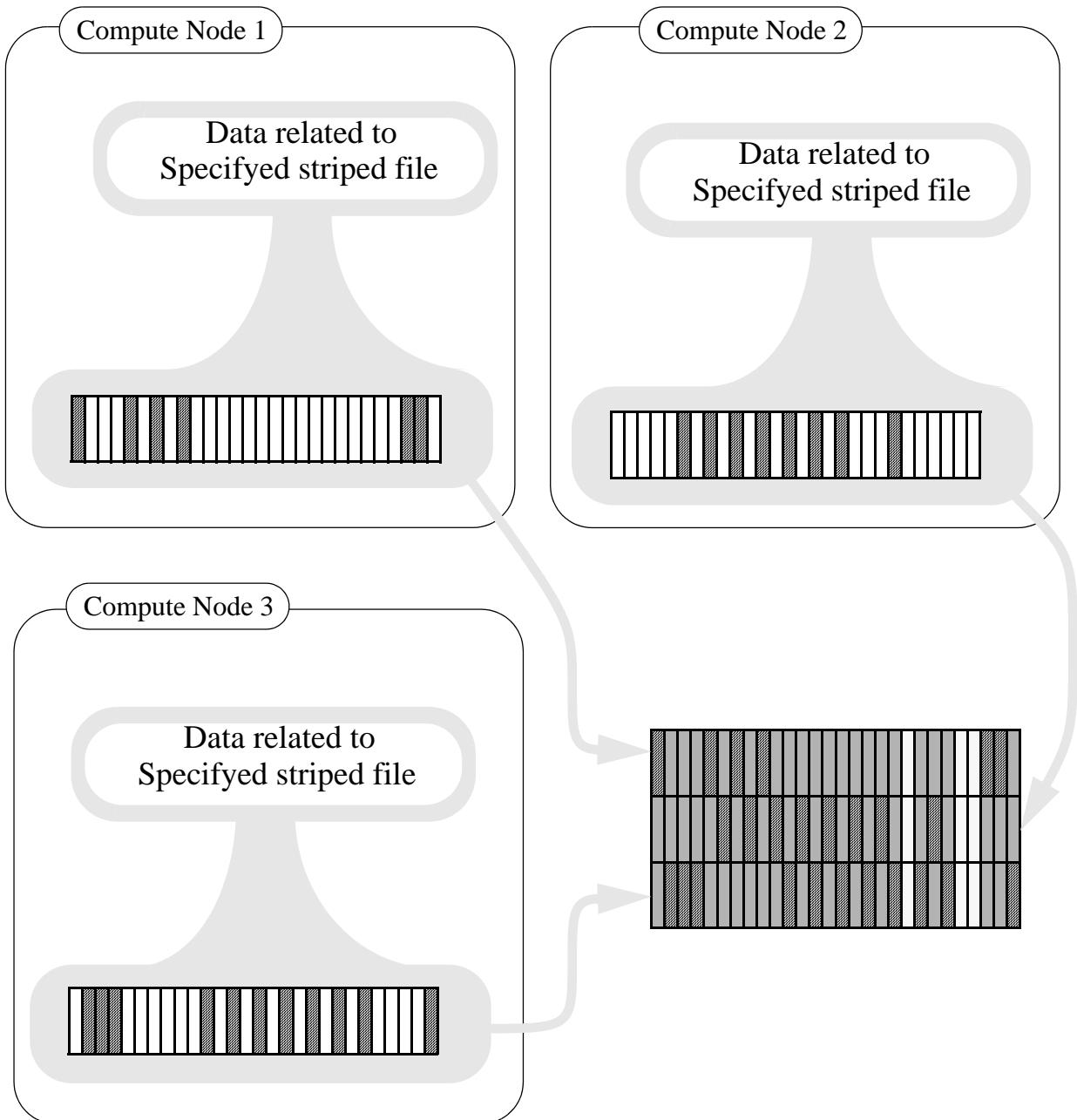


Fig. 11. Another example when disk access optimisation is required at collective



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

2.0 SFIO software architecture

2.1 Encapsulation over MPI

2.2 Functionality

2.3 Data organisation and optimisation

Fig. 12.

SFIO Software Architecture

- Encapsulation over MPI
- Functionality
- Data organisation and optimisation

Fig. 13. Encapsulation over MPI

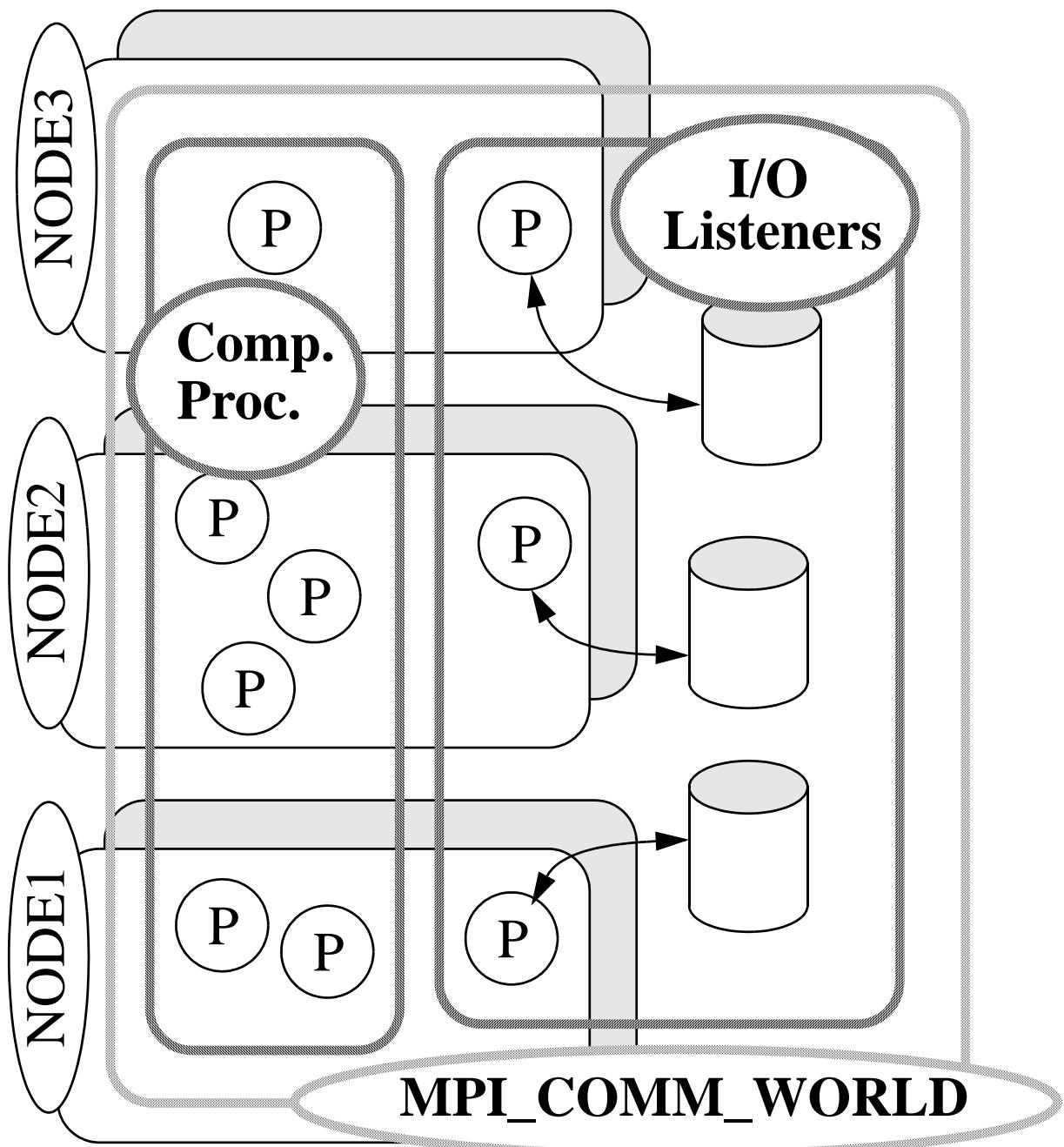
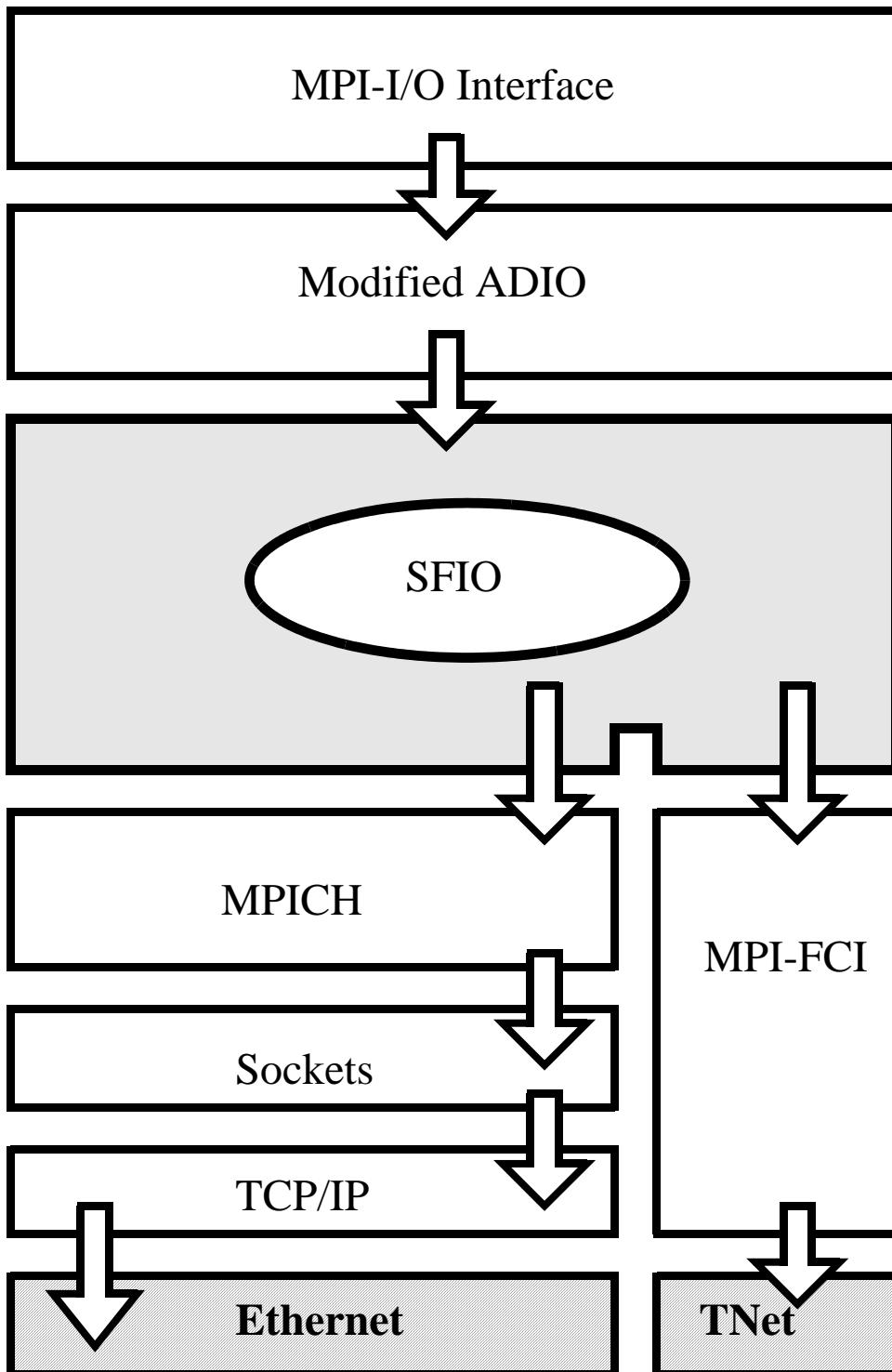
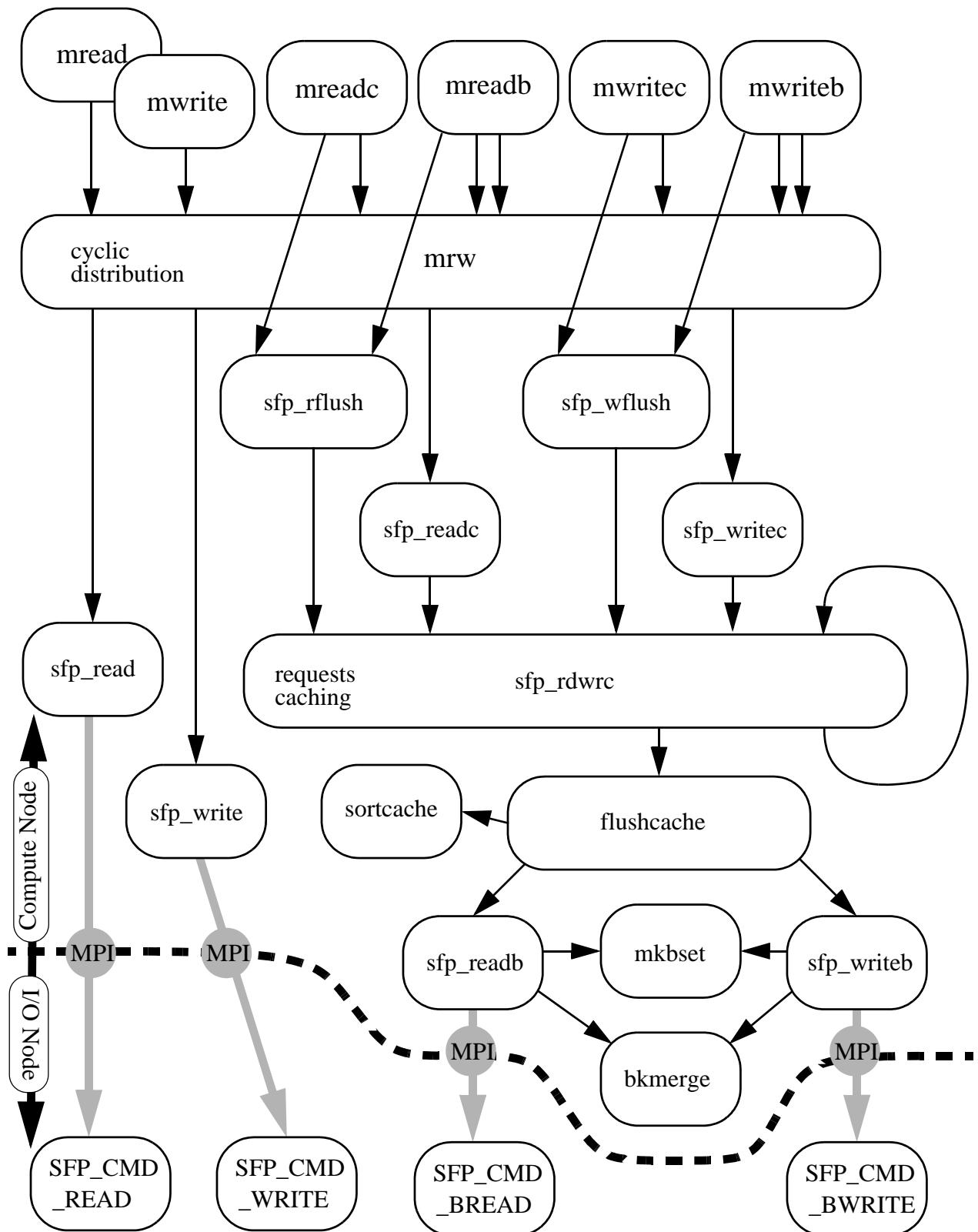


Fig. 14. Functionality. SFIO is implemented and tested with Digital Unix



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

Fig. 15. Functionality, data organisation and optimisation



3.0 SFIO interface

3.1 Single Block Access interface

3.2 Multiple Block Access interface

Fig. 16.

SFIO Interface

- Single Block Access Interface
- Single 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;
}
```

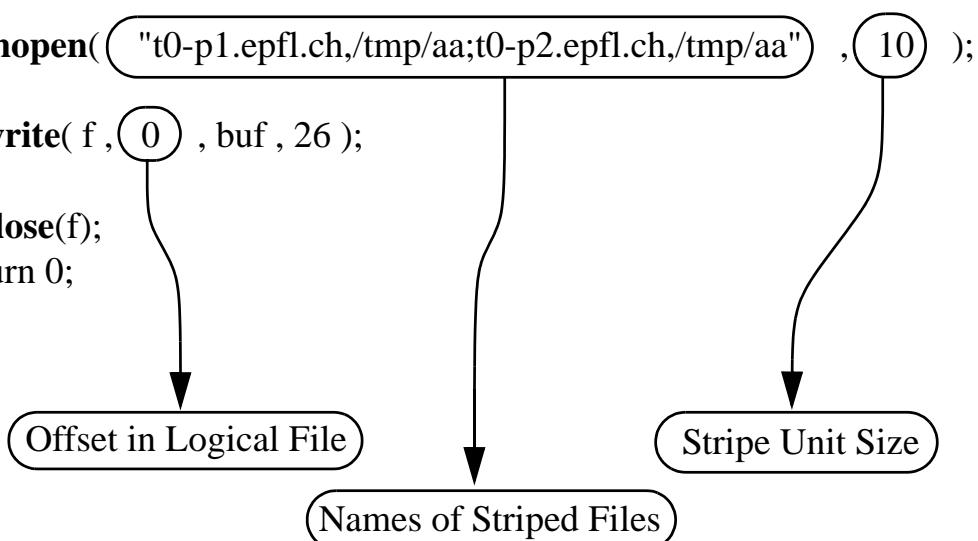


Fig. 18. If there are more than one compute nodes

Single Block Access Interface

```
#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;
}
```

The diagram illustrates the mapping of parameters from the provided C code to the SFIO interface. The parameters are:

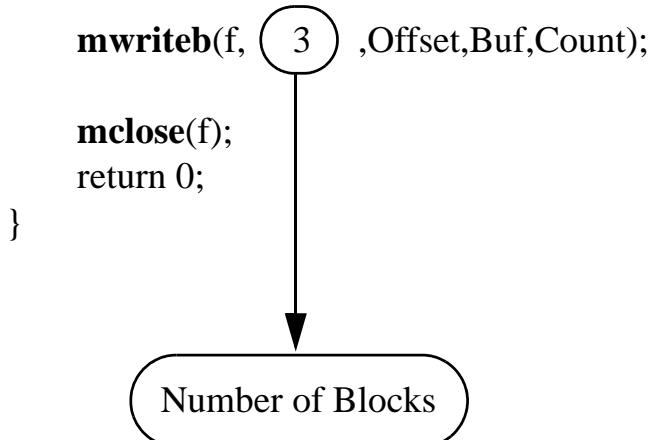
- Names of Striped Files**: Mapped to the file names in the `mopen` call: `"t0-p1.epfl.ch,/tmp/aa;t0-p2.epfl.ch,/tmp/aa"`.
- Stripe Unit Size**: Mapped to the `buf` parameter in the `mwrite` call.
- Offset in Logical File**: Mapped to the `rank*26` parameter in the `mwrite` call.
- Buffer**: Mapped to the `buf` parameter in the `mwrite` call.
- Buffer Size**: Mapped to the `26` parameter in the `mwrite` call.
- 10**: Mapped to the `10` parameter in the `mopen` call.

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);
```



4.0 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 concurrently 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. 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 concurrently reading/writing compute processes.

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

5.0 Integration into MPI-II I/O

Portable MPI Model Implementation of Argonne National Laboratory

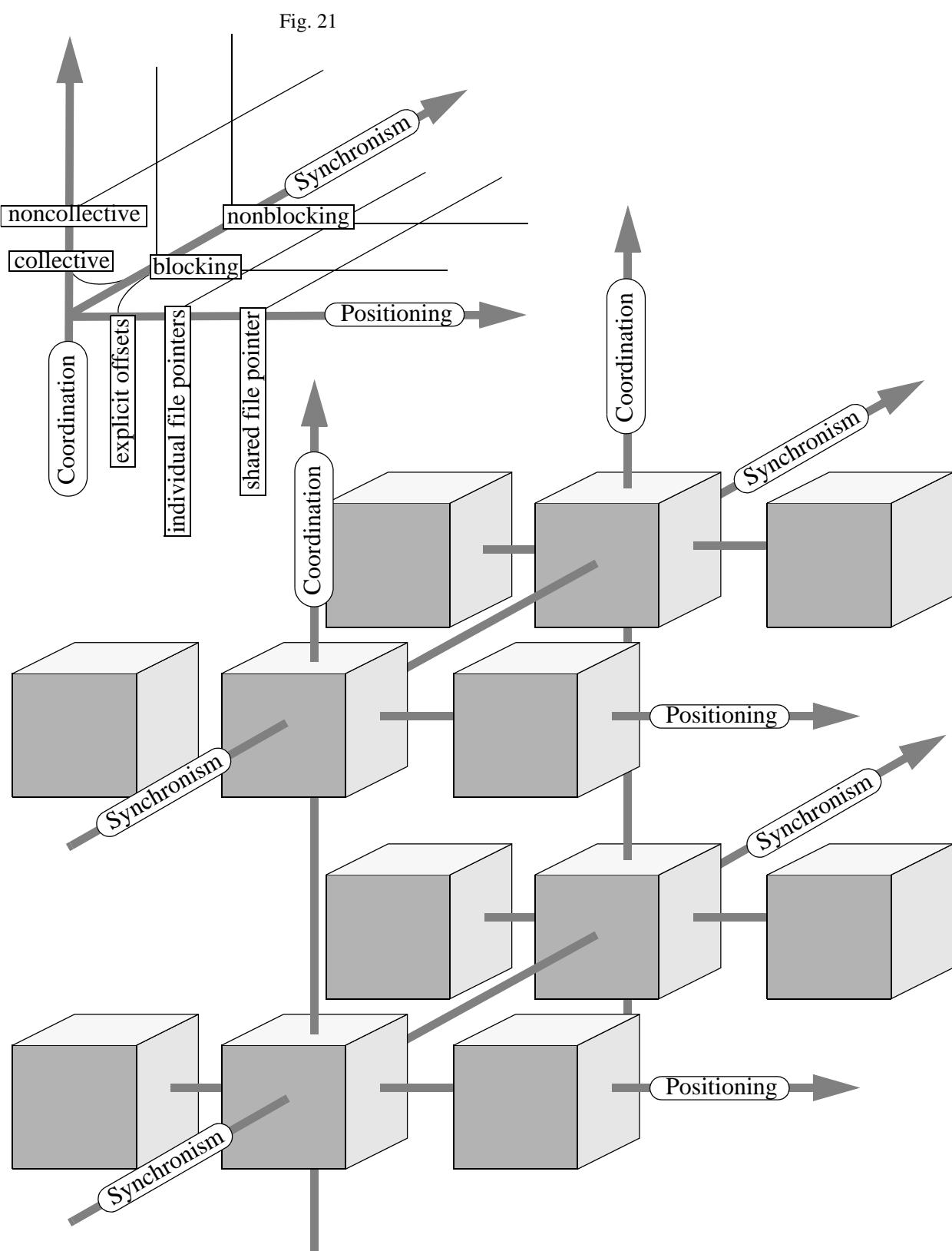
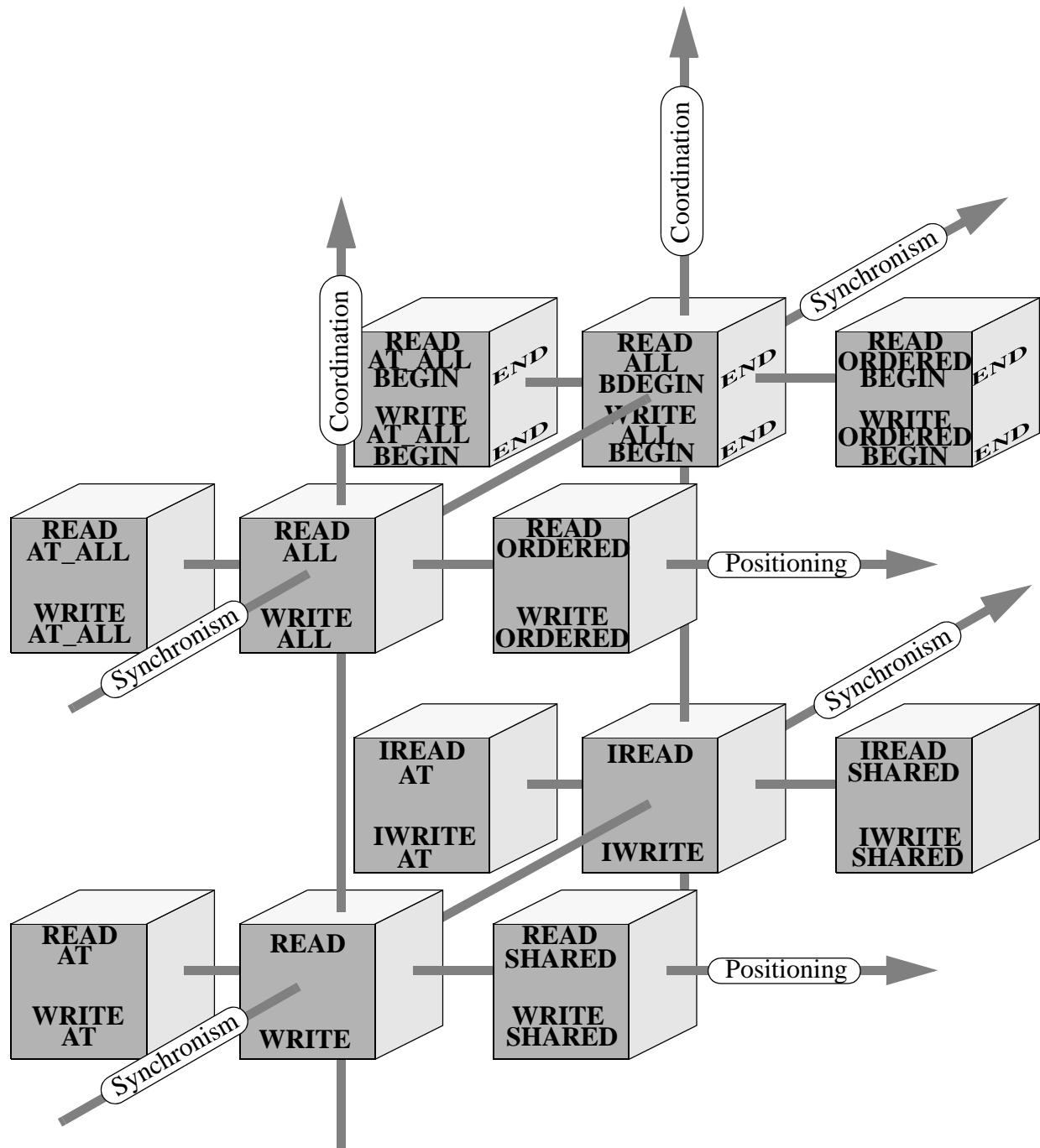


Fig. 22.



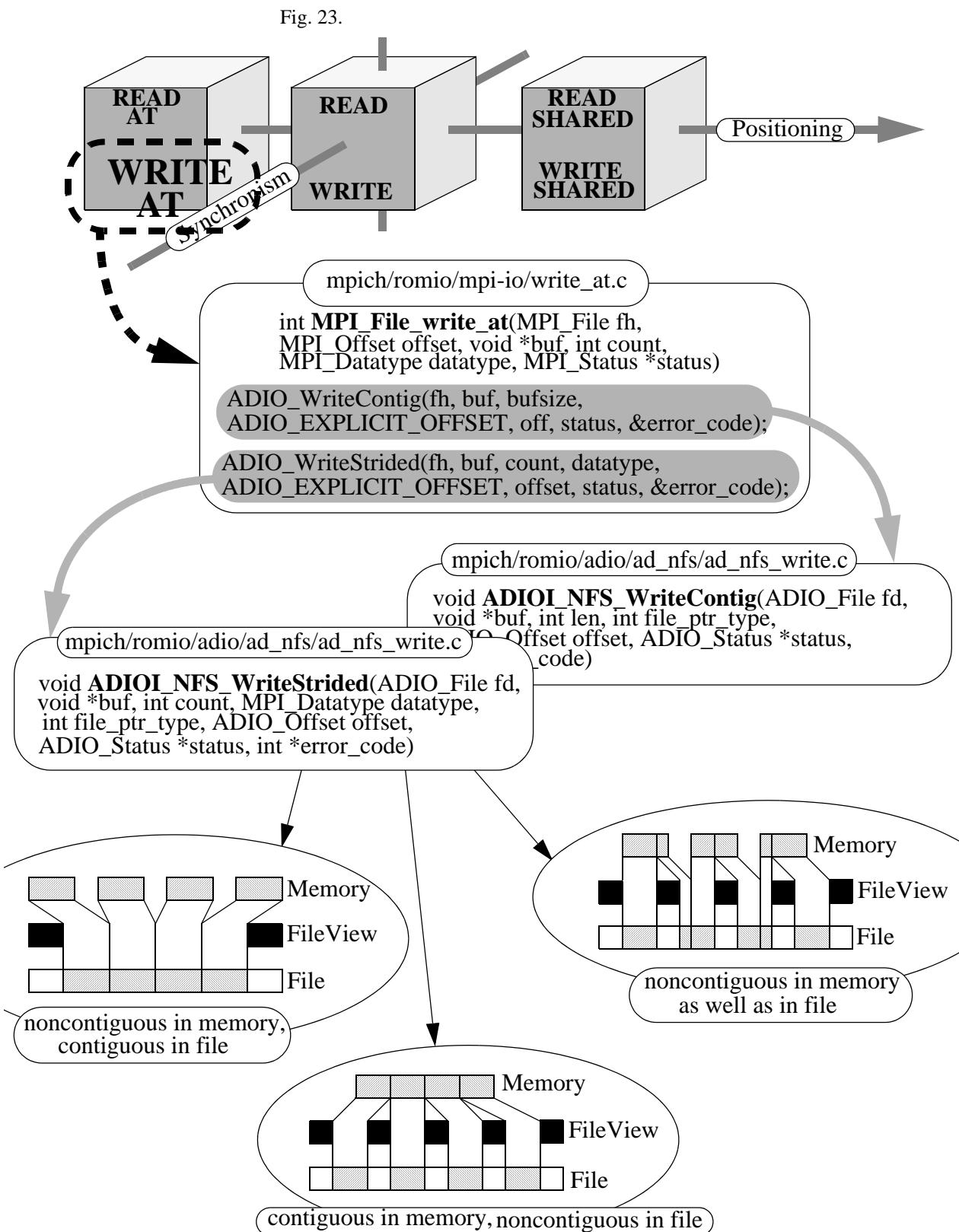


Fig. 24.

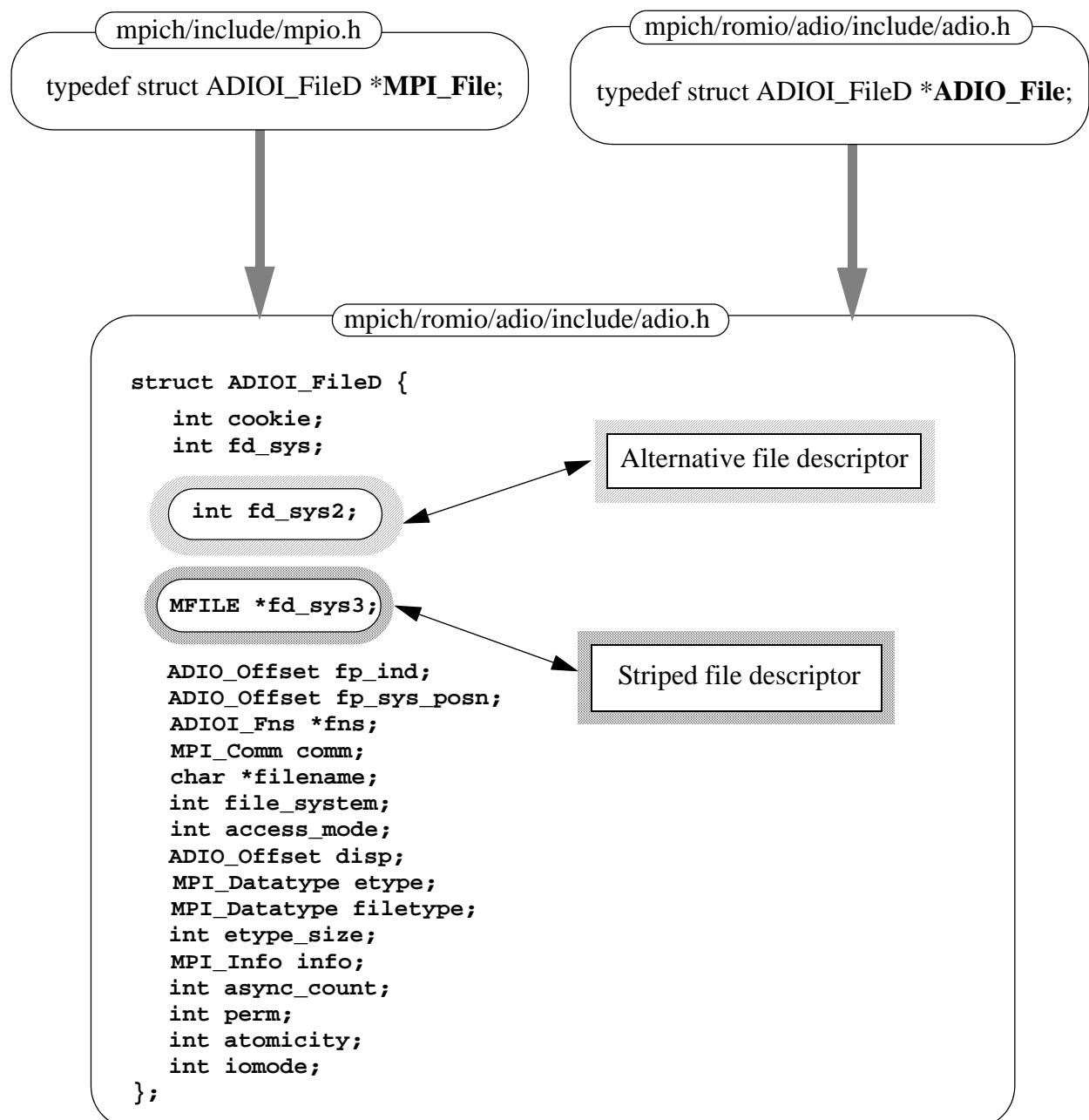


Fig. 25.

mpich/romio/adio/ad_nfs/ad_nfs_open.c

void ADIOI_NFS_Open
(ADIO_File fd, int *error_code)

fd->fd_sys=open(fd->filename,amode,perm);

filename2=mkfilename2(fd->filename);
fd->fd_sys2=open(filename2,amode,perm);
free(filename2);

filename3=mkfilename3(fd->filename);
fd->fd_sys3=mopen(filename3,64);
free(filename3);

char *mkfilename3(char *filename)

{

char *filename3;
const int devnum=8;
char *prefix[8]=
{
 "t0-p1.epfl.ch./scratch/p1/gabriely/sfp:",
 "t0-p2.epfl.ch./scratch/p2/gabriely/sfp:",
 "t0-p3.epfl.ch./scratch/p3/gabriely/sfp:",
 "t0-p4.epfl.ch./scratch/p4/gabriely/sfp:",
 "t0-p5.epfl.ch./scratch/p5/gabriely/sfp:",
 "t0-p6.epfl.ch./scratch/p6/gabriely/sfp:",
 "t0-p7.epfl.ch./scratch/p7/gabriely/sfp:",
 "t0-p8.epfl.ch./scratch/p8/gabriely/sfp:"
};

char *suffix[8]={ ":01;":":02;":":03;":":04;":
 ":05;":":06;":":07;":":08; };

char* name;

int i;

int len;

for(i=0,name=filename;filename[i]!=0;i++)

{

if(filename[i]=='/')

name=filename+i+1;

}

for(i=0,len=0;i<devnum;i++)

{

len+=strlen(prefix[i]);

len+=strlen(name);

len+=strlen(suffix[i]);

}

filename3=(char*)malloc(len+1);

for(i=0,filename3[0]=0;i<devnum;i++)

{

strcat(filename3,prefix[i]);

strcat(filename3,name);

strcat(filename3,suffix[i]);

}

return filename3;

}

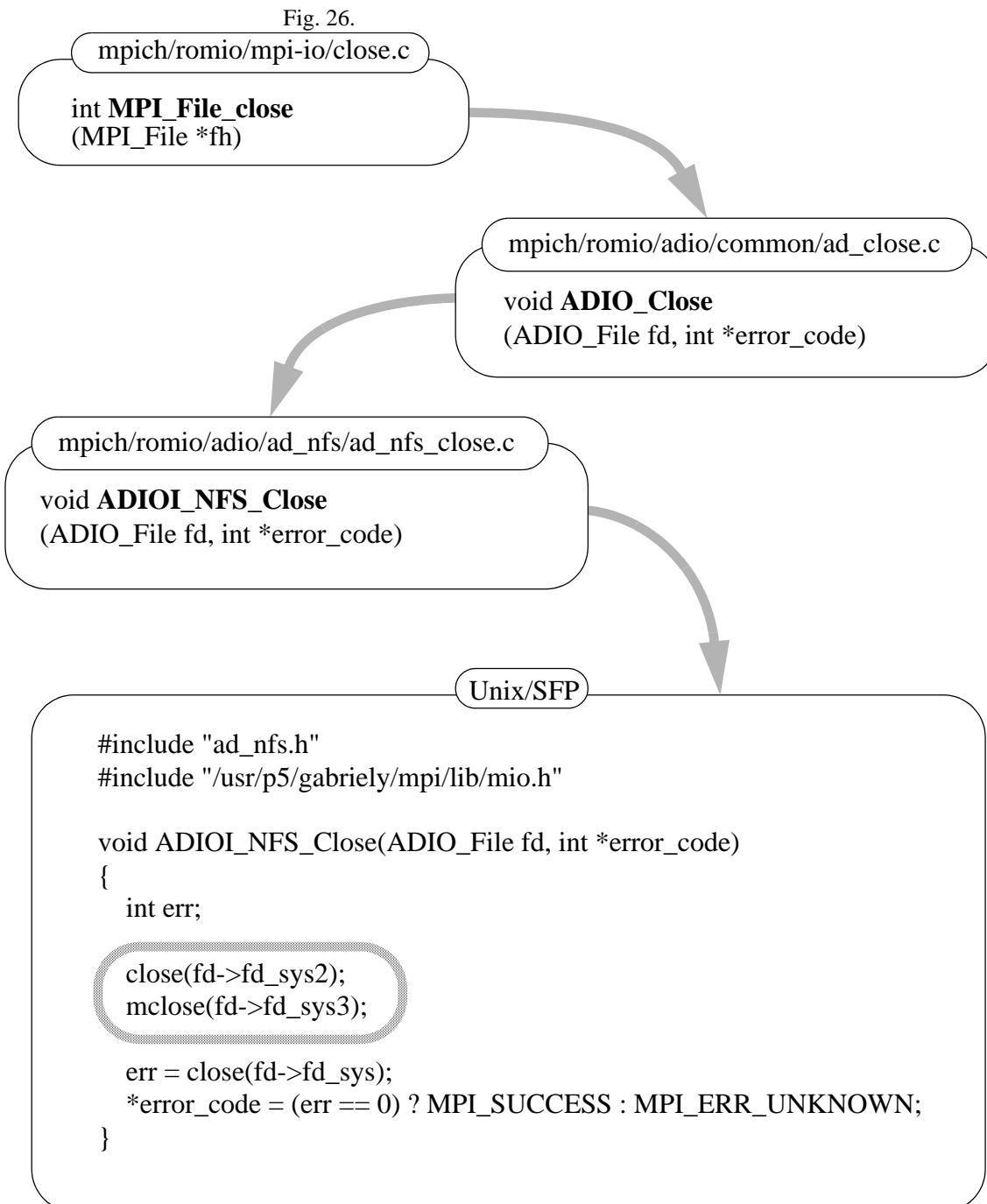
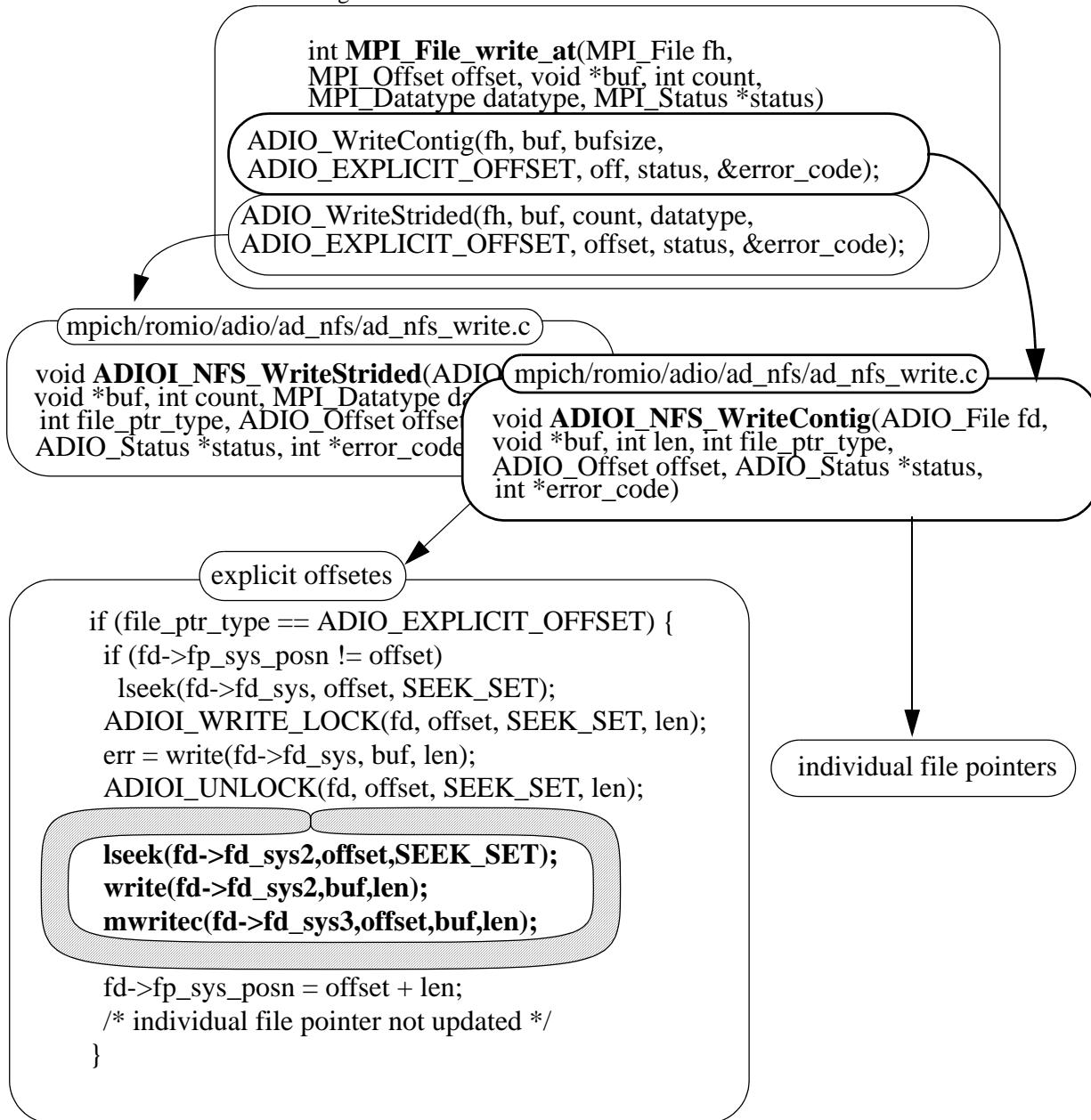


Fig. 27.



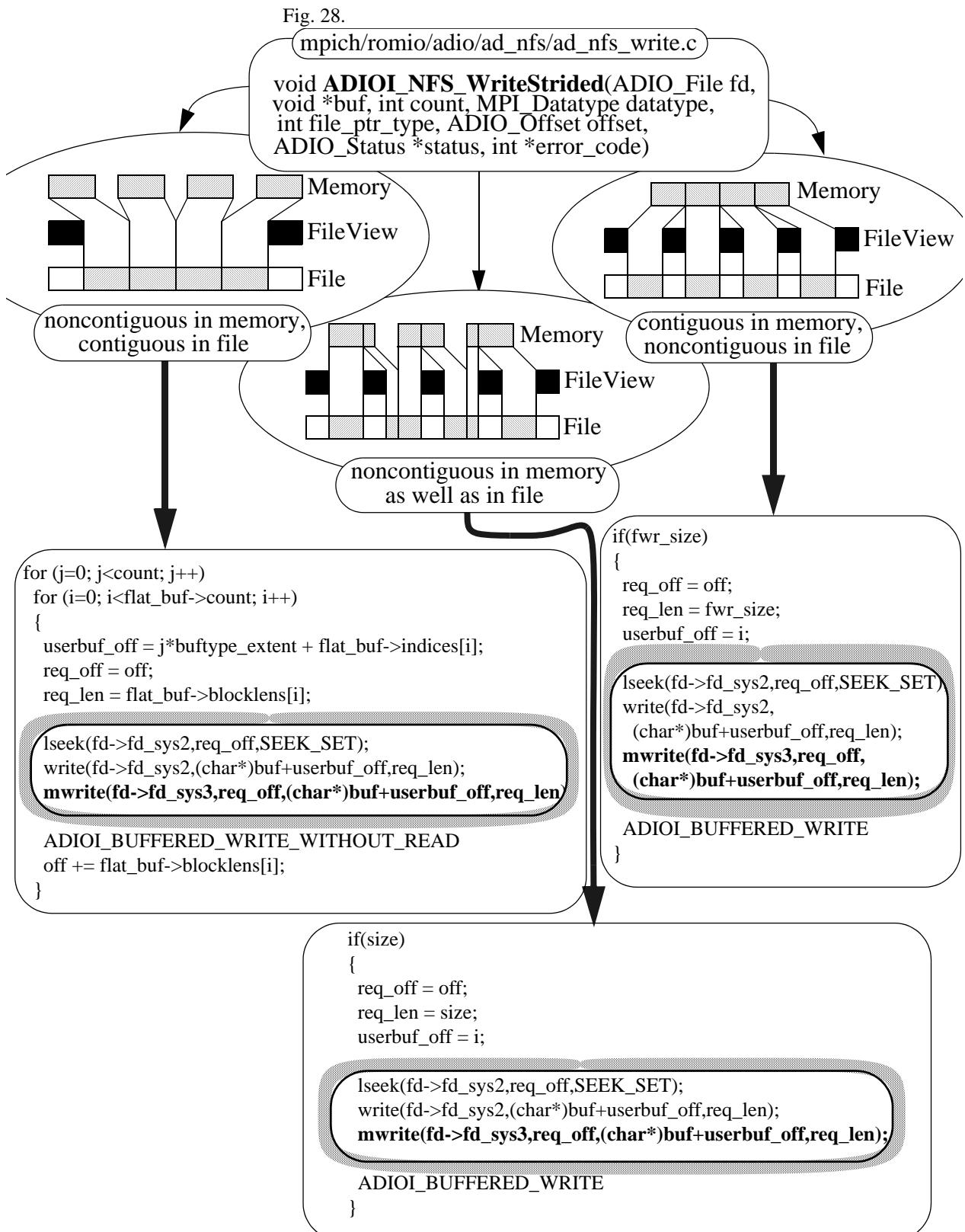


Fig. 29.

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.

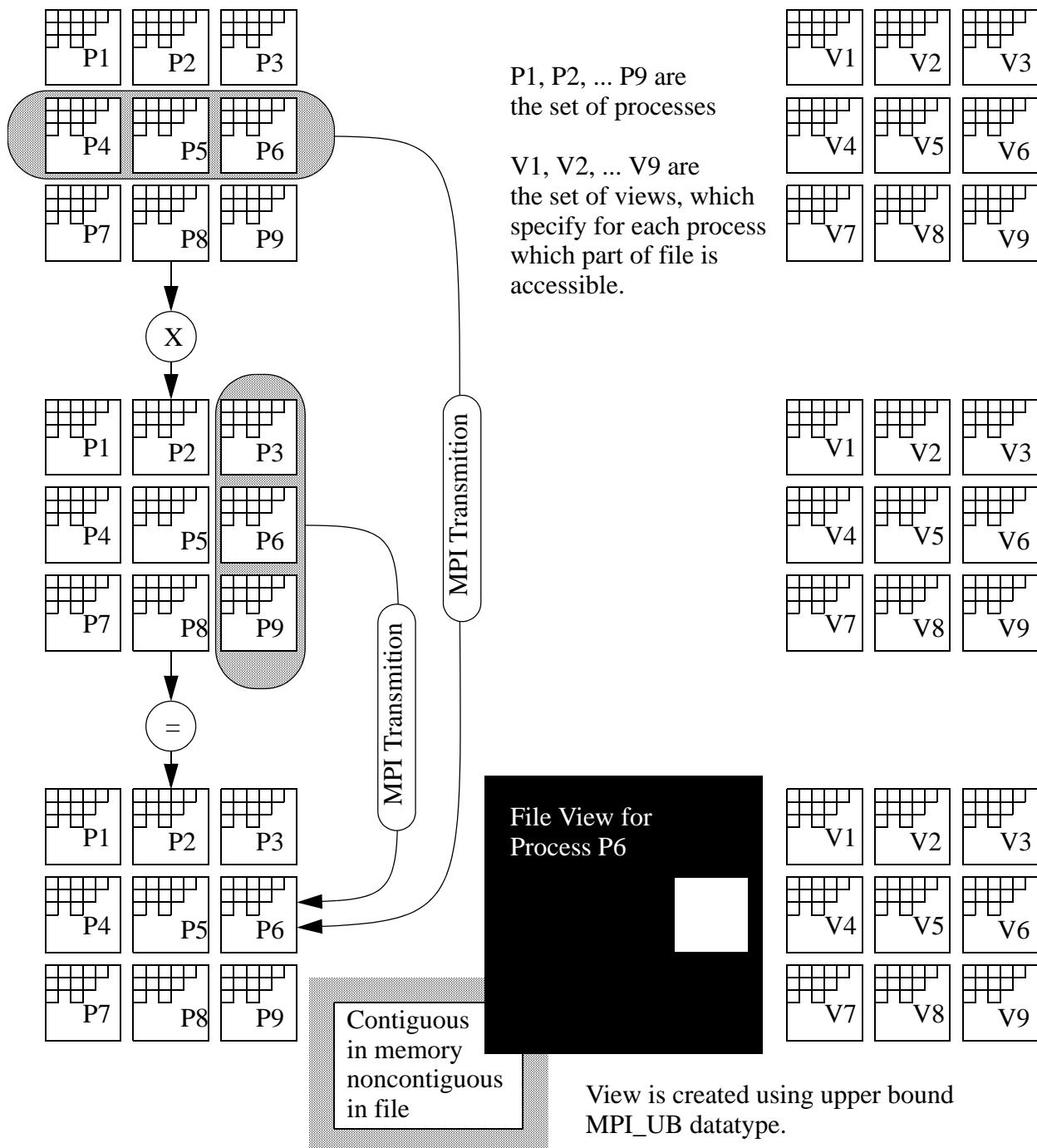
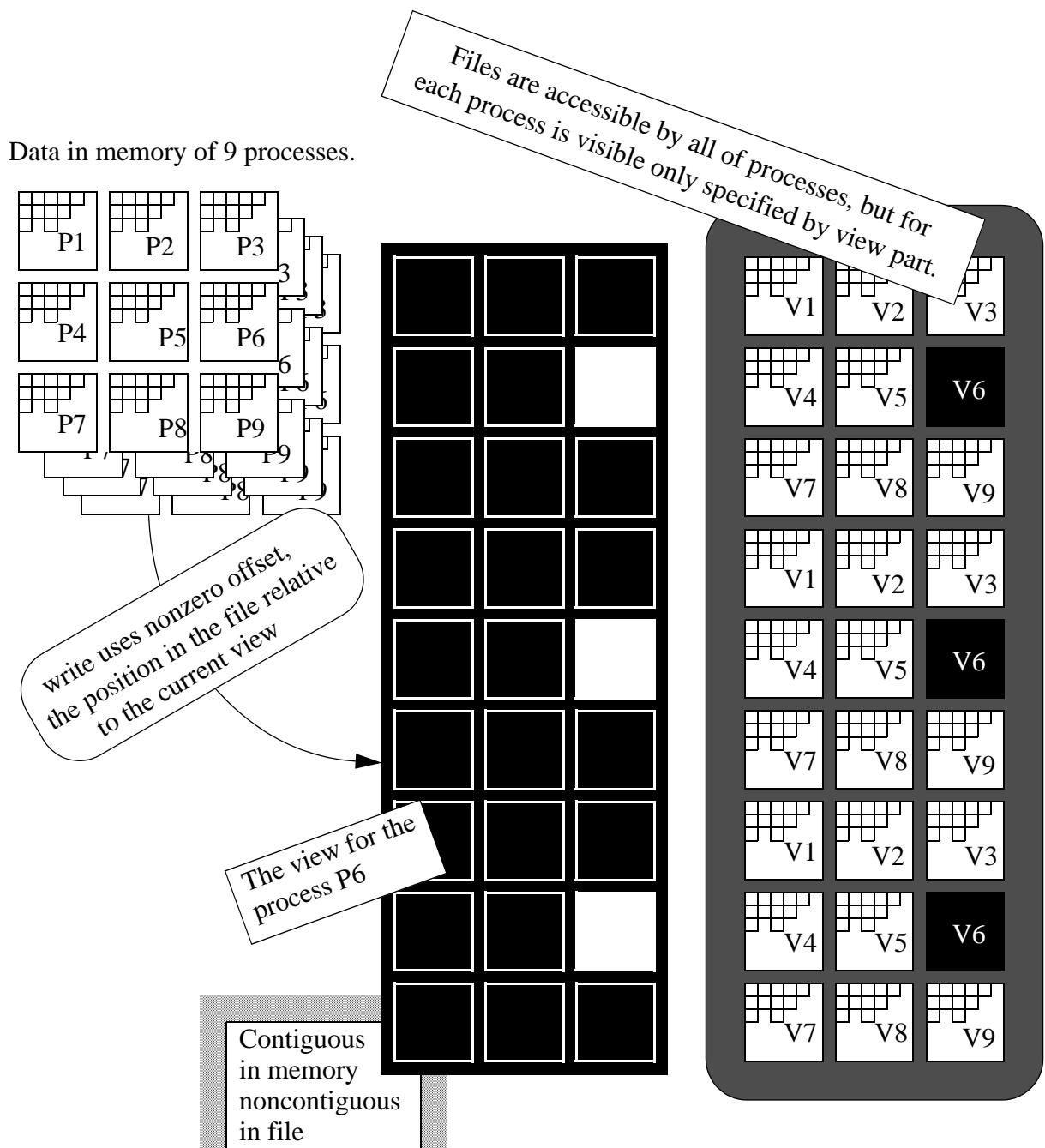


Fig. 30.



6.0 Conclusion

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