Fig. 00.

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

● Introduction
- ● Requirements for Parallel I/O
- ● Striped Files Parallel Access Solution

● Optimisations
- ● Network Optimisation
- ● Noncollective Access Disk Optimisation
- ● Collective Access Disk Optimisation

● Architecture
- ● Encapsulation over MPI
- ● Functionality
- ● Interface

● Performance
- ● Measurements on Windows NT cluster
- ● Measurements on Swiss-Tx TNET
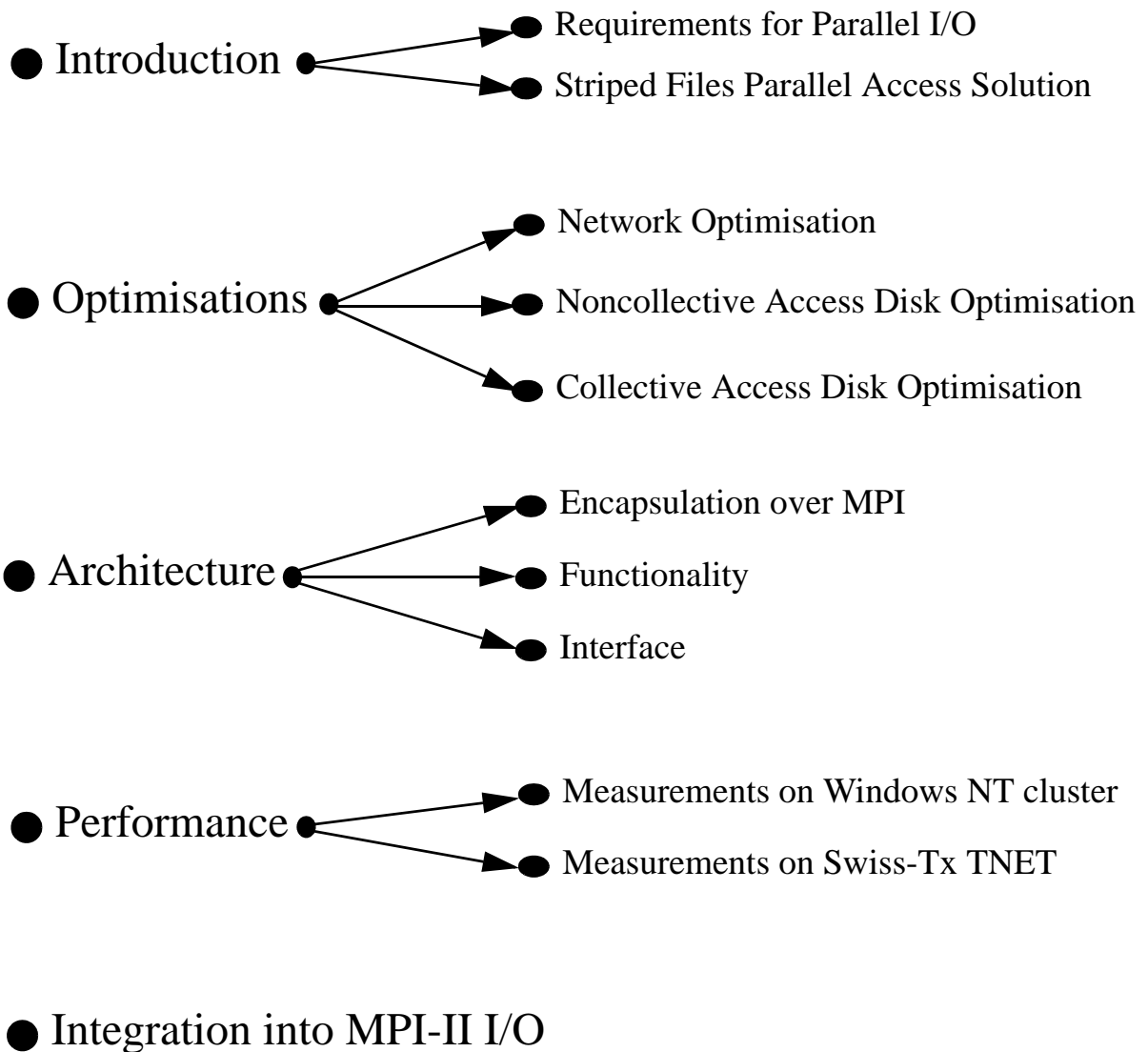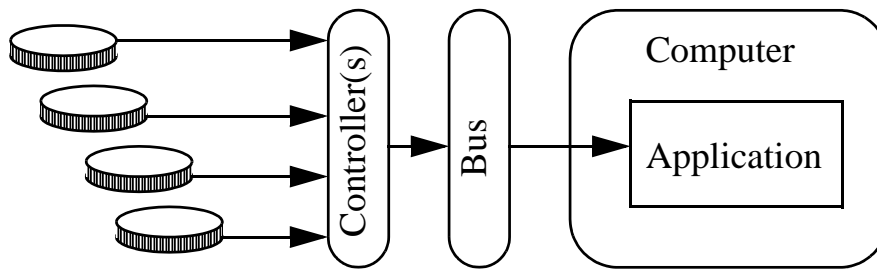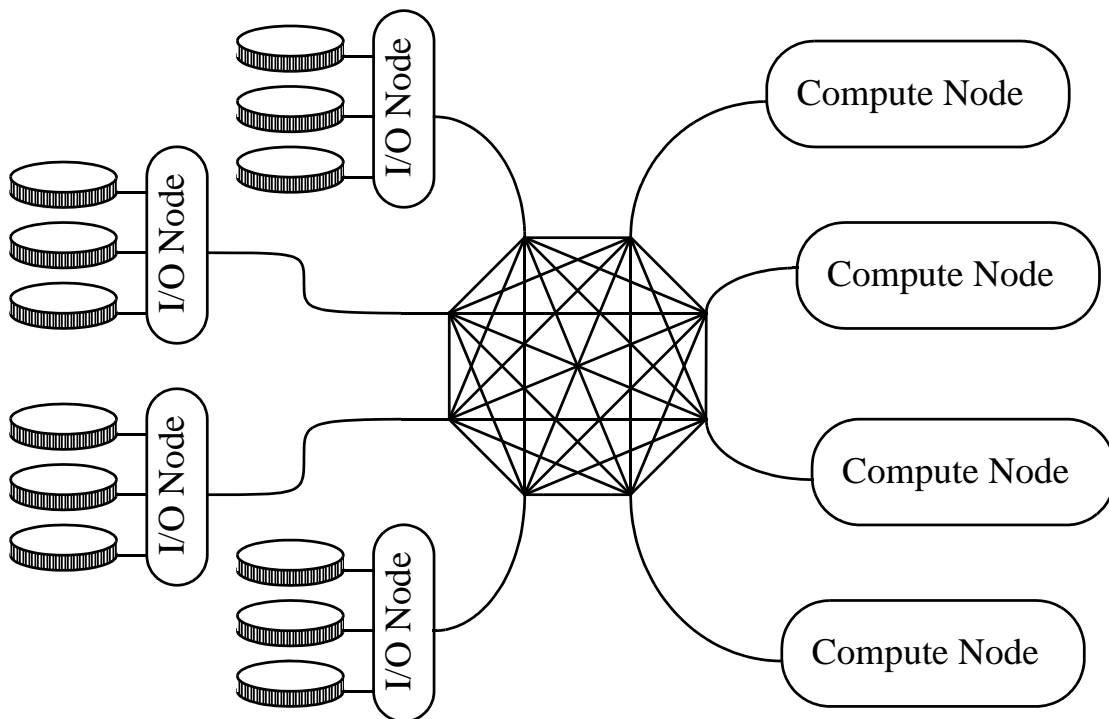
● Integration into MPI-II I/O

Fig. 01.

# Requirements for Parallel I/O

● Scalable Throughput



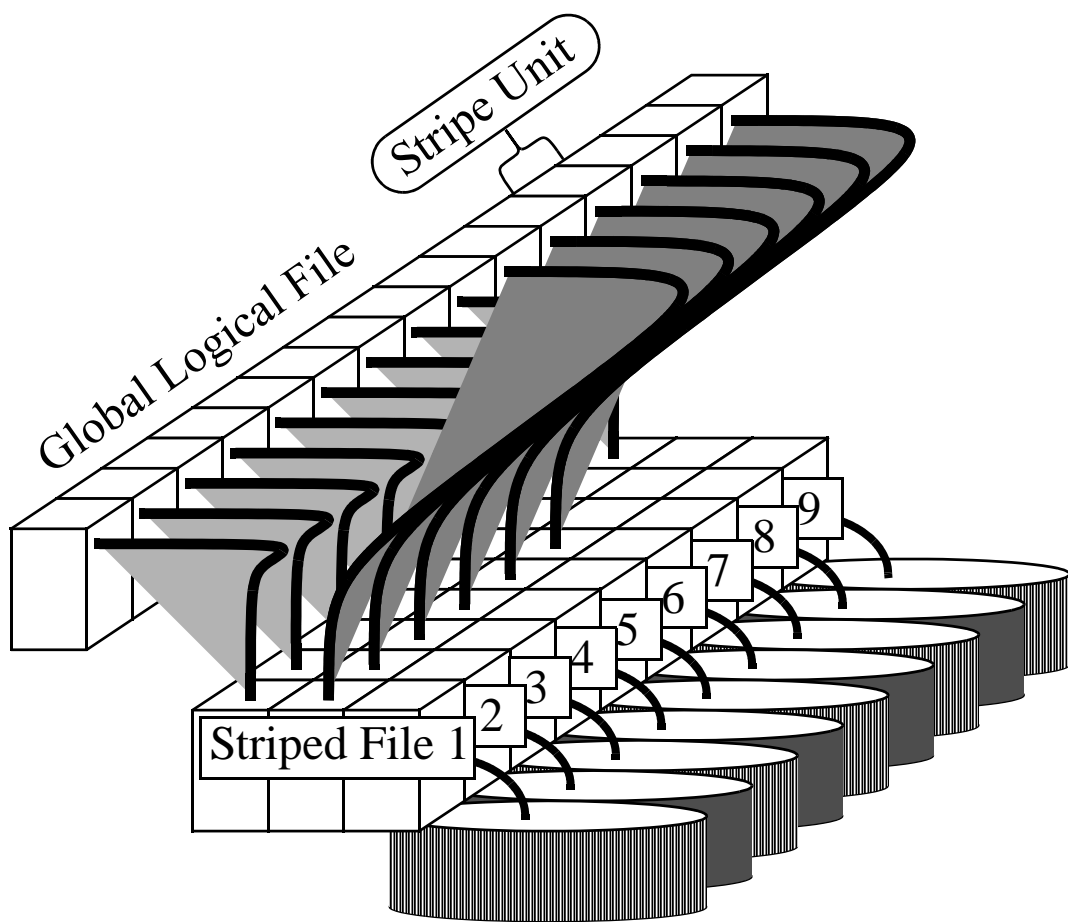● Concurent Access



The system should be scalable. System should be adapted to concurent access

Fig. 02.

# Striped Files Parallel Access Solution



Stripe Unit

Global Logical File

Striped File 1   2   3   4   5   6   7   8   9

The basic idea of our implementation is cyclical distribution of logical file accross
the set of striped files. Access to single part of logical file require number of
accesses to set of striped files, here is the space to paralelise access.

# Network Transmission Optimisation

Merged Data



UserBlock2

UserBlock1

Logical File

Compute Node

I/O Node

Derived Data
Zero Copy
Transmission

Striped Files

Fig. 04.

# 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

Two local I/O
oprations in place of 7

I/O Node side
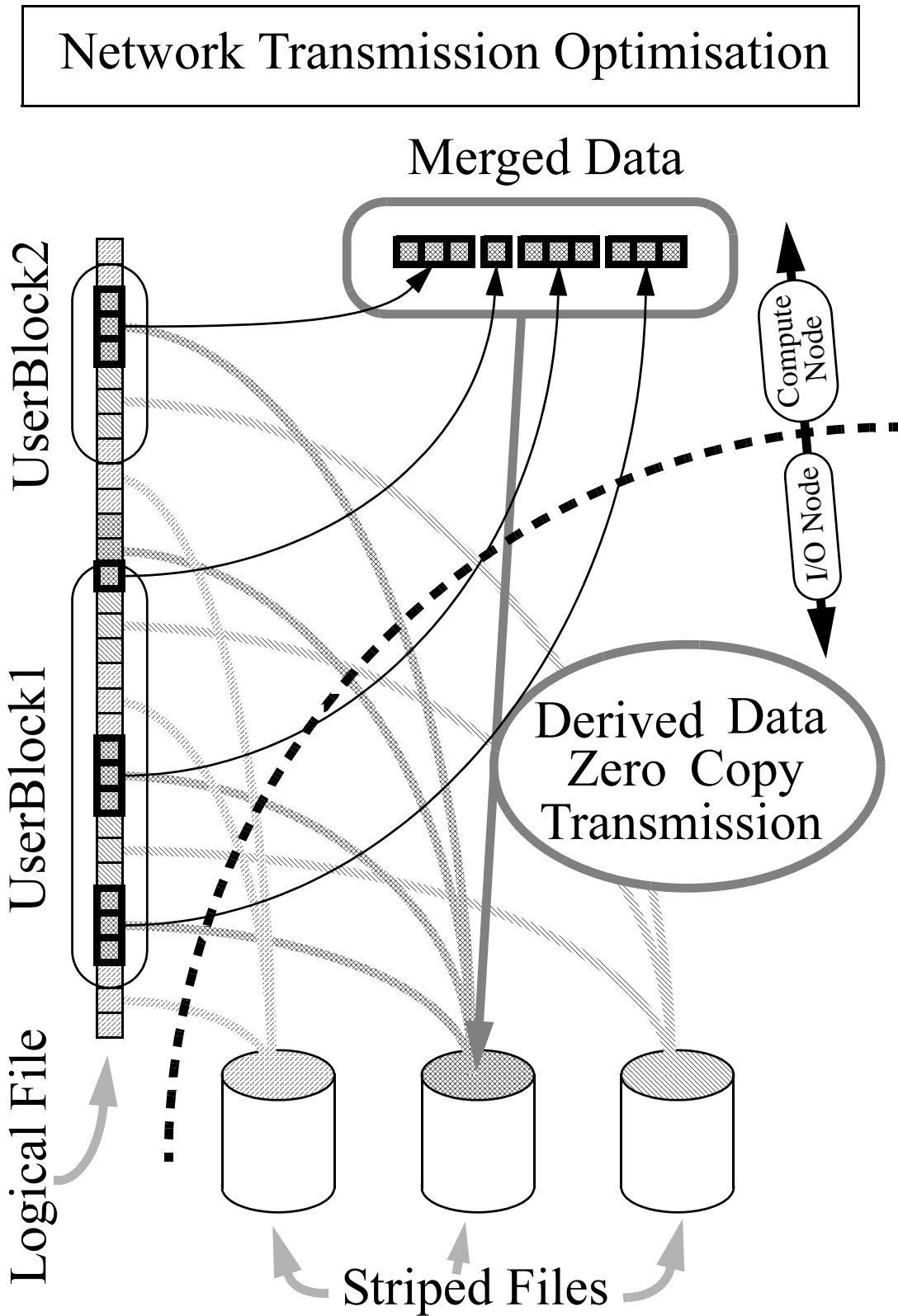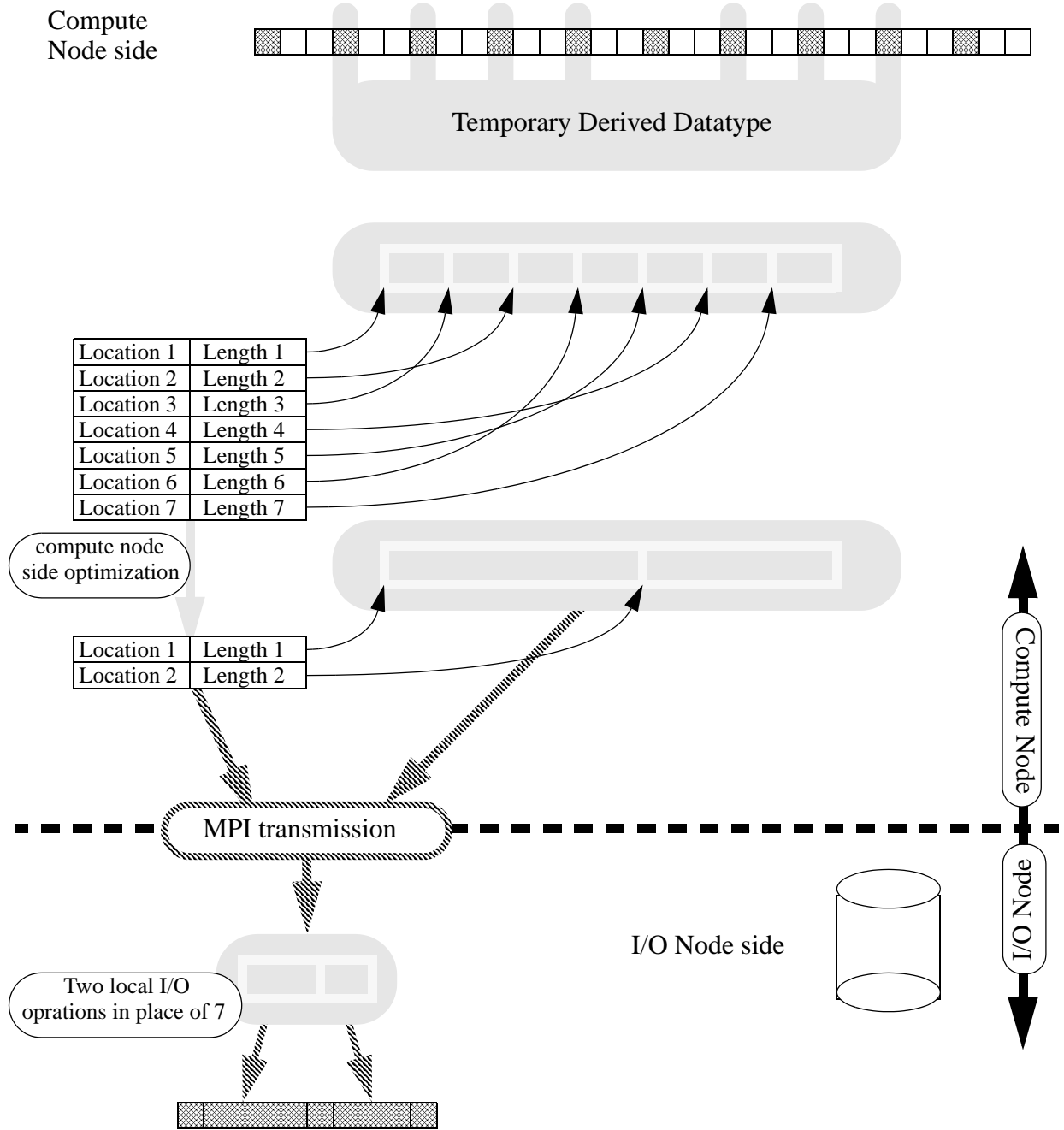
The compute node before accessing I/O node try to reorganise data so as to have
minimum number of contiguous pieces
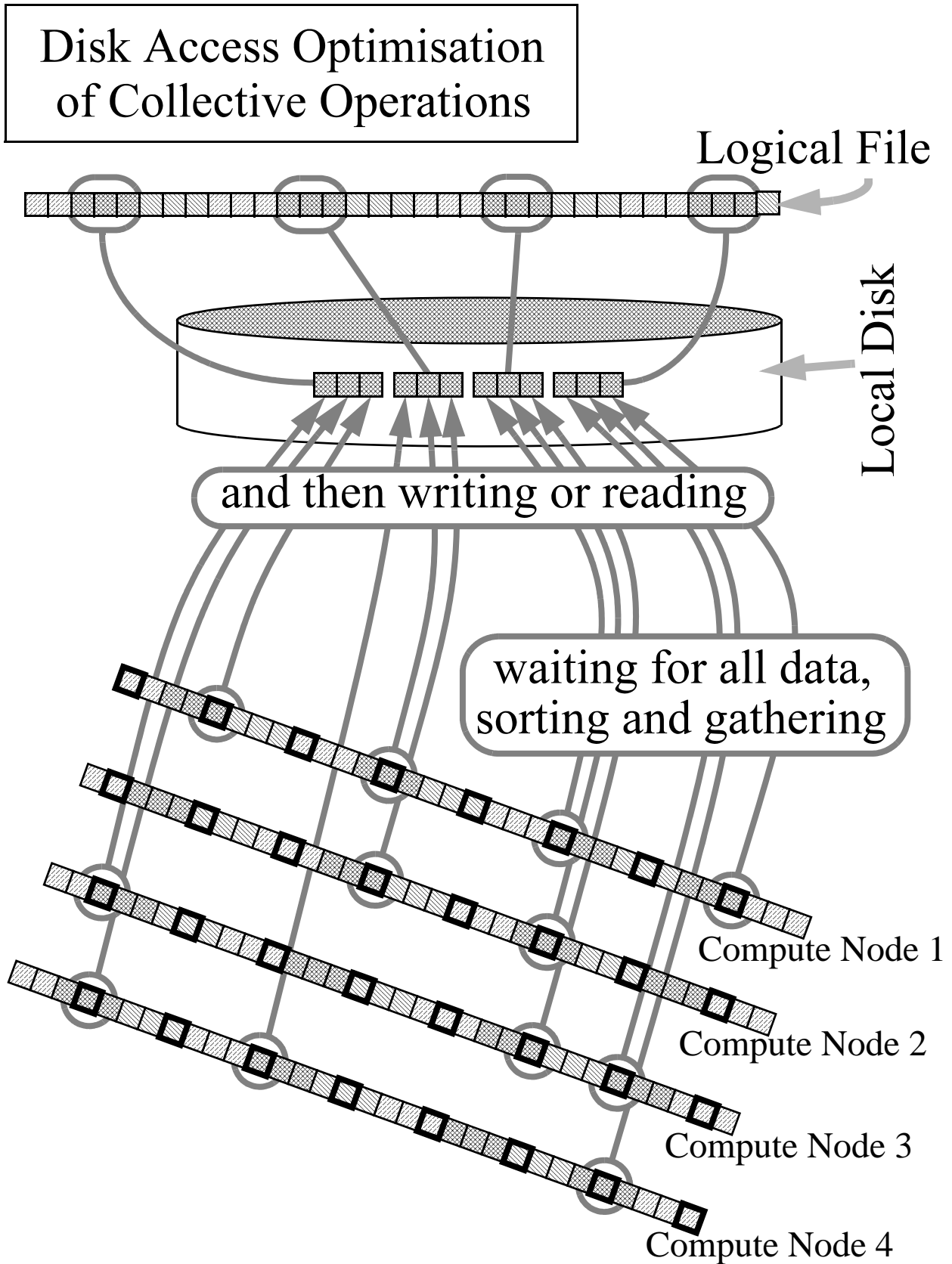
Fig. 05. Disk Access Optiisation of Collective Operations

# Disk Access Optimisation of Collective Operations

Logical File

Local Disk

and then writing or reading

waiting for all data, sorting and gathering

Compute Node 1

Compute Node 2

Compute Node 3

Compute Node 4

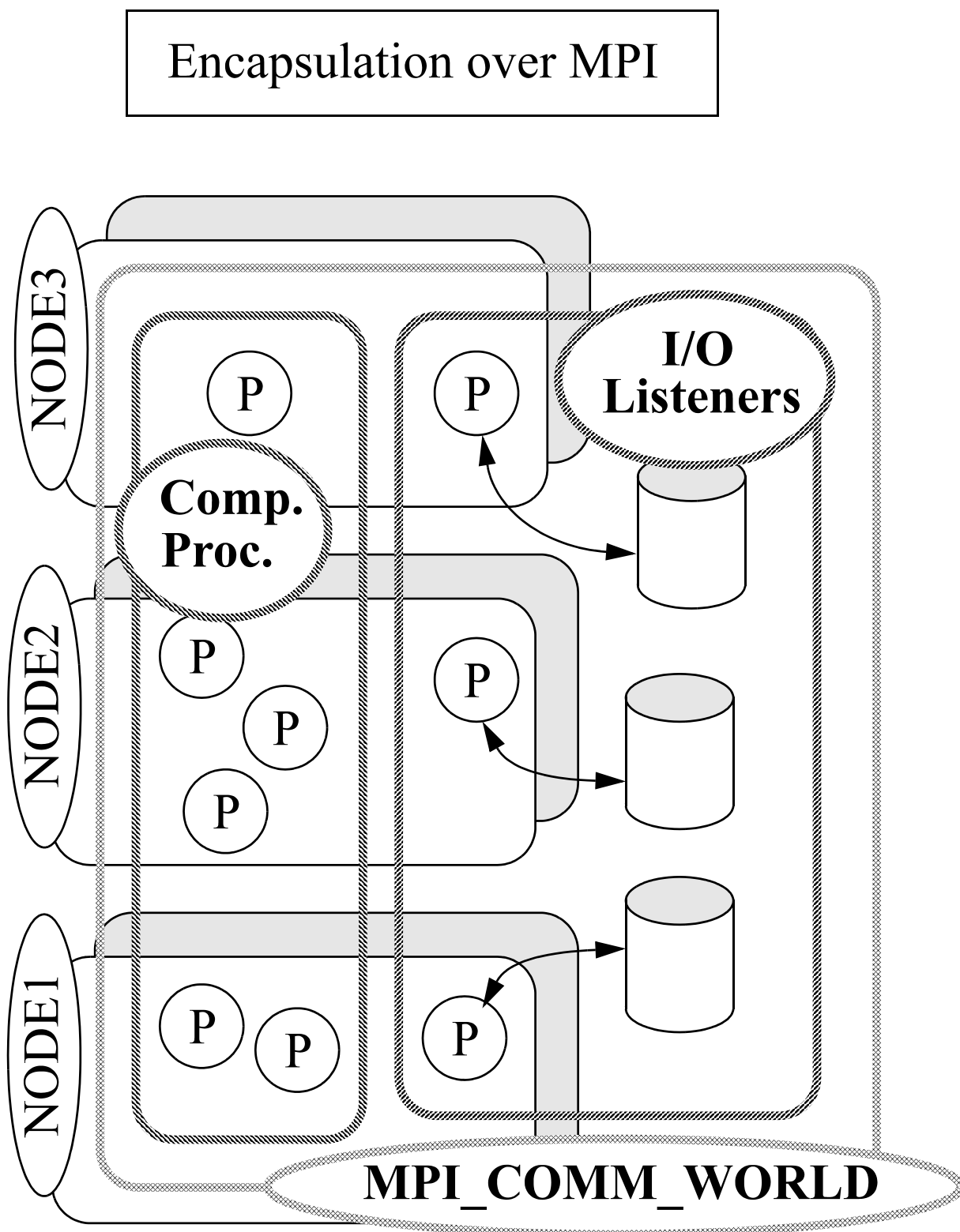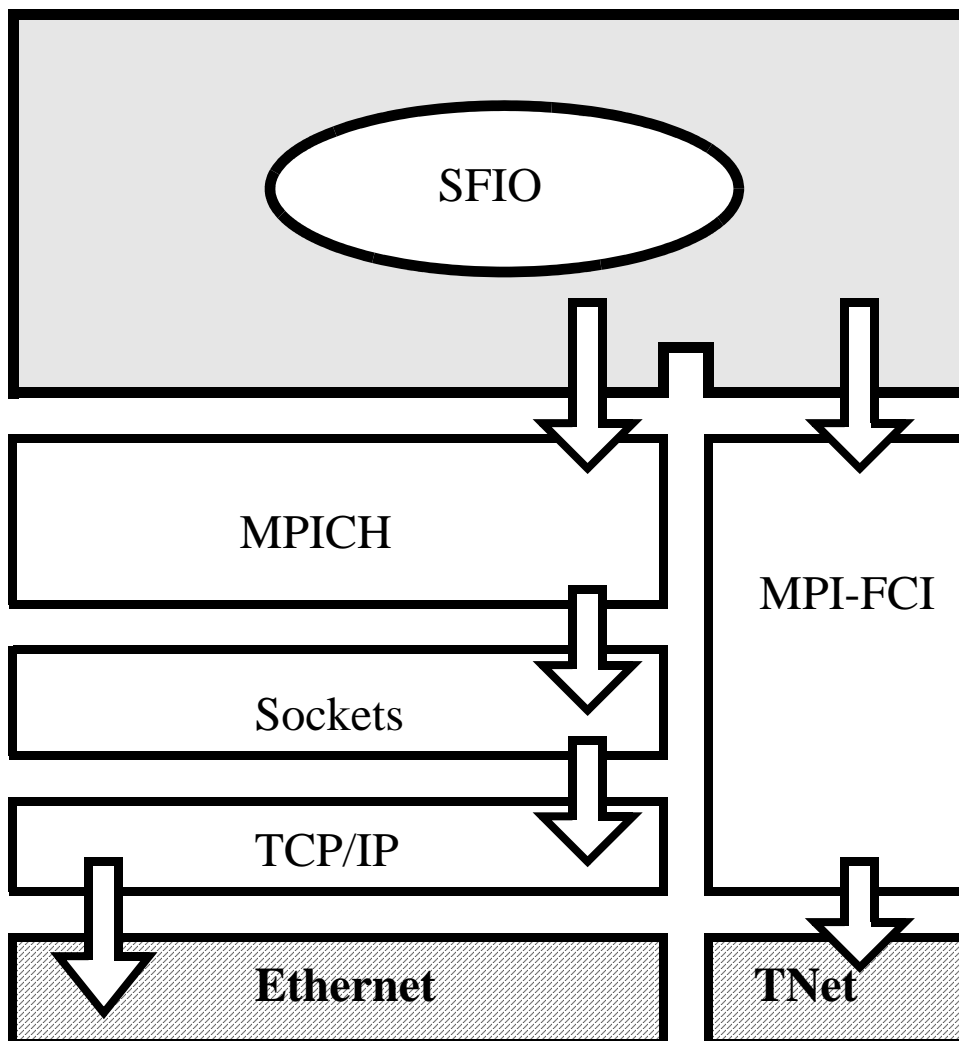Fig. 06.  Encapsulation over MPI

# Encapsulation over MPI

Fig. 07A.  Functionality.



SFIO is implemented and tested with Digital Unix (MPICH, FCI) and Intel Windows NT (MPICH).
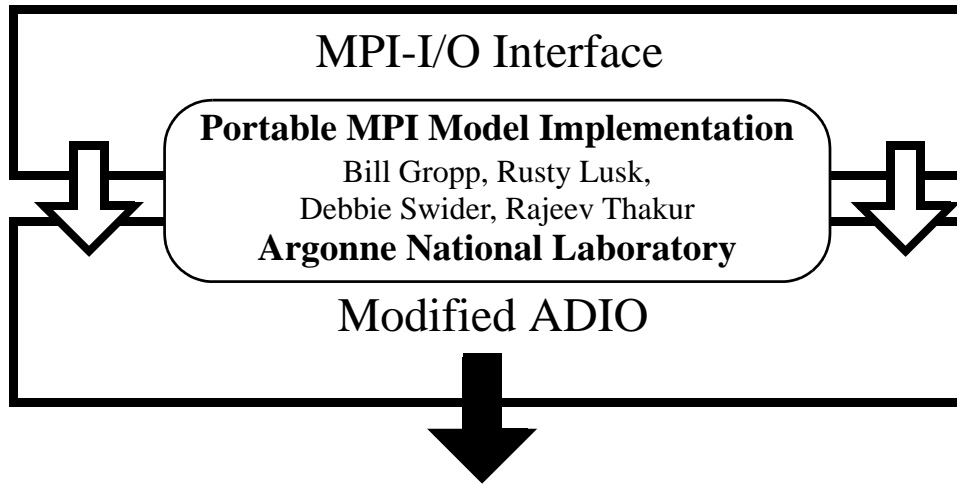
Fig. 07B. Functionality.

MPI-I/O Interface

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

Modified ADIO

Fig. 08.

# 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 |

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

Fig. 09. Optimisation on WinNT (Fast Ethernet Switch)



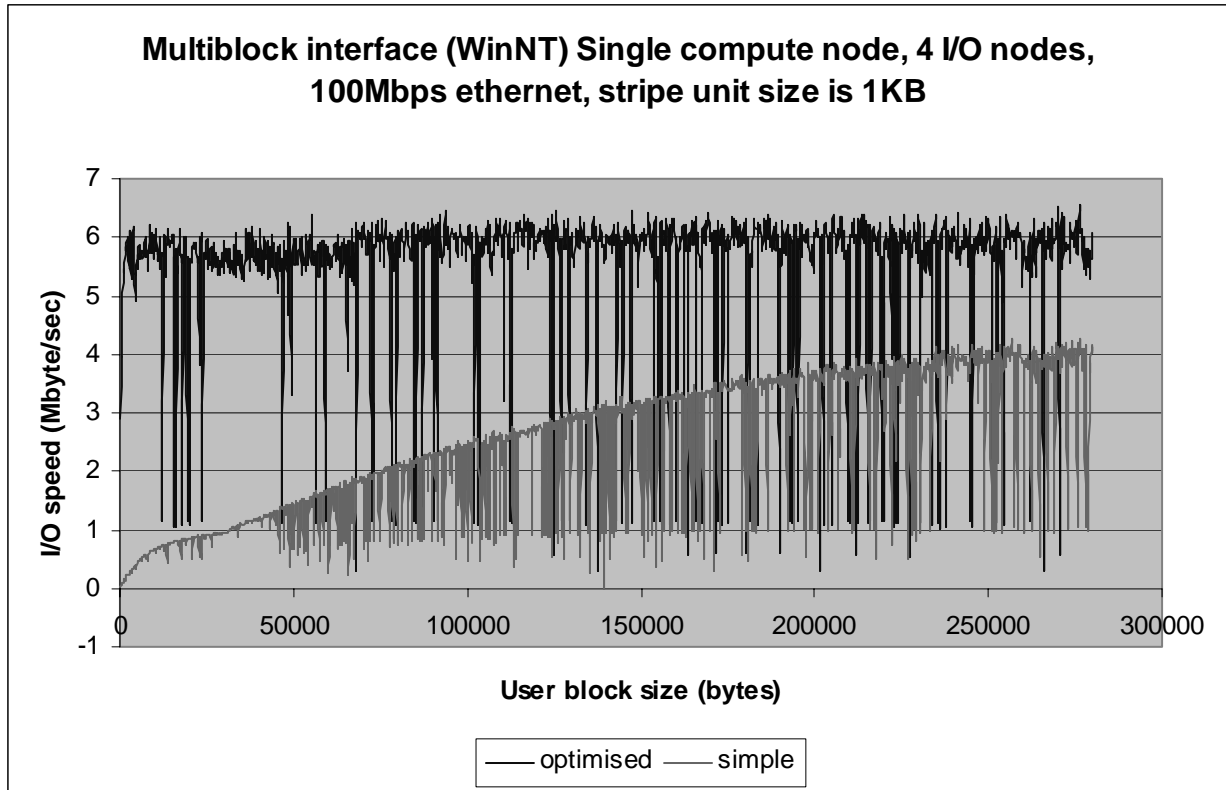**Multiblock interface (WinNT) Single compute node, 4 I/O nodes, 100Mbps ethernet, stripe unit size is 1KB**

Fig. 10. Optimisation on Swiss-Tx (Fast Communication Interface)

Fig. 11.

Integration
into MPI/IO

noncollective

nonblocking

collective

blocking

Synchronism

Coordination

explicit offsets

individual file pointers

shared file pointer

Positioning

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

Positioning

Synchronism

Synchronism

IREAD
AT

IWRITE
AT

IREAD

IWRITE

IREAD
SHARED

IWRITE
SHARED

Synchronism

READ
AT

WRITE
AT

READ

WRITE

READ
SHARED

WRITE
SHARED

Positioning

Synchronism