# Parallel File Striping
# Design trade-offs

**R. D. Hersch, B. Gennart, EPFL**

{hersch,gennart}@di.epfl.ch
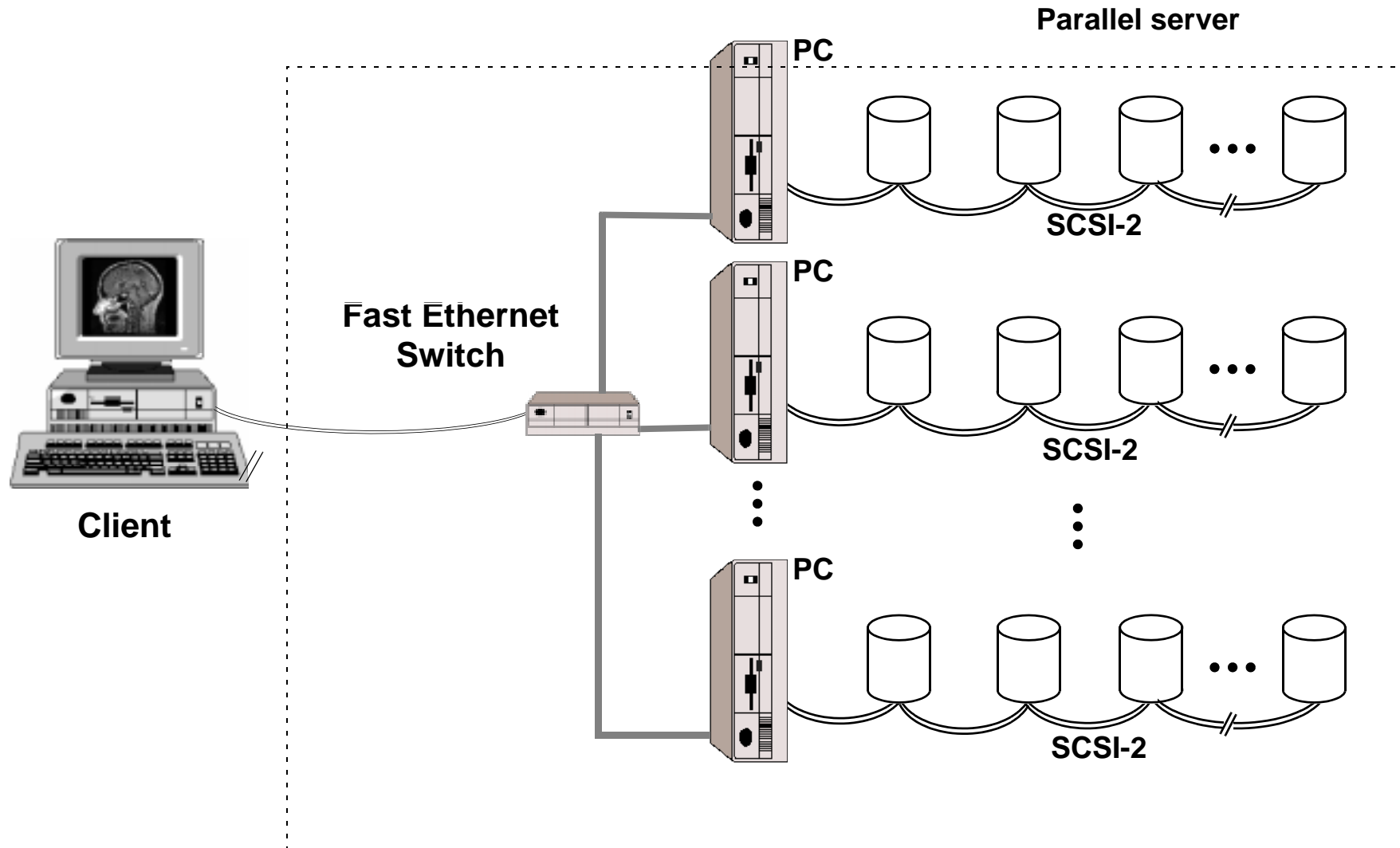http://visiblehuman.epfl.ch

**Content**

1. Context

2. I/O Interface at the striped file level

3. High-level striped file interface

4. Performances of the striped file interface

5. Conclusions

# 1. Context

- Parallel distributed memory servers

  for I/O - & compute-intensive applications

- Assumptions:
  - Parallel file striping library linked
    to application
  - Local files are standard files (NT, Unix)

- Goal: make maximal use of underlying
  hardware and software :
  => Simultaneous I/O accesses,
       communications and computations
  => Generate asynchronous schedule of
       operations

# Considered Architecture

**Parallel server**

**PC**

**PC**

**PC**

**SCSI-2**

**SCSI-2**

**SCSI-2**

**Fast Ethernet Switch**

**Client**

· · ·

· · ·

· · ·

# 2. I/O Interface at the striped file level

- Generating asynchronous schedules with CAP (Computer-Aided Parallelization)

- Example: The Visible Human Slice Server

  http://visiblehuman.epfl.ch

  see next slices:
  - Pipelined parallel slice extraction
  - Performances

- Mapping specified by application:
  Global file position
  -> LocalFileIndex
  -> LocalExtentIndex
  System library specified mapping
  LocalFileIndex -> DiskServerIndex
  DiskServerIndex -> ComputeServerIndex

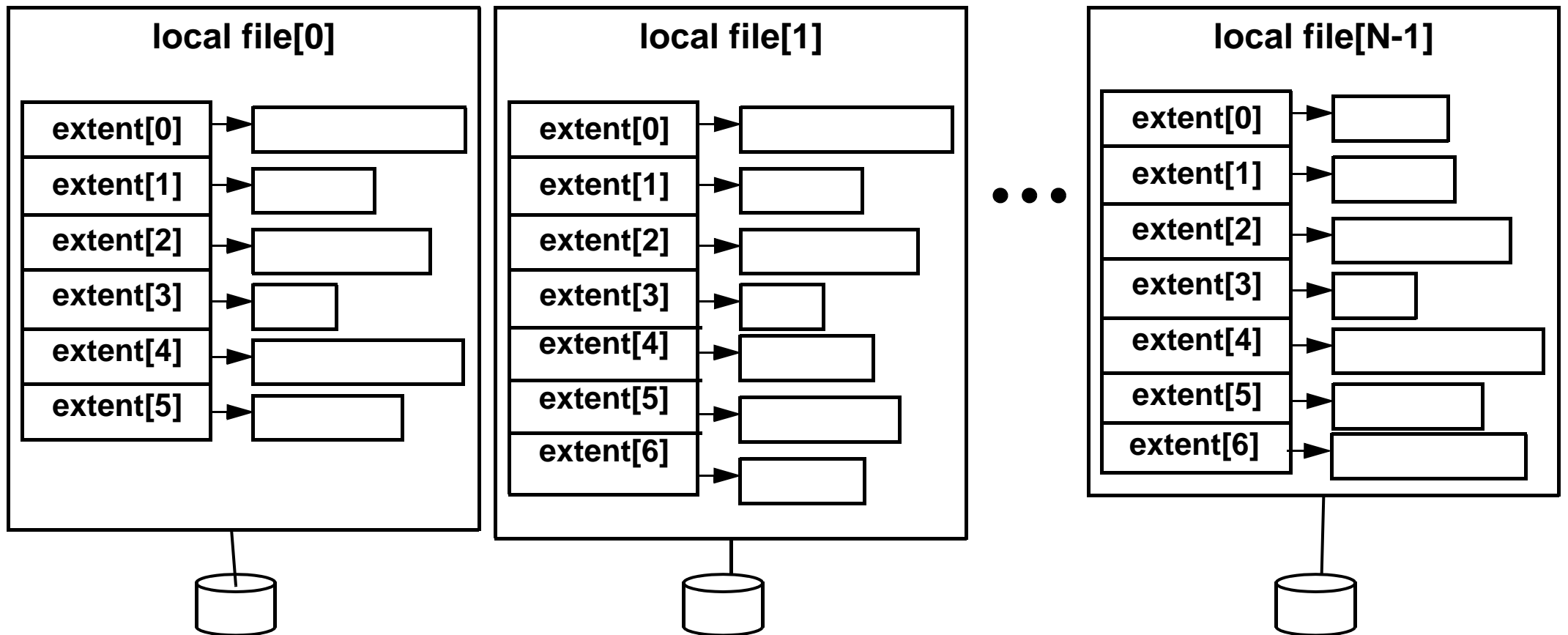- Programmer must deal with 4 indices
  Advantages:
  - When possible, computations on same
    node as disk data
  - Simultaneous computations & disk accesses
    (programmable prefetching strategy)
  - Disk directed I/O feasible, i.e. disk server
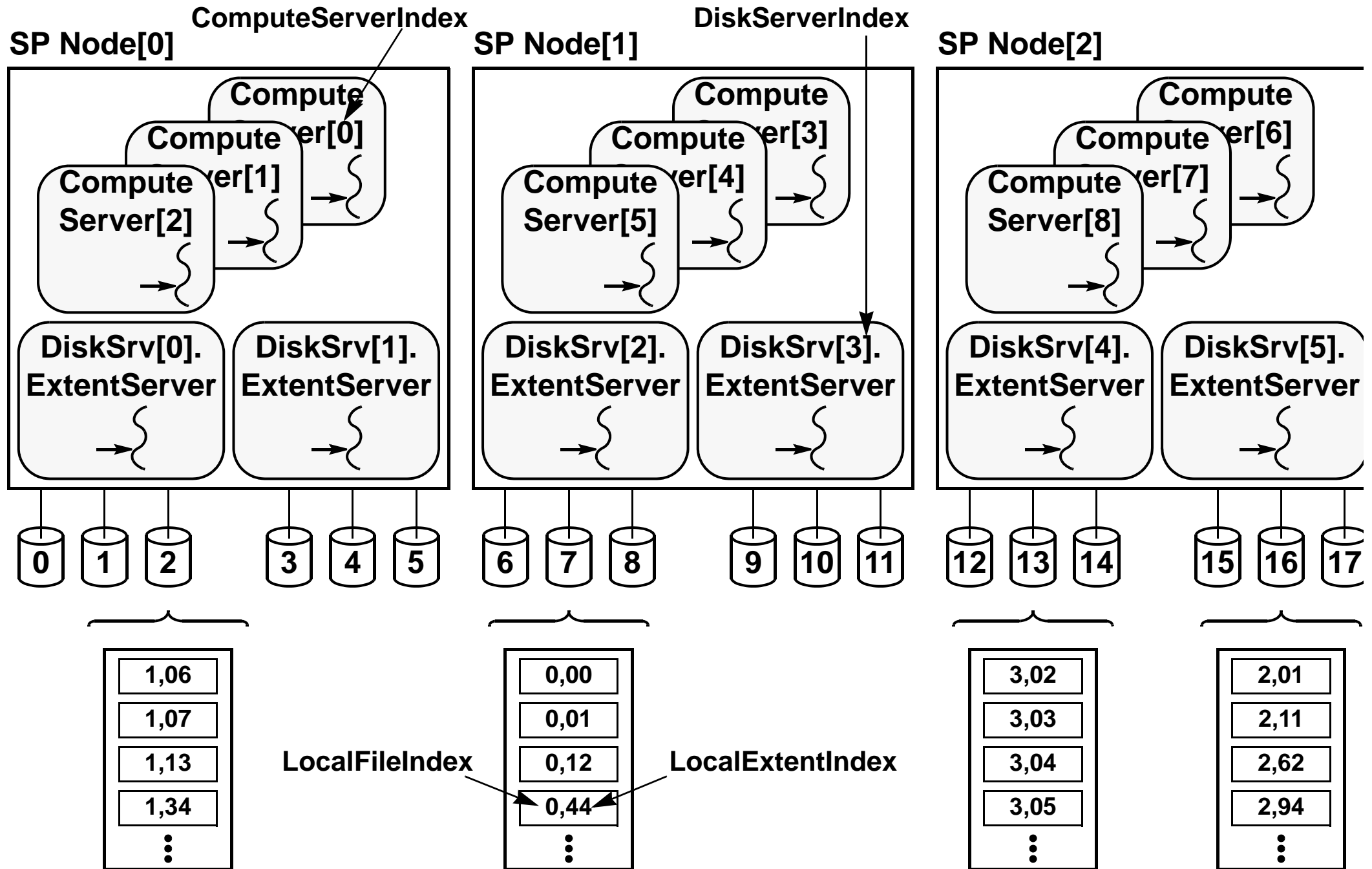    may request data from appl. processes

# Two-Dimensional Parallel File Structure

**extent location given by <LocalFileIndex, LocalExtentIndex>**

**parallel file declustered across N local files (disks)**

| local file[0] | local file[1] | local file[N-1] |
|---|---|---|
| extent[0] | extent[0] | extent[0] |
| extent[1] | extent[1] | extent[1] |
| extent[2] | extent[2] | extent[2] |
| extent[3] | extent[3] | extent[3] |
| extent[4] | extent[4] | extent[4] |
| extent[5] | extent[5] | extent[5] |
|  | extent[6] | extent[6] |

• • •

# LocalFileIndex, LocalExtentIndex, DiskServerIndex & ComputeServerIndex

**ComputeServerIndex**

**DiskServerIndex**

**SP Node[0]**

**SP Node[1]**

**SP Node[2]**

Compute Server[0]

Compute Server[1]

Compute Server[2]

Compute Server[3]

Compute Server[4]

Compute Server[5]

Compute Server[6]

Compute Server[7]

Compute Server[8]

**DiskSrv[0]. ExtentServer**

**DiskSrv[1]. ExtentServer**

**DiskSrv[2]. ExtentServer**

**DiskSrv[3]. ExtentServer**

**DiskSrv[4]. ExtentServer**

**DiskSrv[5]. ExtentServer**

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17

| 1,06 |
| 1,07 |
| 1,13 |
| 1,34 |
⋮

| 0,00 |
| 0,01 |
| 0,12 |
| 0,44 |
⋮

| 3,02 |
| 3,03 |
| 3,04 |
| 3,05 |
⋮

| 2,01 |
| 2,11 |
| 2,62 |
| 2,94 |
⋮

**LocalFileIndex**

**LocalExtentIndex**

# Example: the Visible Human Slice Server

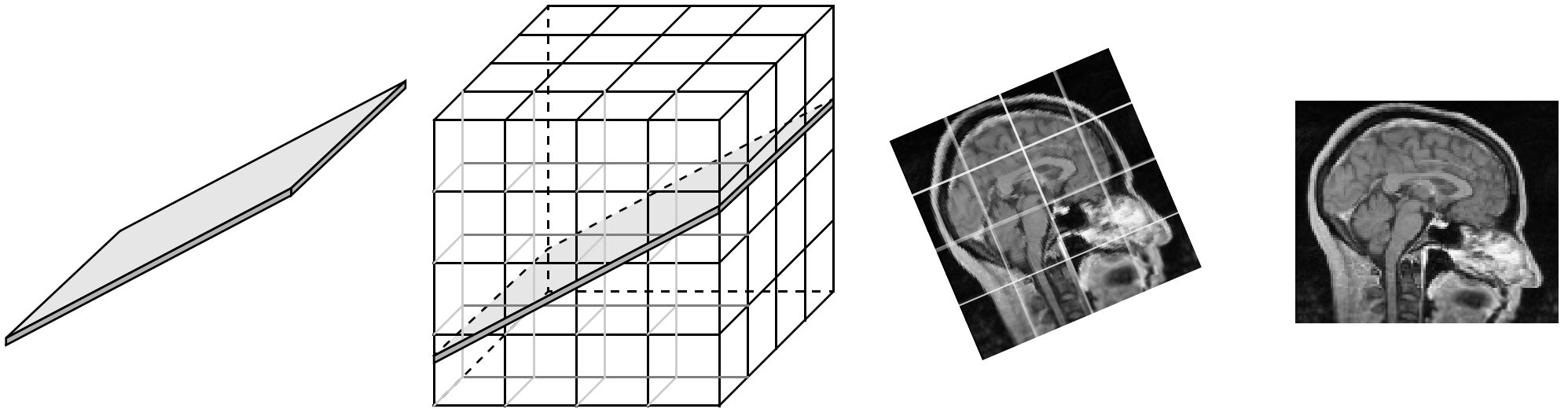(http://visiblehuman.epfl.ch)

- Human Male: 13 GB of data striped
  over 60 disks (sub-volumes: file extents)

- Resolution:  one anatomical slice/mm
  in each anatomical slice: 3 samples/mm

- Size: 1840 horizontal slices,
          each 2048x1216 pixels

- Software :

  - client sends slice position & orientation

  - server interface asks for subslices

  - server processors read sub-volumes,
    extract and project slices (resampling)

  - server interface assembles resulting full
    slice, compresses it and sends it to client

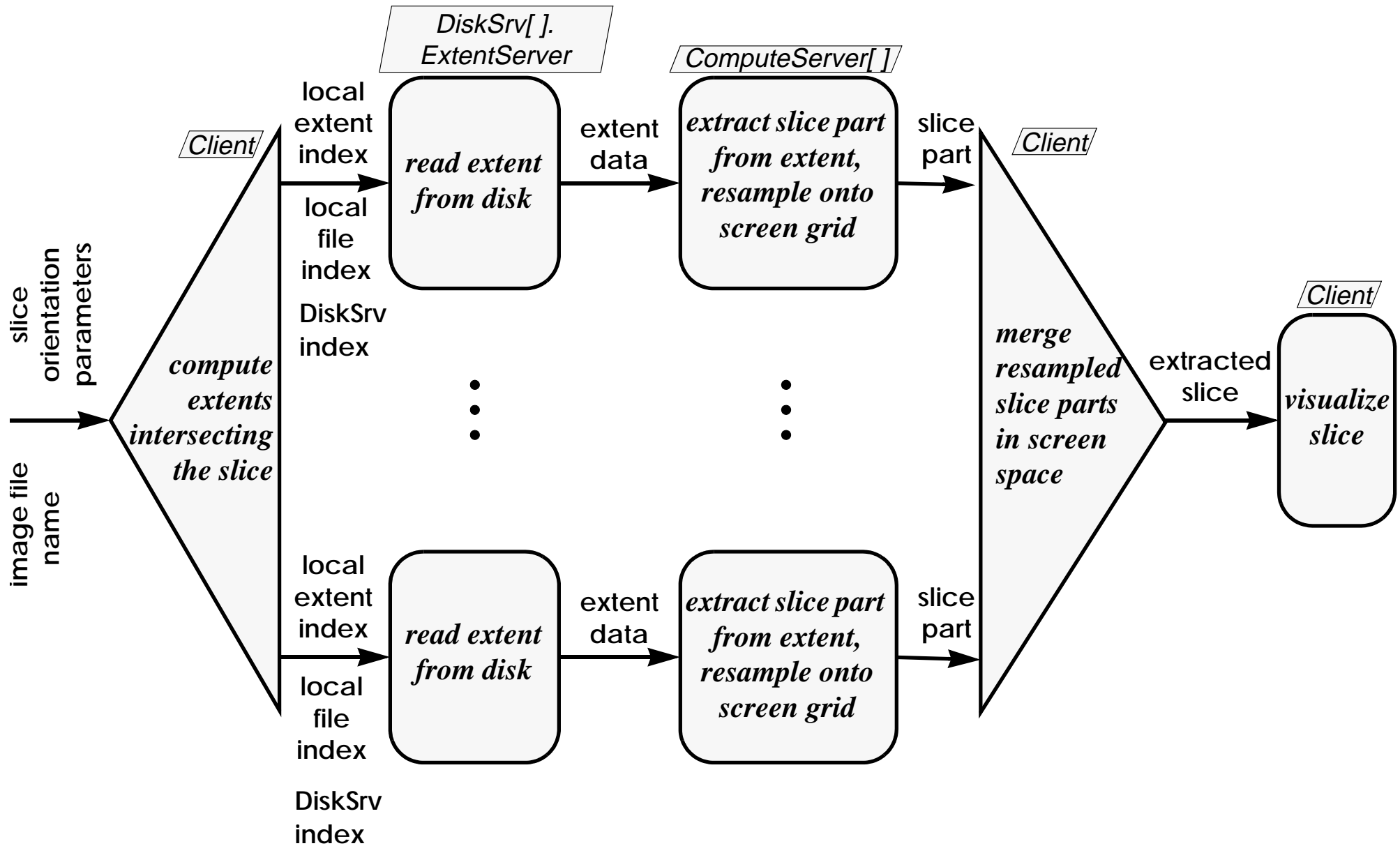# Extraction of Slice Parts from Volumic File Extents



Extraction slice specification

Extraction of the digital slice from the 3D image

Extracted slice parts

Resampled slice parts merged into the final displayable slice

# Pipelined Parallel Slice Extraction

slice orientation parameters

image file name

*Client*

**compute extents intersecting the slice**

local extent index

local file index

DiskSrv index

*DiskSrv[ ]. ExtentServer*

**read extent from disk**

extent data

*ComputeServer[ ]*

**extract slice part from extent, resample onto screen grid**

slice part

*Client*

local extent index

local file index

DiskSrv index

**read extent from disk**

extent data

**extract slice part from extent, resample onto screen grid**

slice part

**merge resampled slice parts in screen space**

extracted slice

*Client*

**visualize slice**

# CAP Specification of the Pipelined Parallel Slice Extraction

```
1   leaf operation Ps2ExtentServerT::ReadExtent
2     in ExtentReadingRequestT* InputP
3     out ExtentT* OutputP;
4
5   leaf operation Ps2ComputeServerT::ExtractAndResampleSlicePart
6     in ExtentT* InputP
7     out SlicePartT* OutputP
8   {  ...C++ sequential code  }
9
10  bool SplitSliceRequest(SliceExtractionRequestT* FromP,
11                         ExtentReadingRequestT* PreviousP,
12                         ExtentReadingRequestT* &ThisP)
13  {
14    ...C++ sequential code
15    return (IsNotLastExtentReadingRequest);
16  }
17
18  void MergeSlicePart(SliceT* IntoP, SlicePartT* ThisP)
19  {
20    ...C++ sequential code
21  }
22
23  operation Ps2ServerT::ExtractSlice
24    in SliceExtractionRequestT* InputP
25    out SliceT* OutputP
26  {
27    parallel while (SplitSliceRequest, MergeSlicePart, Client, SliceT Output)
28    (
29      DiskSrv[thisTokenP->DiskServerIndex].ExtentServer.ReadExtent
30      >-->
31      ComputeServer[thisTokenP->ComputeServerIndex].ExtractAndResampleSlicePart
32    )
33  }
```

see http://diwww.epfl.ch/w3lsp/pub/publications/gigaview/captutorial.pdf

# Visible Human Slice Extraction Performances
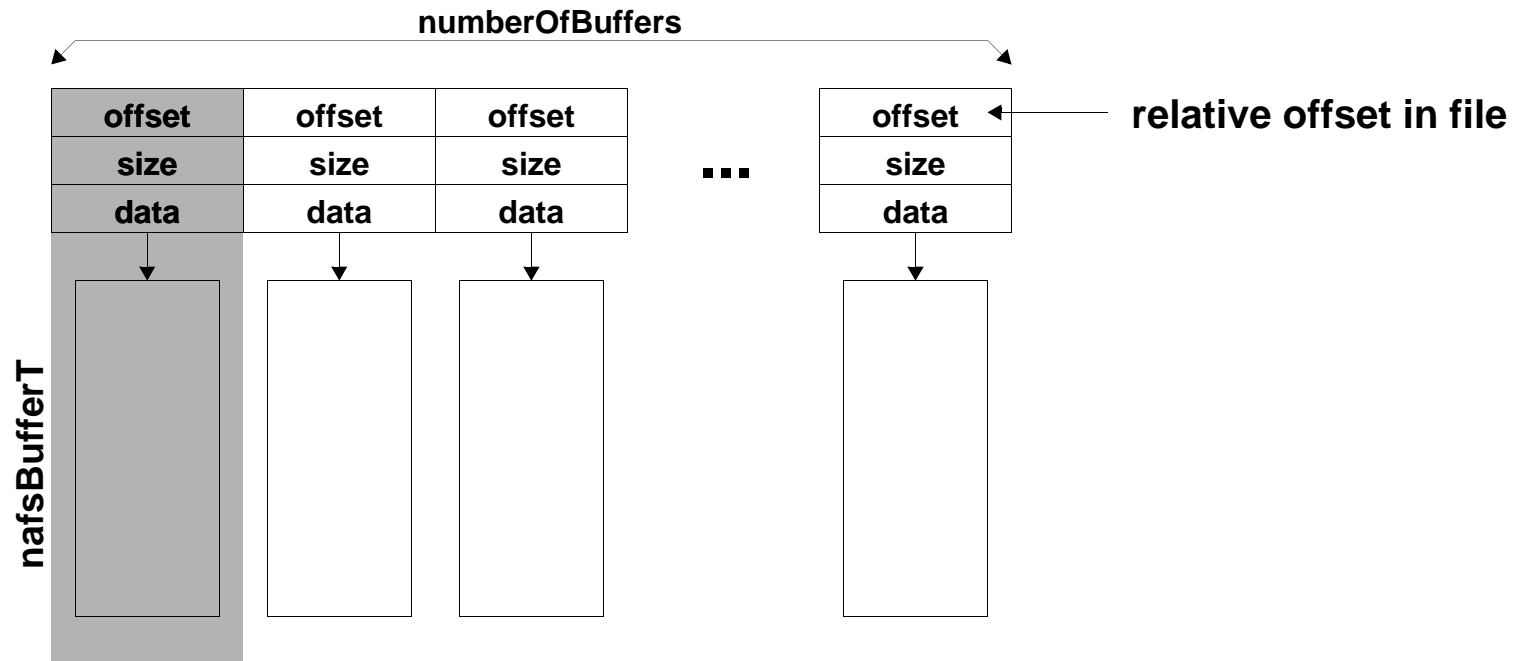
- Max. configuration, caches disabled
  5 server PC's, each with 12 disks

  4.8 extracted slices/s
  client processor utilization: 85%
  => 104 MB/s read from 60 disks
  => mean throughput per disk 1.74 MB/s
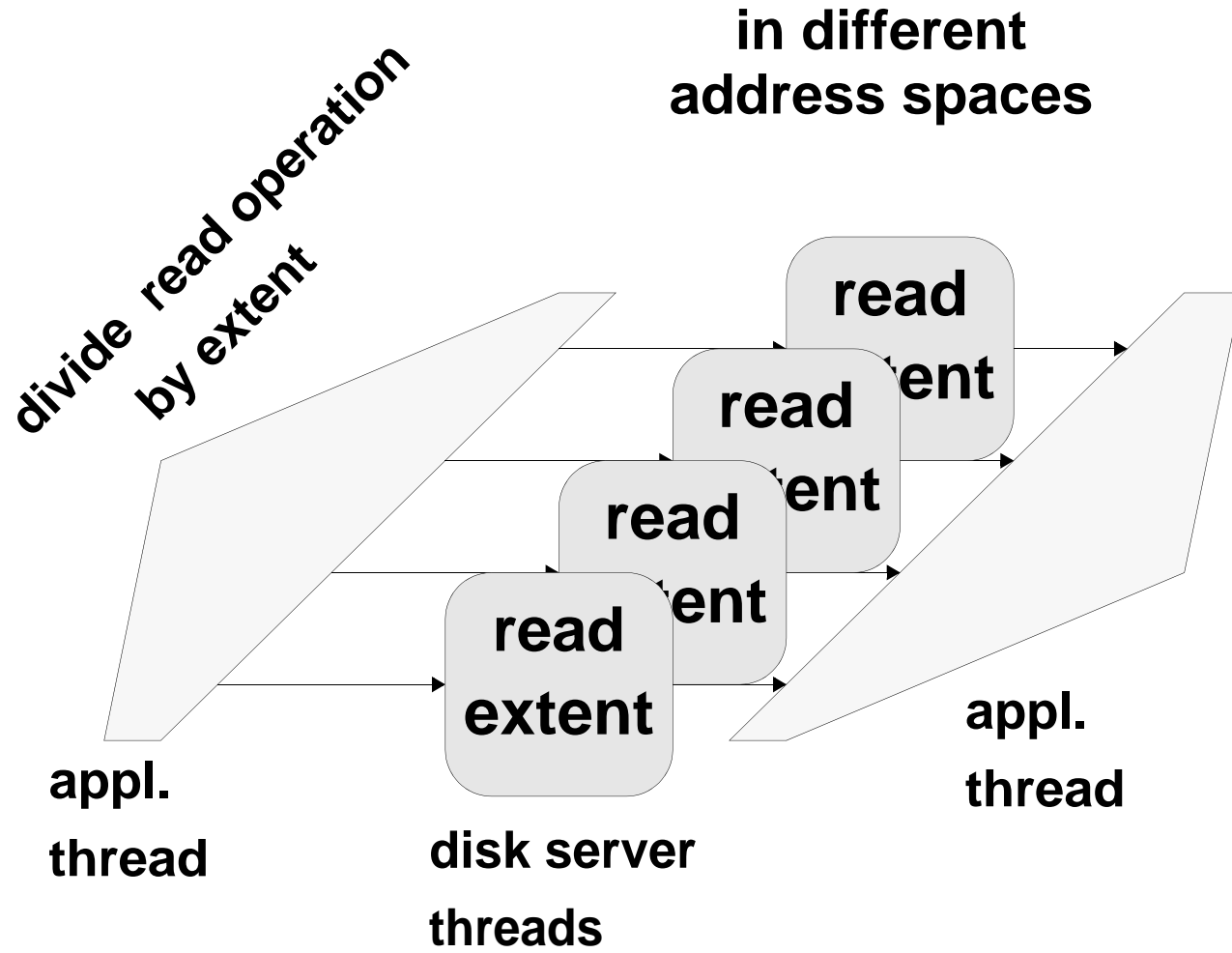   => disks are bottleneck

# 3. High-level Striped File Interface (NAFS)

- Interface translates global byte positions to LocalFileIndices, LocalExtentIndices and local byte positions

- User specified stripe unit, striping over a set of user specified paths (disks).

- Read, Write, CollectiveRead, CollectiveWrite of a set of contiguous data buffers

- Assessment: Simple interface, but no direct knowledge about locality of data
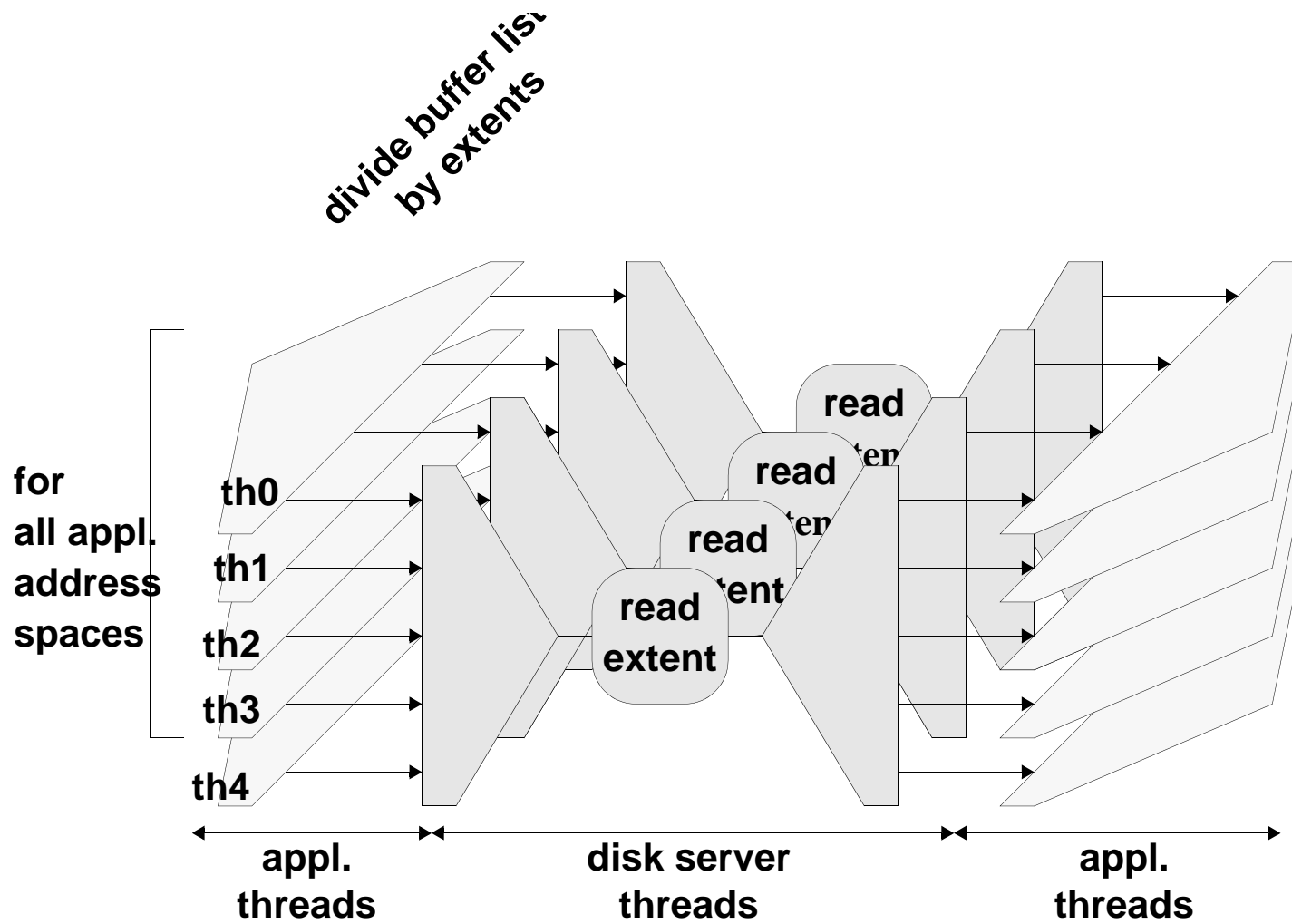
# High-level striped file interface (NAFS)

```
nafsFileT* nafsOpen (char* fileNameS, char openMode) ;
int nafsClose (nafsFileT* fileP) ;
int nafsCreate (char* fileNameS) ;
int nafsDelete (char* fileNameS) ;
int nafsWrite (nafsFileT* fileP, UINT64_T absoluteOffsetInFile,
            UINT32_T numberOfBuffers, nafsBufferT* buffersP) ;
int nafsRead  (nafsFileT* fileP, UINT64_T absoluteOffsetInFile,
            UINT32_T numberOfBuffers, nafsBufferT* buffersP) ;
int nafsCollectiveWrite (nafsFileT* fileP, UINT64_T absoluteOffsetInFile,
                UINT32_T numberOfBuffers, nafsBufferT* buffersP) ;
int nafsCollectiveRead (nafsFileT* fileP, UINT64_T absoluteOffsetInFile,
                UINT32_T numberOfBuffers, nafsBufferT* buffersP) ;
```

# Read operation

divide  read operation
by extent

in different
address spaces

**read**
~~ent~~

**read**
~~ent~~

**read**
~~ent~~

**read
extent**

**appl.
thread**

**disk server
threads**

**appl.
thread**

# Collective read operation



divide buffer list
by extents

for
all appl.
address
spaces

th0
th1
th2
th3
th4

read
extent

read
extent

read
extent

read
extent

appl.
threads

disk server
threads

appl.
threads

# 4.  Performance of the high-level Striped File Interface

- Configuration:
  4 Bi-PentiumPro PC's, with 8 disks (5400 rpm), interconnected by Fast Ethernet switch

- Writes to local disks: 8 x 2.5MB/s => 20MB/s (32KB data chunks), no much processing power used

- Collective write: each processor writes to each disk many 2KB data chunks

  max sustainable Fast Ethernet throughput per PC: 4MB/s (100% processor utilization)

  => 2MB/s for sending out data from application
     2 MB/s for receiving data to the disks

  => 2MB/s per node global write throughput
     half the processing power used

# 5. Conclusions

- As long as TCP-IP is used, minimize accesses to remote disks (low-level striped file interface)

- Even for I/O
  fast communication network is required
  => Socket based interface, with flow
  control but without TCP-IP processing
  overhead

- High-Performance File Server:
  TCP-IP data transfers considerably reduce
  effective RAID access bandwidth
  => during computations, parallel I/O
  across a large number of nodes is
  preferable (SwissTx:  3MB/s
  per node expected)

- After computation session,
  save striped files to file server
  pmv, pcp: parallel move or parallel copy

- In case of processor/disk failure, current
  results are lost; relaunch computation

- Option: delayed generation of a parity file
  RDH/March 23[rd], 1999, SOS Workshop, Villard, Switzerland