

Dynamic optimization techniques to enhance scalability and performance of MPI-based applications on Multi-core cluster.

Author: Rosa Filgueira Vicente
University: University Carlos III of Madrid

Summary

1. Problem description
2. Main Objectives
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance communication operations.
5. Evaluation tools.

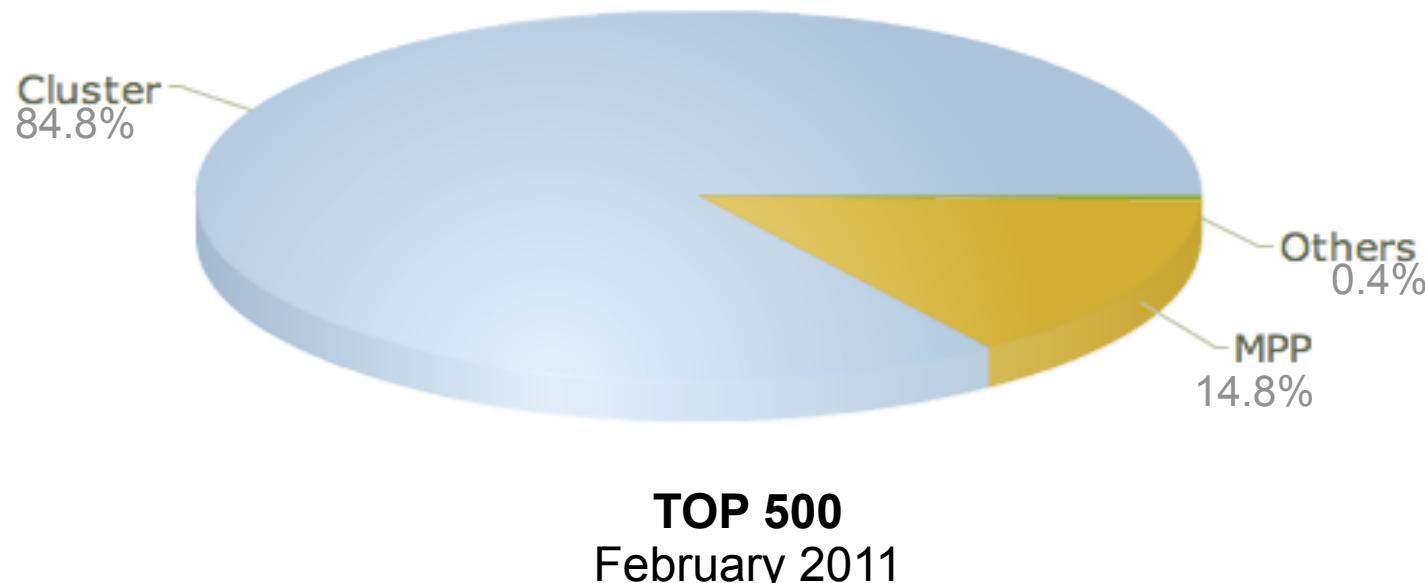
Summary

1. Problem description
2. Main Objectives
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance communication operations.
5. Evaluation tools

Problem description

4

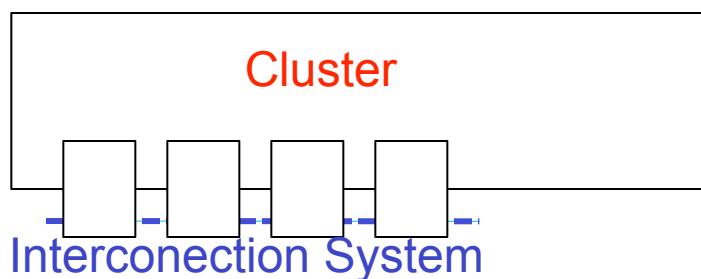
- Parallel computation on cluster has become the most common solution for HPC application.



Problem description

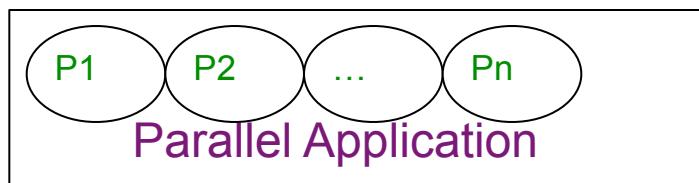
5

- Cluster is a group of linked computers, working together closely thus in many respects forming a single computer. The components of a cluster are commonly, but not always, connected to each other through fast local area networks

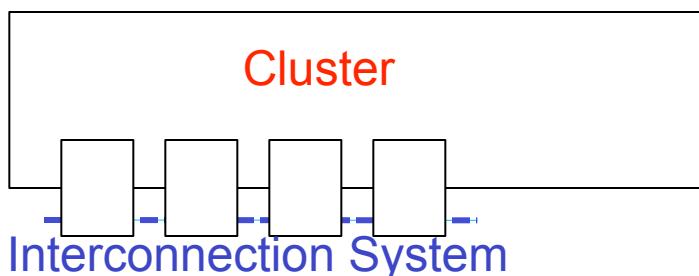


Problem description

6



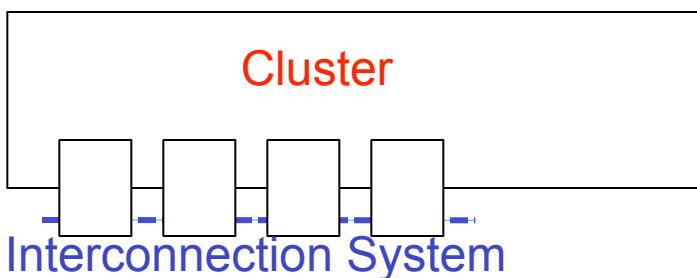
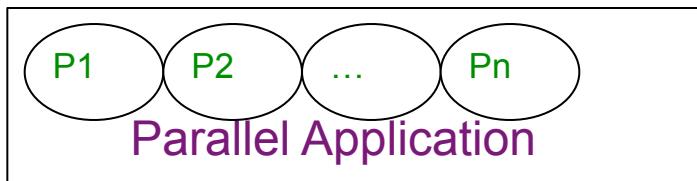
- On a cluster, one/many parallel applications could be running.



Problem description

7

- One of most communication middleware used is the standard **MPI** (Message Passing Interface)



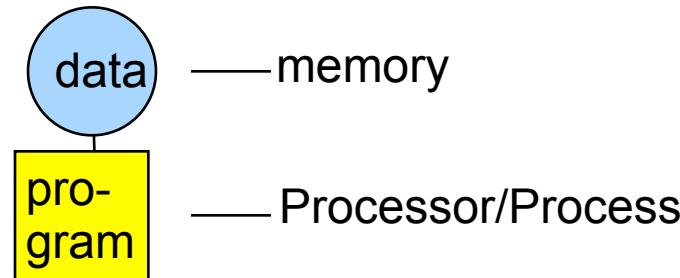
What is MPI (Message Passing Interface)?

- Standardized message passing library specification (IEEE)
 - for parallel computers, clusters and heterogeneous networks
 - not a specific product, compiler specification etc.
 - many implementations, MPICH, LAM, OpenMPI ...
- Portable, with Fortran and C/C++ interfaces.
- Many functions
- Real parallel programming

The Message-Passing Programming Paradigm

9

- Sequential Programming Paradigm



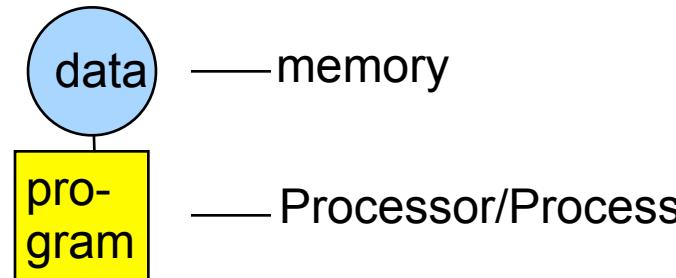
□

9

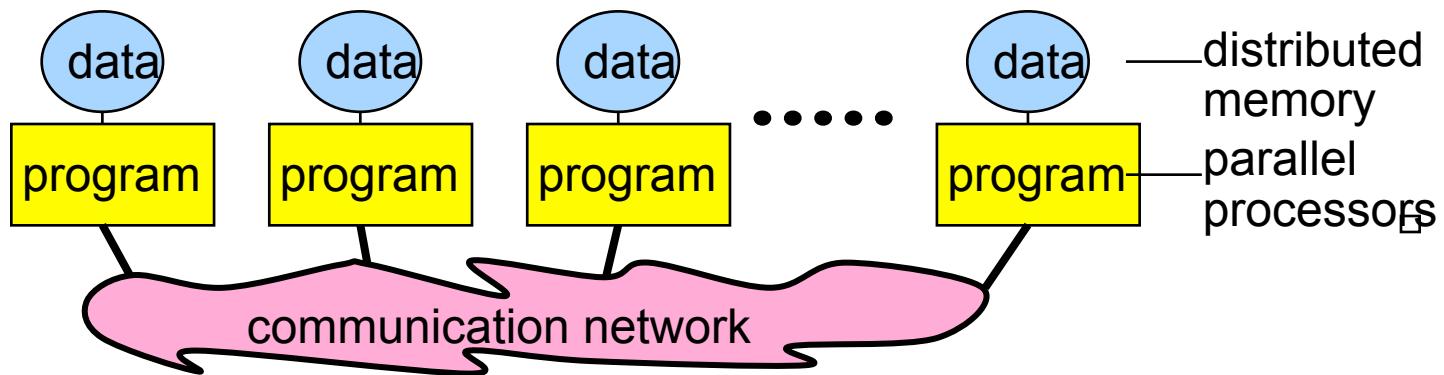
The Message-Passing Programming Paradigm

9

- Sequential Programming Paradigm



- Message-Passing Programming Paradigm

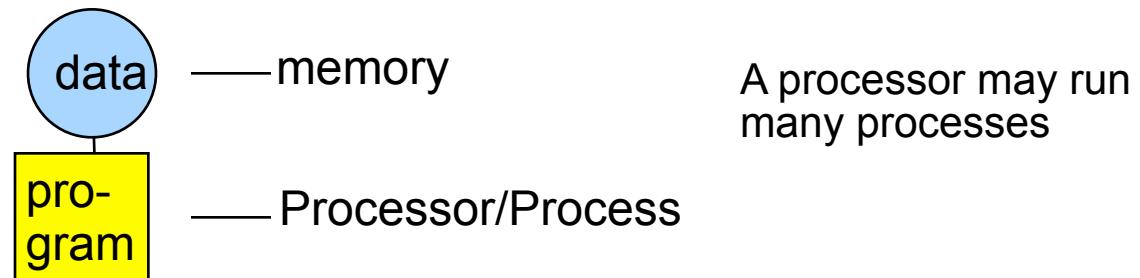


9

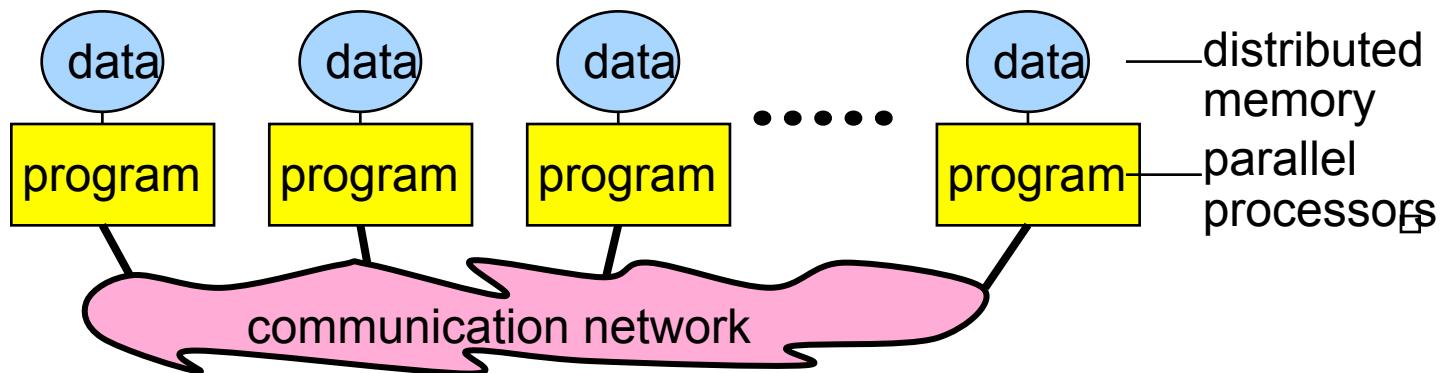
The Message-Passing Programming Paradigm

9

- Sequential Programming Paradigm



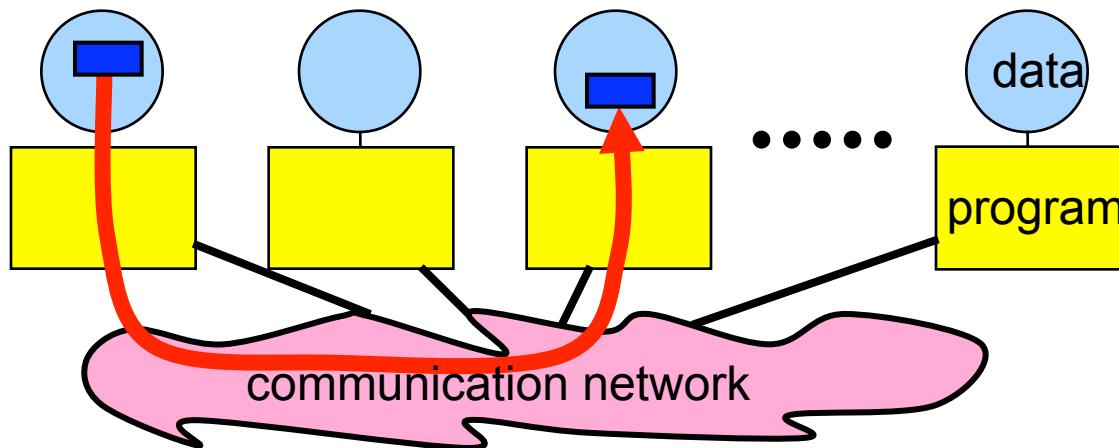
- Message-Passing Programming Paradigm



9

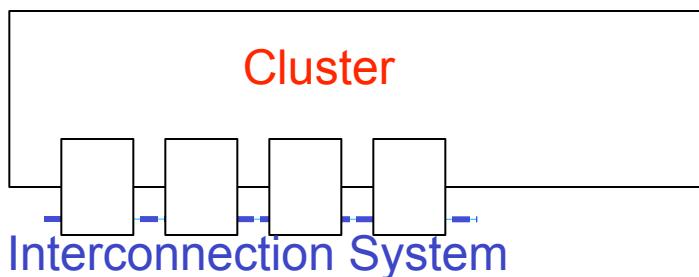
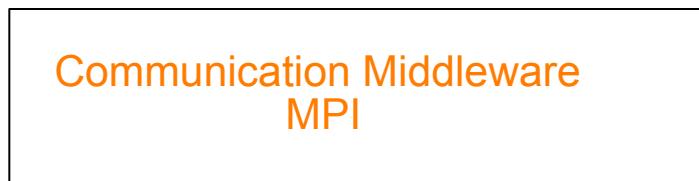
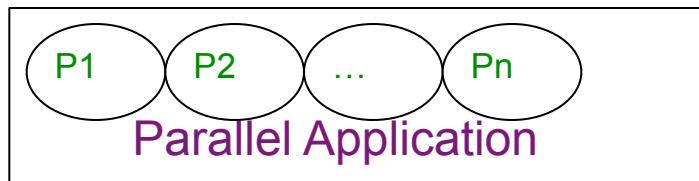
Message passing

- Messages are packets of data moving between sub-programs
- Necessary information for the message passing system:
 - sending process
 - source location
 - source data type
 - source data size
 - receiving process
 - destination location
 - destination data type
 - destination buffer size

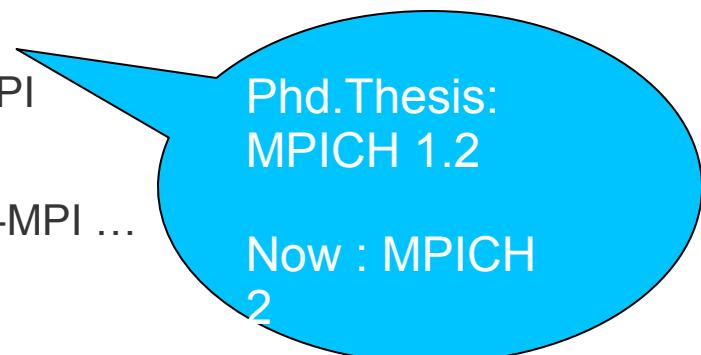


Problem description

11

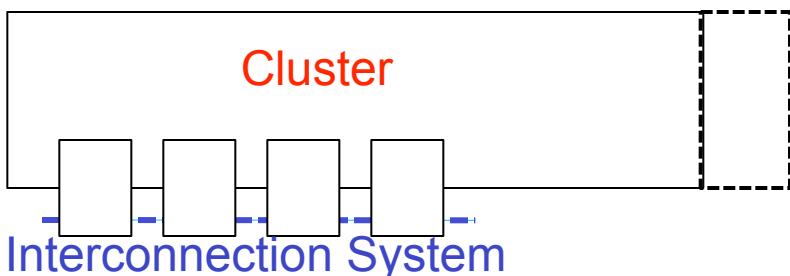
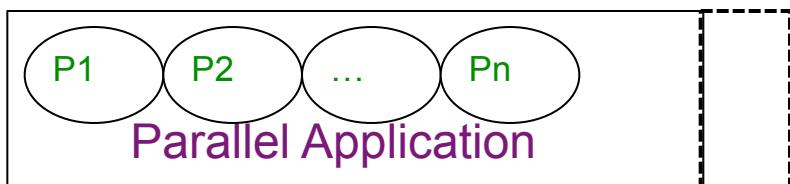


- One of most communication middleware used is the standard **MPI** (Message Passing Interface)
- Different implementations:
 - MPICH
 - OpenMPI
 - LAM
 - CHIMP-MPI ...



Problem description: trend

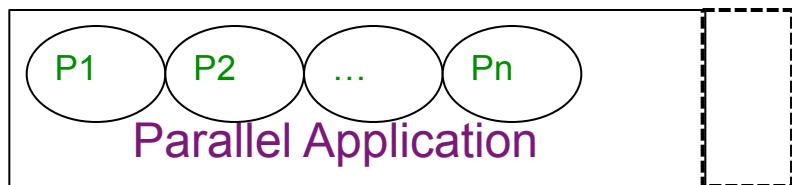
12



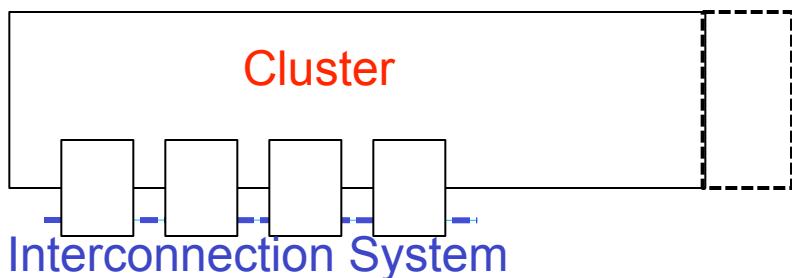
- Multicore processors provides a flexible way to increase the computation capability of cluster
- System performance may be improved with multicore **but** bottleneck from other components could reduce the scalability.

Problem description: Bottleneck

13

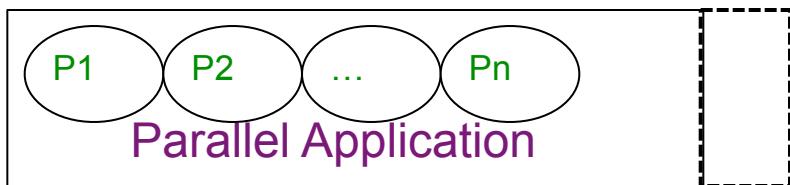


- **Bottleneck :**
 - I/O subsystem
 - Communication subsystem

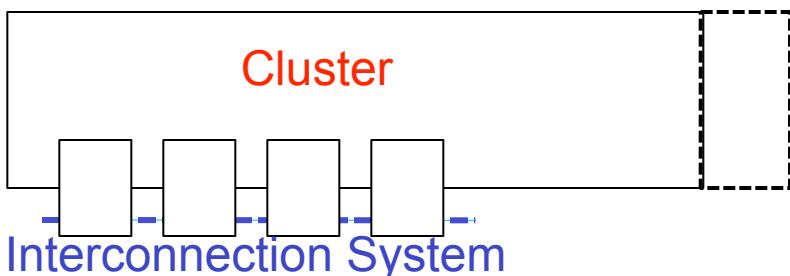


Problem description: Bottleneck

14

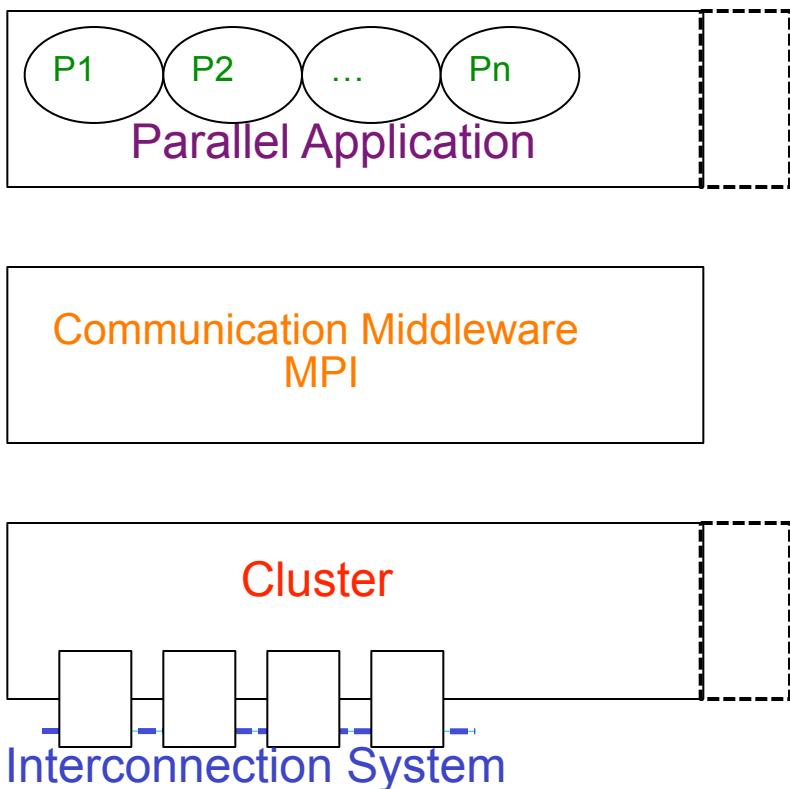


- I/O Bottleneck:
 - I/O requests initiated by multiple cores and besides when non-contiguous disk accesses is used.



Problem description: Bottleneck

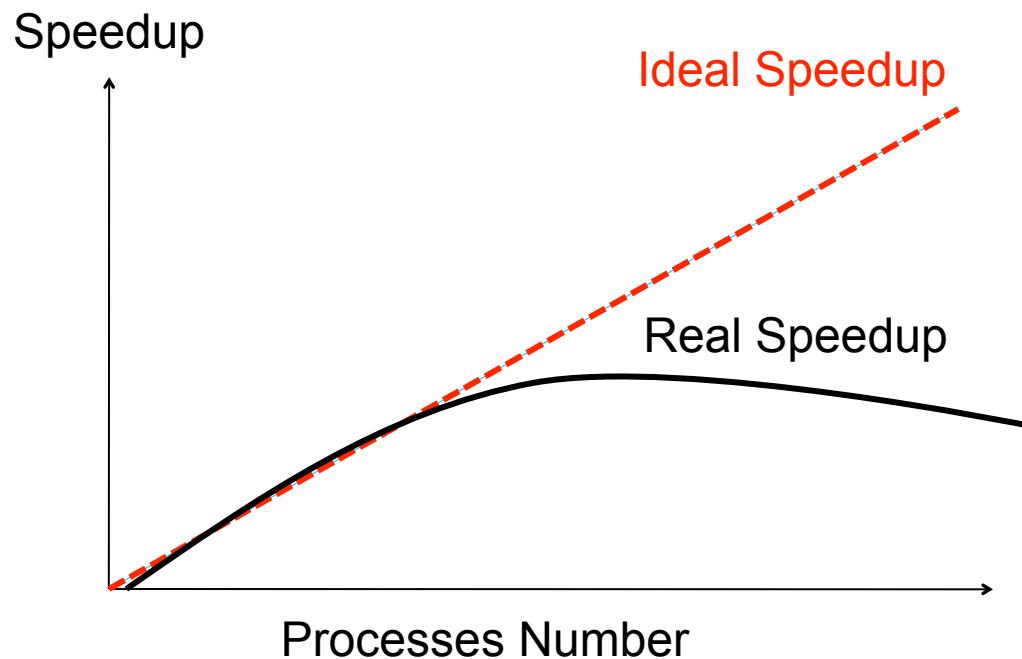
15



- **Communication Bottleneck:**
 - Network used are very fast and low latency.
 - Computational capability in multicore very high.
 - The frequency of message increase a lot
→ insufficient the increase of the bandwidth and latency

Problem description: Bottleneck

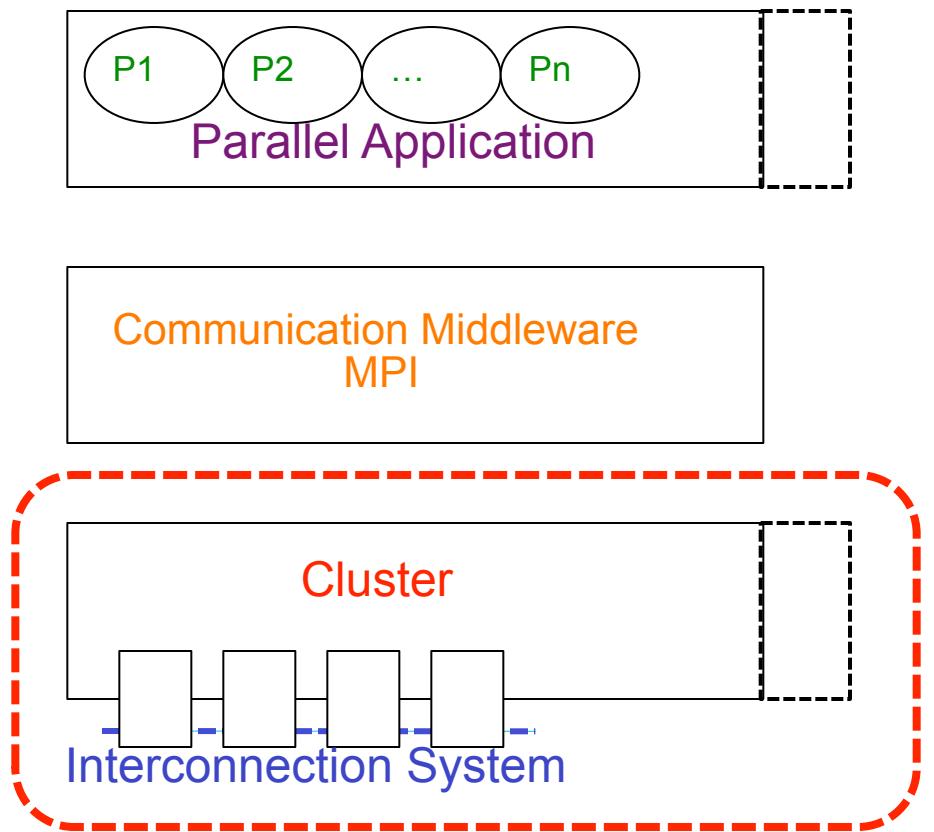
16



- Communication and I/O saturation:
 - Scalability problem
 - Performance problem

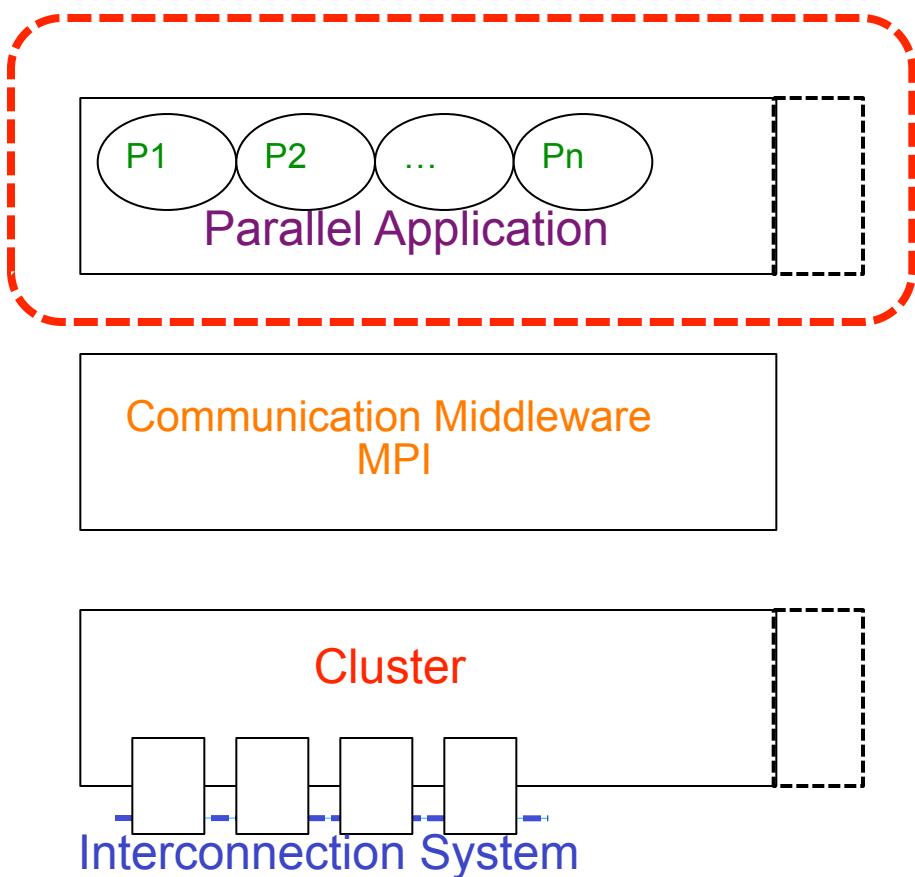
Problem description: possible solutions

17



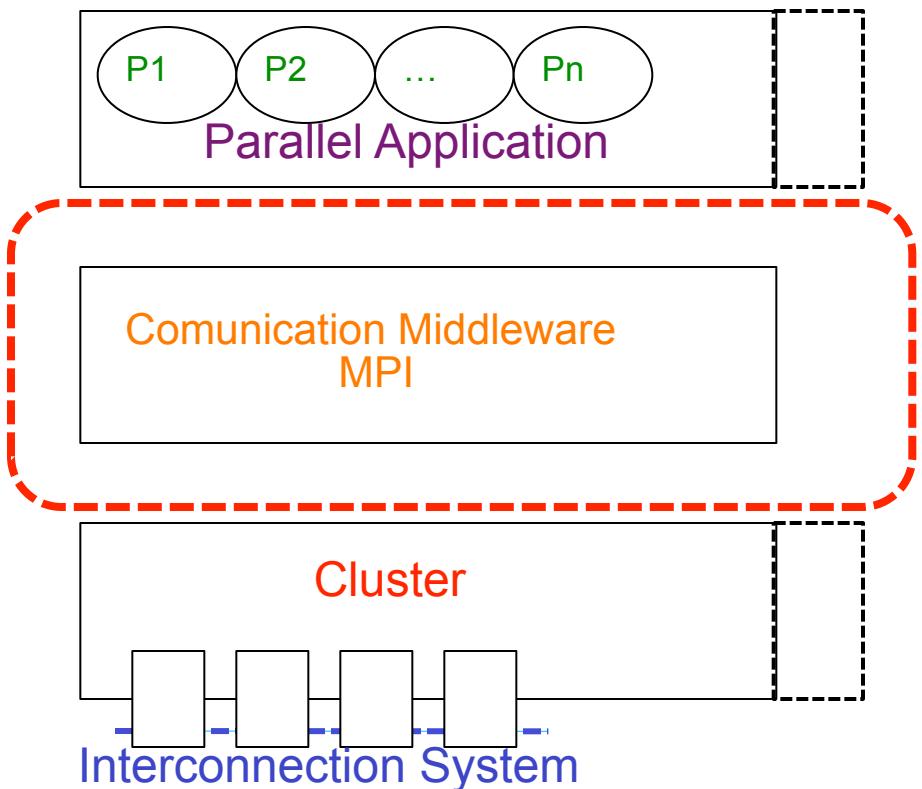
- 1) Improve network
 - Expensive.
 - Limited to current technology

Problem description: possible solutions



- 2) Improve the applications:
 - More effort in the design
 - Not always possible
 - The improvement affects few users

Problem description: possible solutions



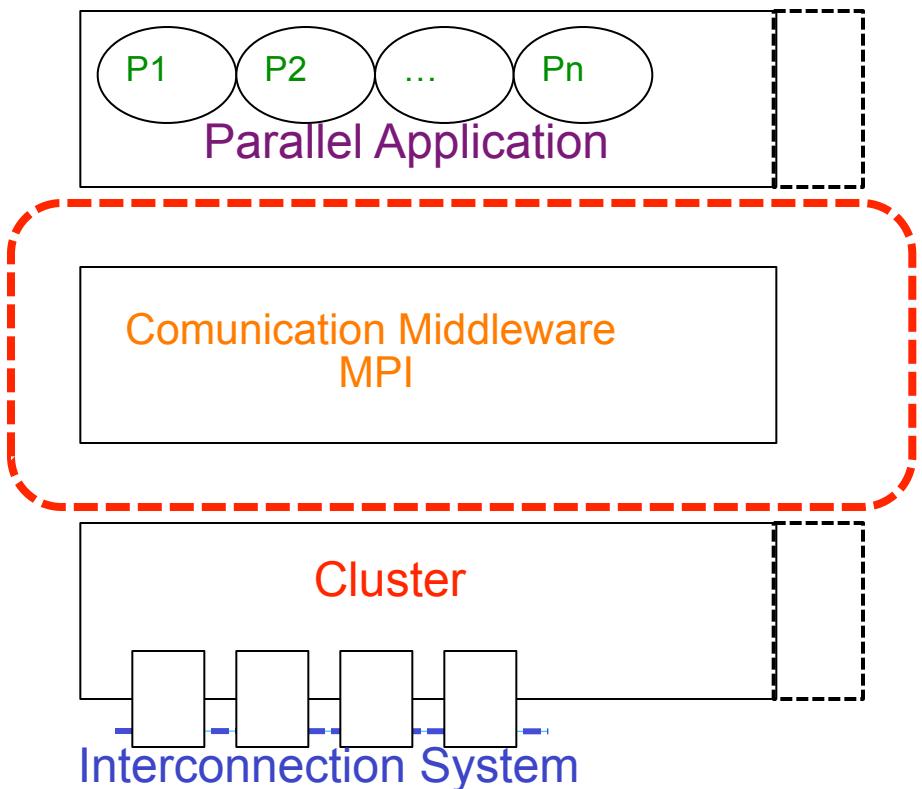
- 3) Improve the communication middleware
 - Portability
 - Greater user benefited
 - Lower Cost
 - Transparent:
 - Users
 - Applications

Summary

20

1. Problem description
2. **Main Objectives**
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance of I/O and communication operations.
5. Evaluation

Main objectives



- Improve the scalability and the performance of MPI based applications executed on Multicore cluster
- **How? Improving the Middleware MPI**

Specific objectives

22

1. Reduction of the number of communications in collective I/O operations.
2. Reduction of transferred data volume in communications.

Summary

23

1. Problem description
2. Main Objectives
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance of I/O and communication operations.
5. Evaluation tools

Strategies for improve the performance of collective I/O operation

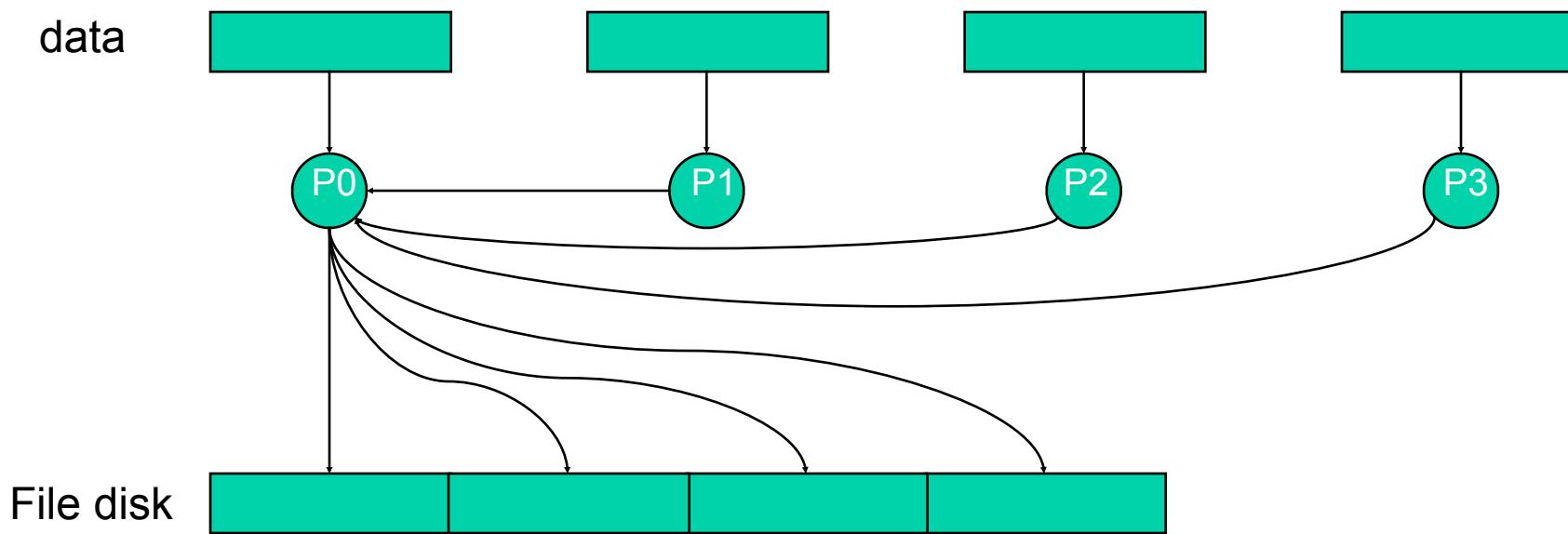
24

- Parallel scientific applications generate a lot of data to write/read in/from disk.
- Access pattern:
 - Sequential access.
 - Individual access.
 - Collective access

Sequential I/O

25

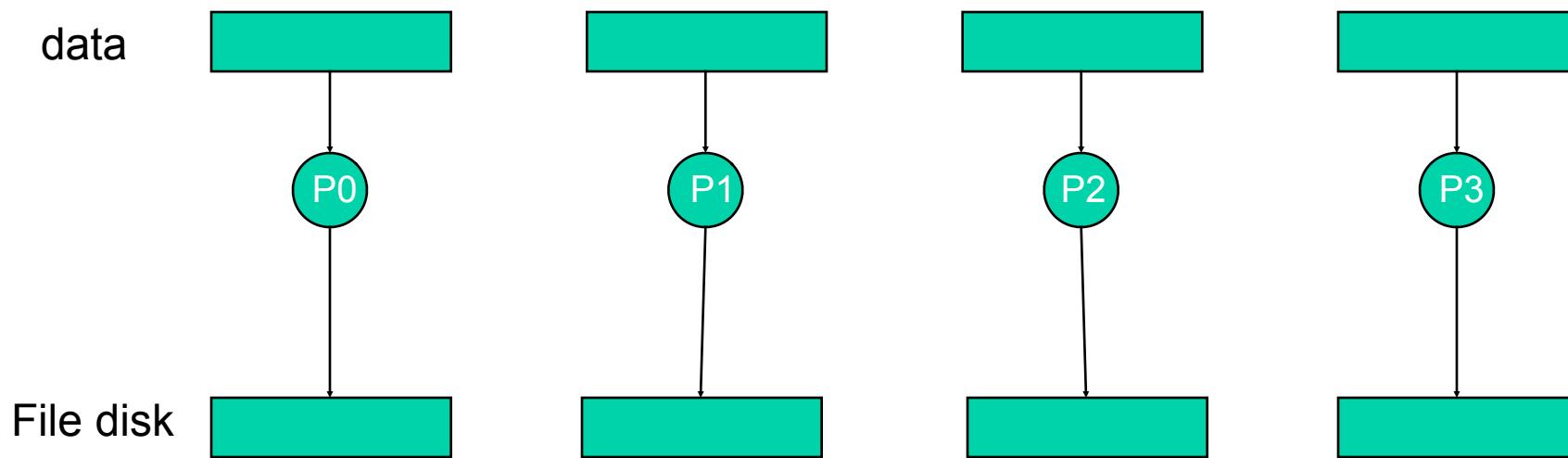
- All processes send data to one process (usually process 0), and this process writes it to the file



Individual I/O

26

- Each process writes to a separate file

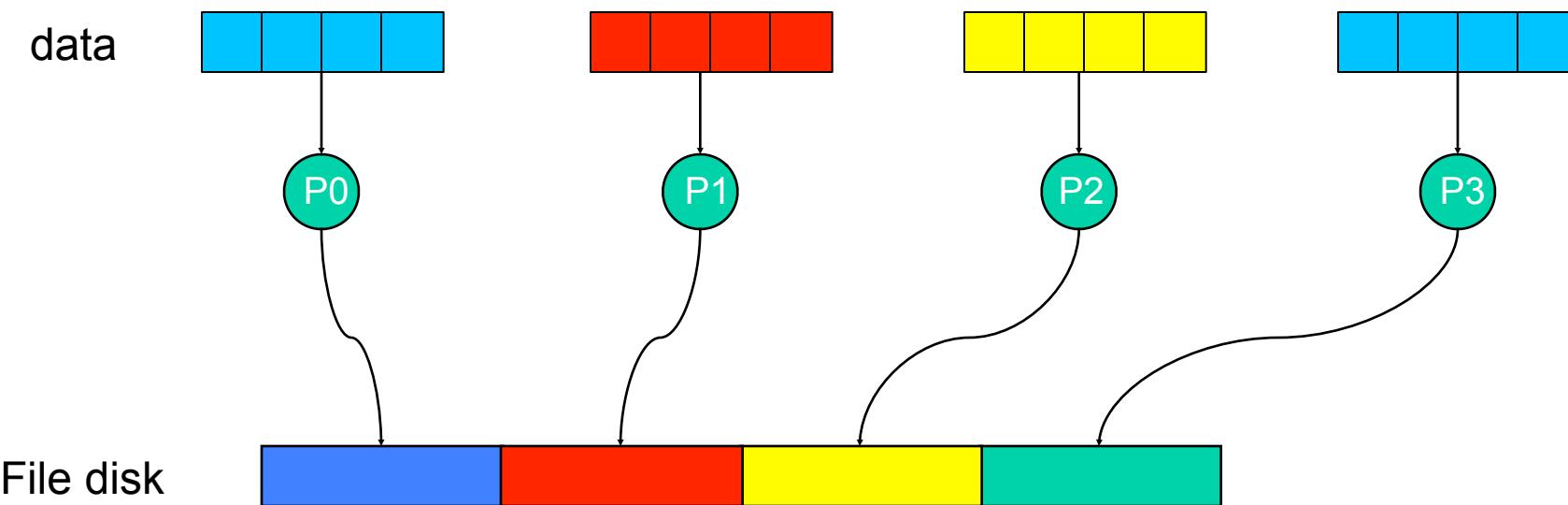


26

Collective I/O

27

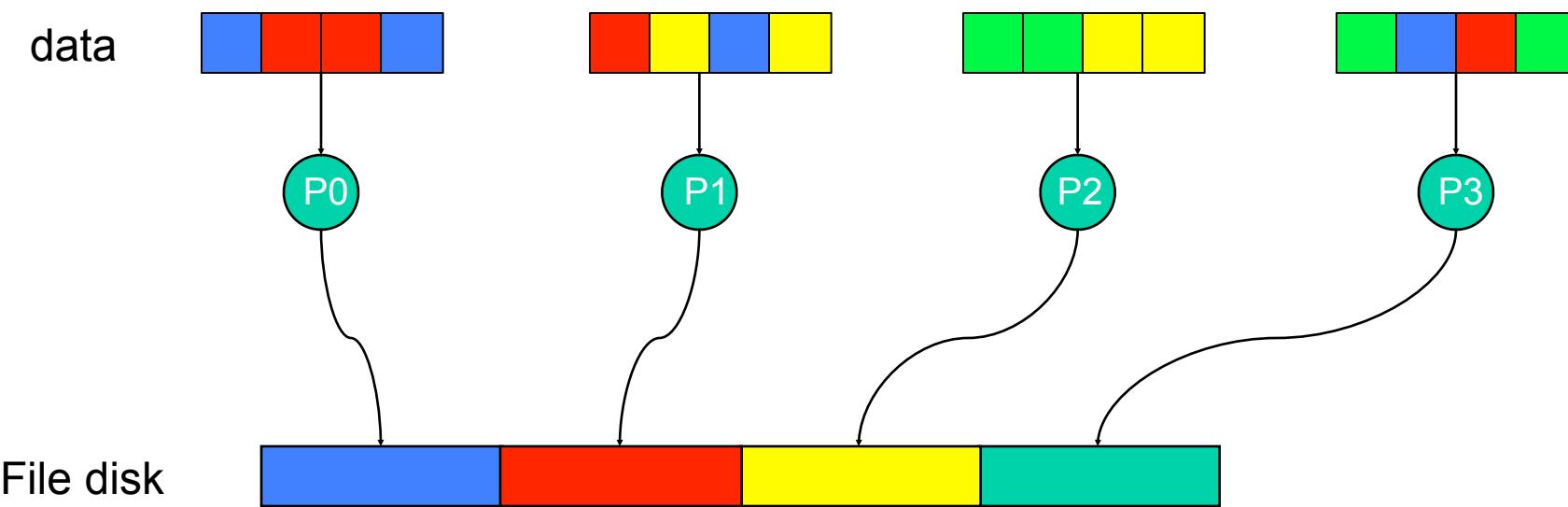
- Processes write to shared file
- Contiguous data
- Two Phase I/O technique.



Collective I/O

28

- Processes write to shared file
- Non-Contiguous data
- Two Phase I/O technique.



Two_Phase I/O

29

- Two-Phase I/O phases:
 - **Shuffle**: aggregate data into contiguous buffers.
 - **I/O**: transfer contiguous buffer to file system.
- Before these two phases:
 - File region is divided into equal contiguous regions
 - Called File Domains (FD).
 - Each FD is assigned to a subset of compute nodes (aggregators).
 - Each aggregator is responsible for transferring all data from its FD to the file system.
- Cause of inefficiency: The assignment of FD to aggregators is independent of data distribution.

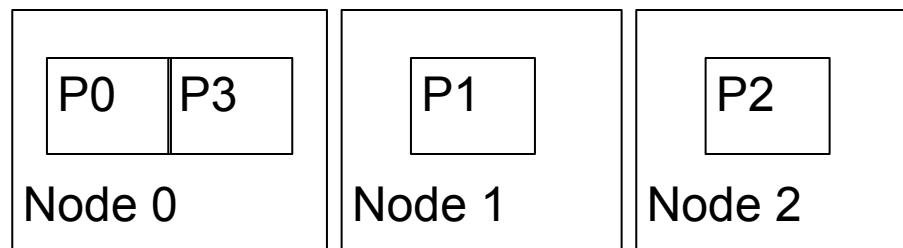
Two_Phase I/O “problem”

30

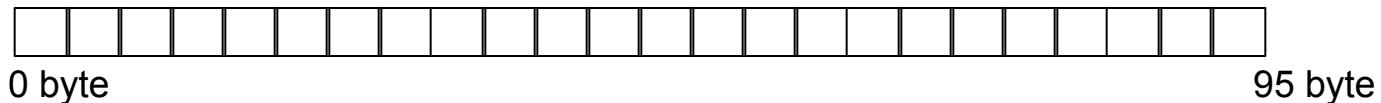
- The assignment of FD to aggregators is fixed.
- Independent of data distribution.
- Default aggregator Pattern:
 - So many aggregator as nodes.
 - Create an array of nodes, ordered by rank (process identifier in MPI), and assign each aggregator one process.
 - When there are more than process per node:
 - The process with the smallest rank is chosen as aggregator.

Example of Two-Phase I/O

31

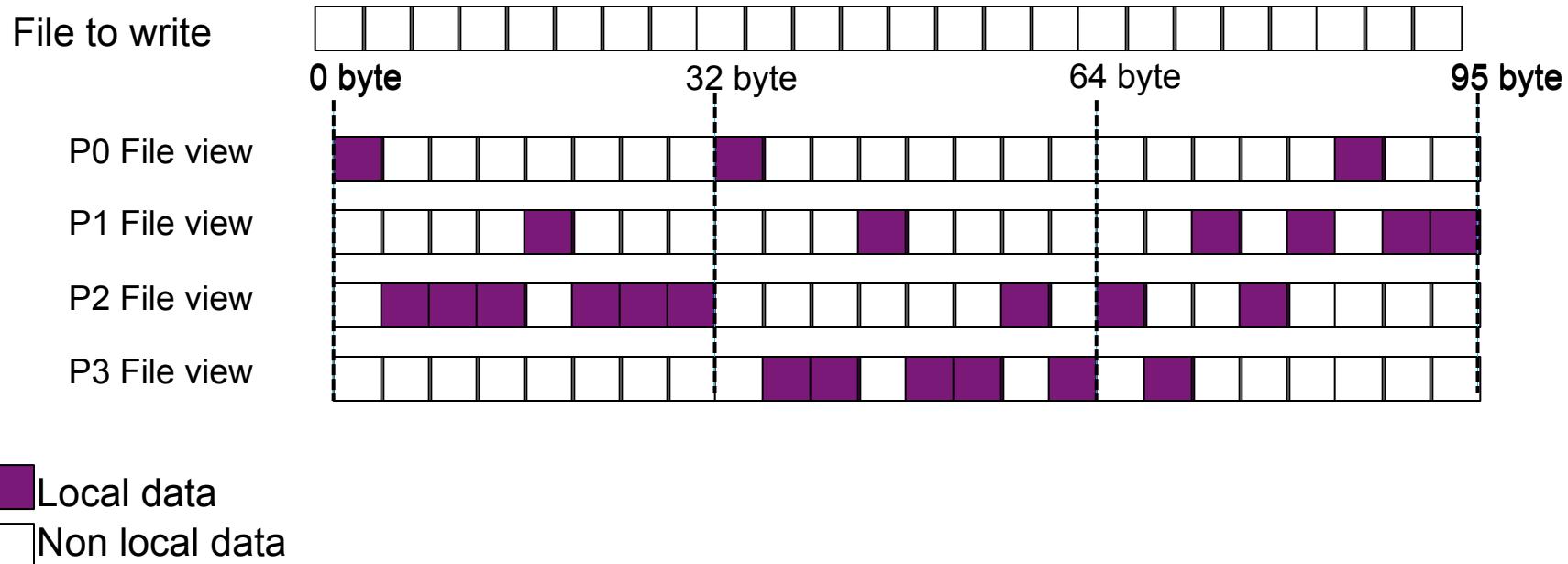
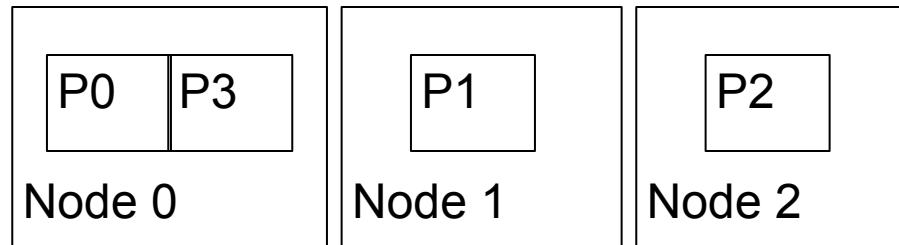


File to write



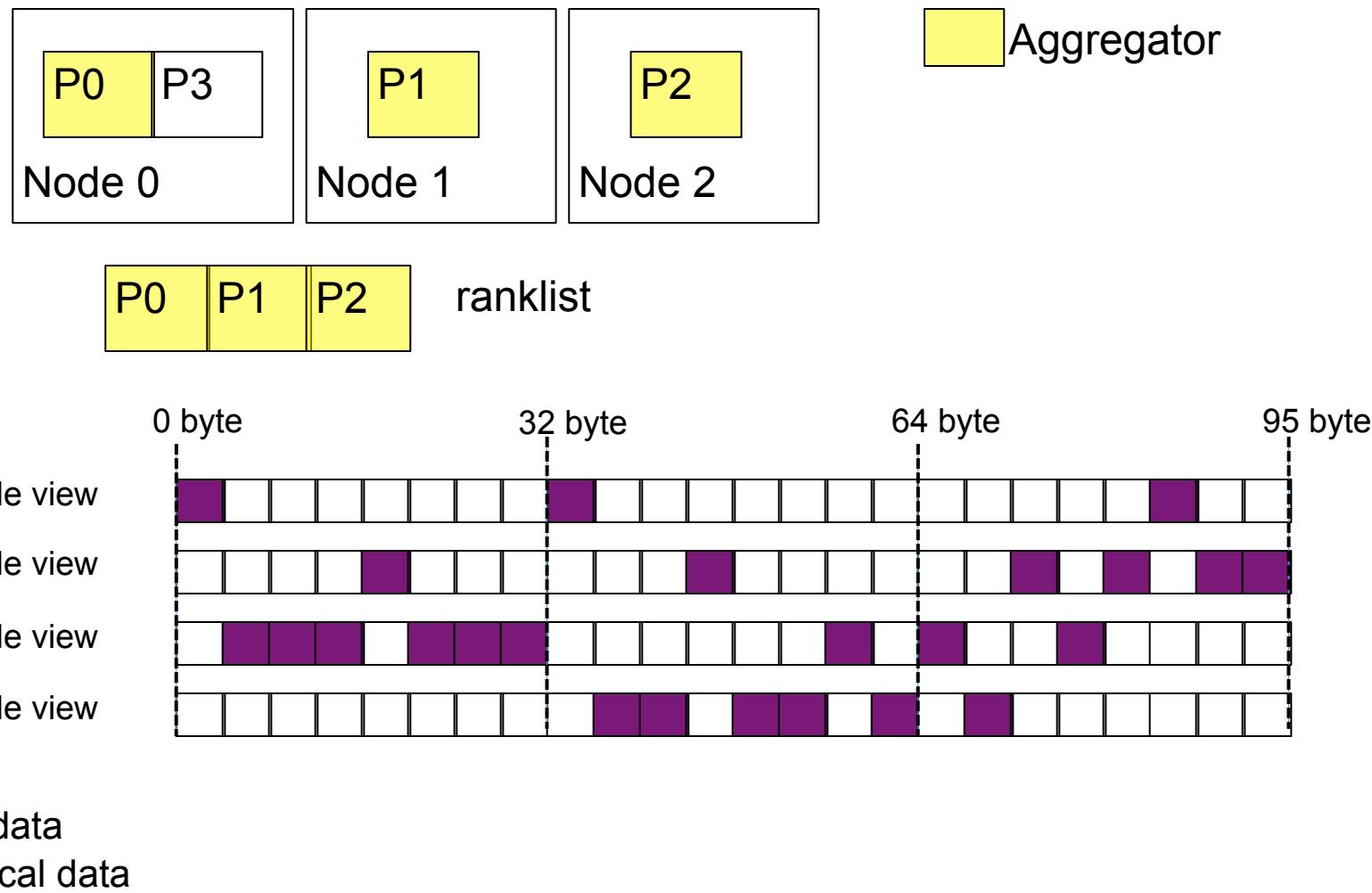
Example of Two-Phase I/O

32



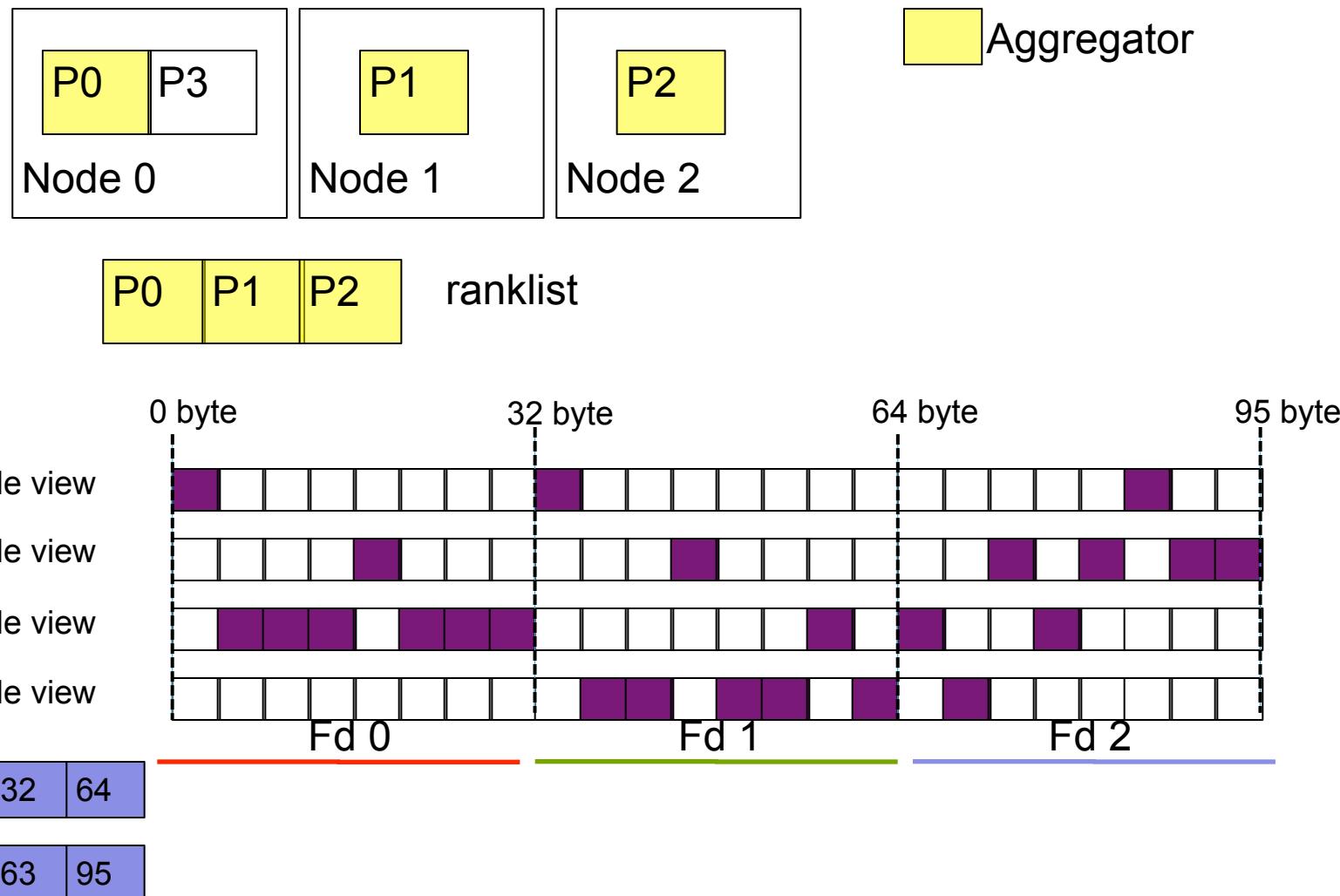
Example of Two-Phase I/O

33



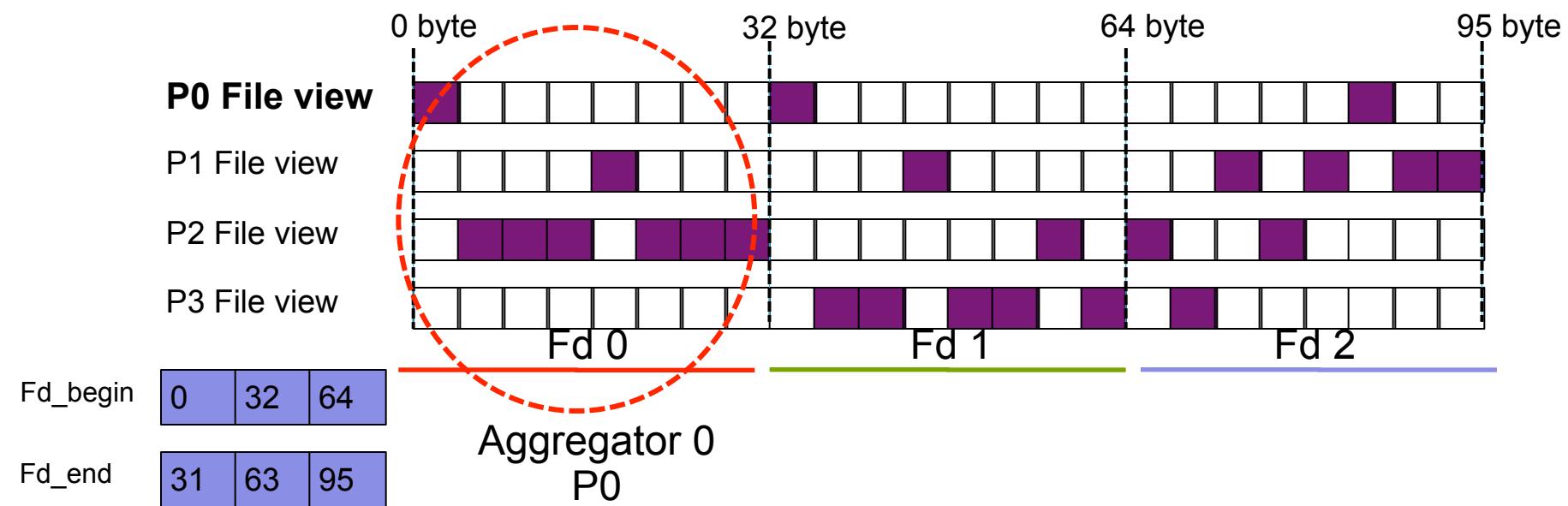
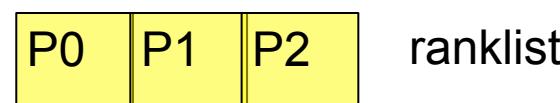
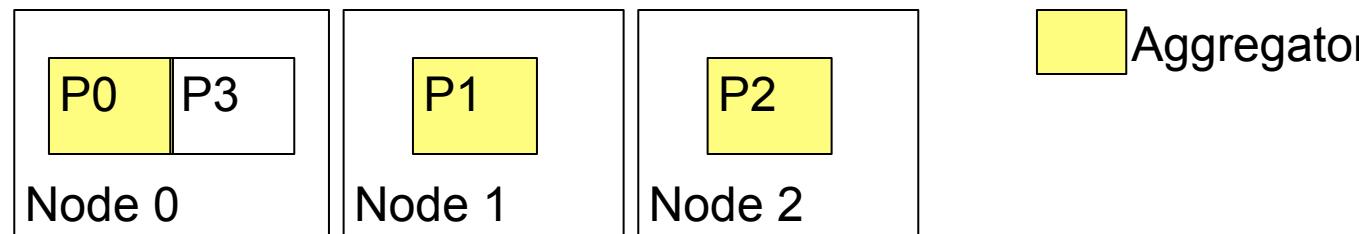
Example of Two-Phase I/O

34



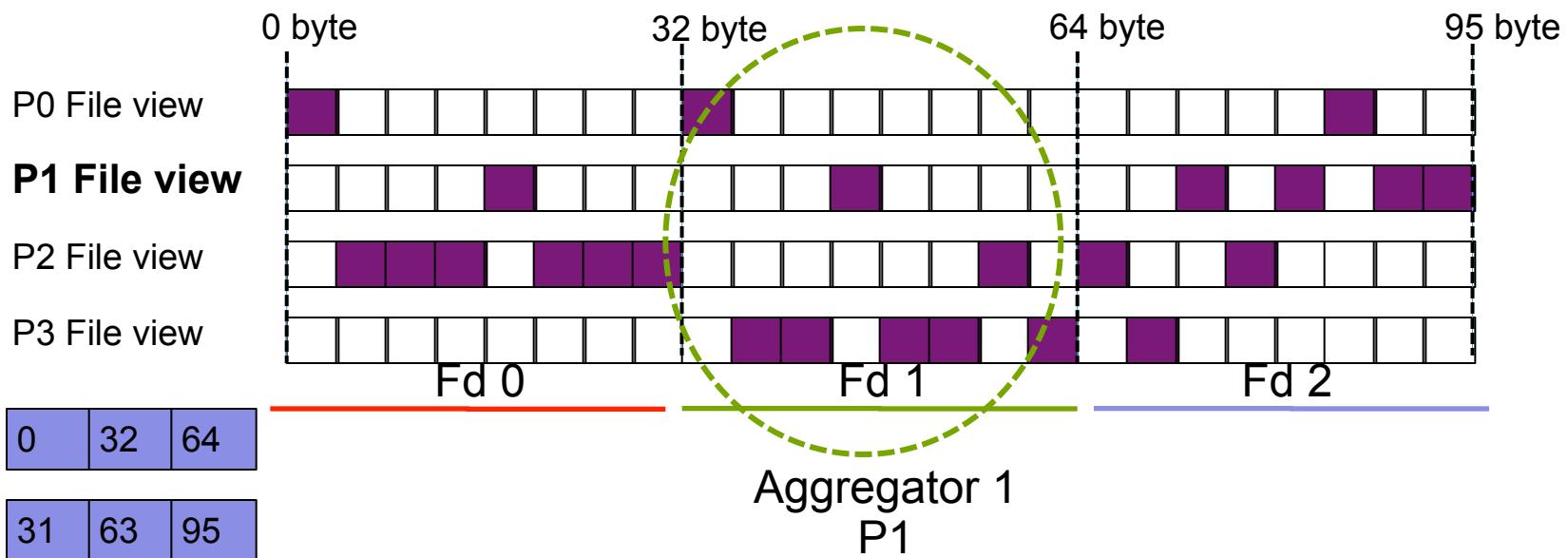
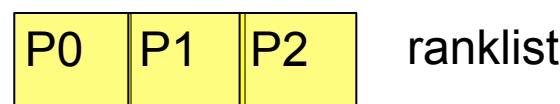
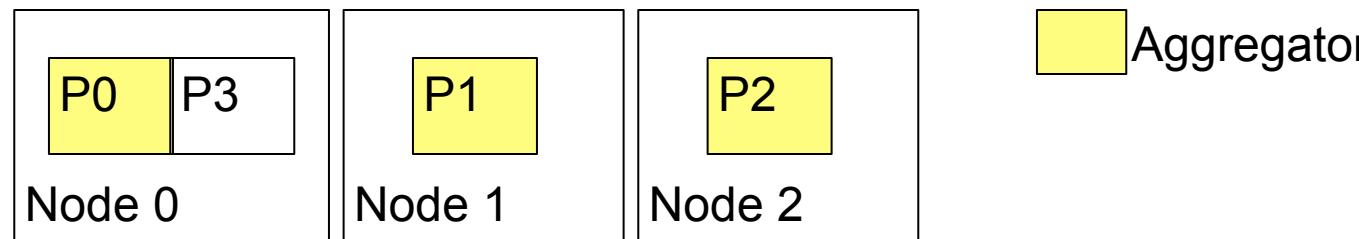
Example of Two-Phase I/O

35



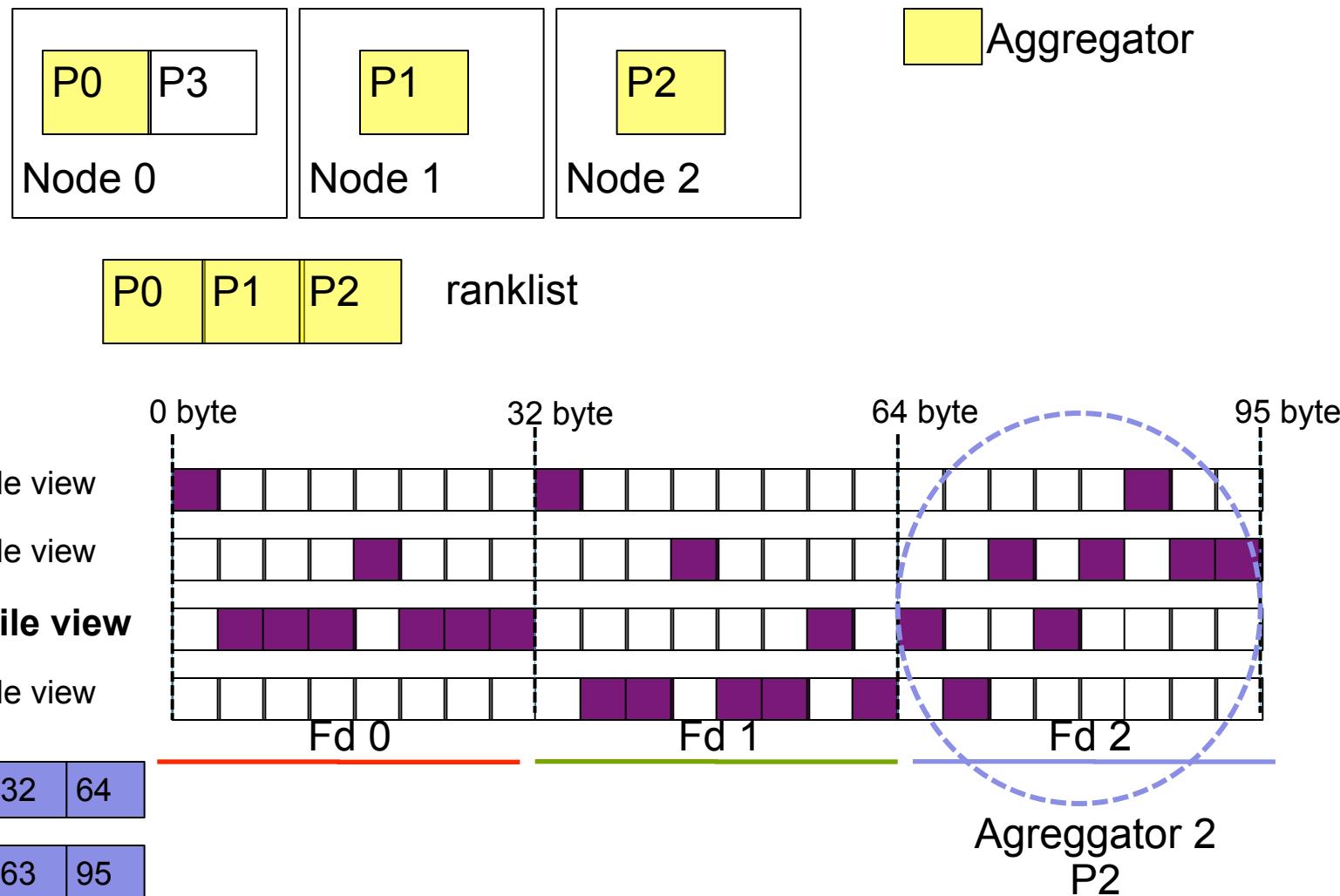
Example of Two-Phase I/O

36



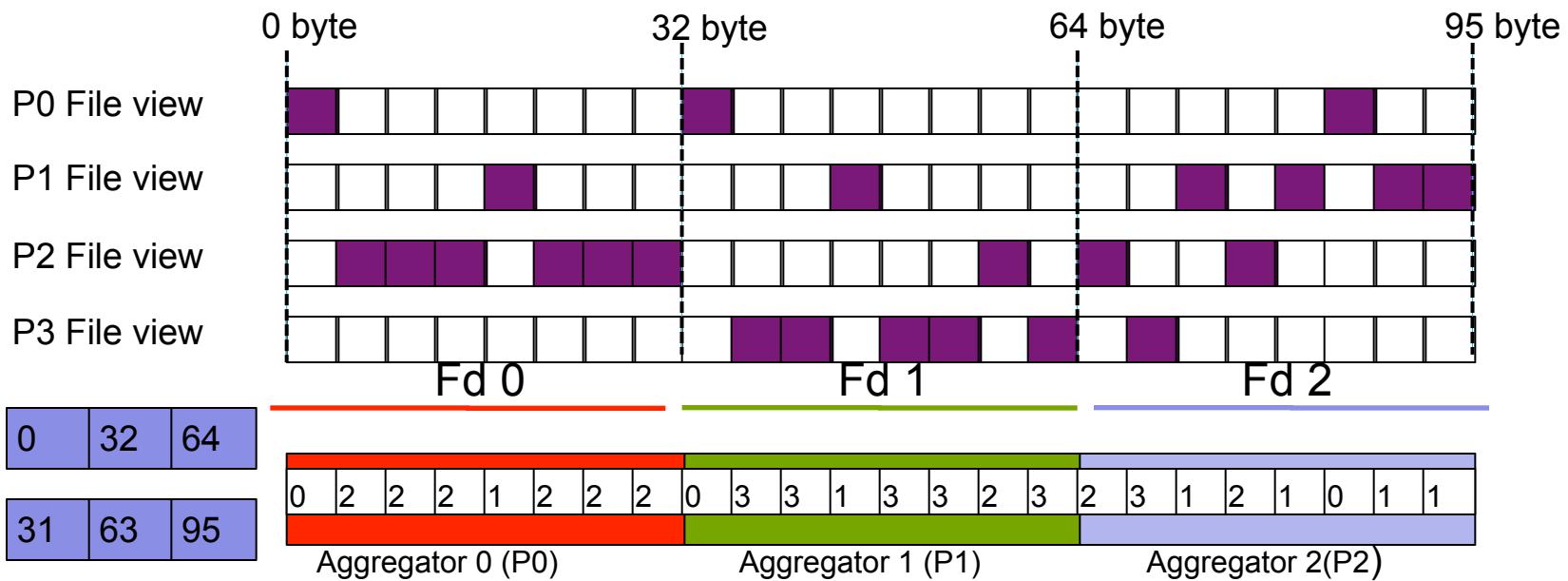
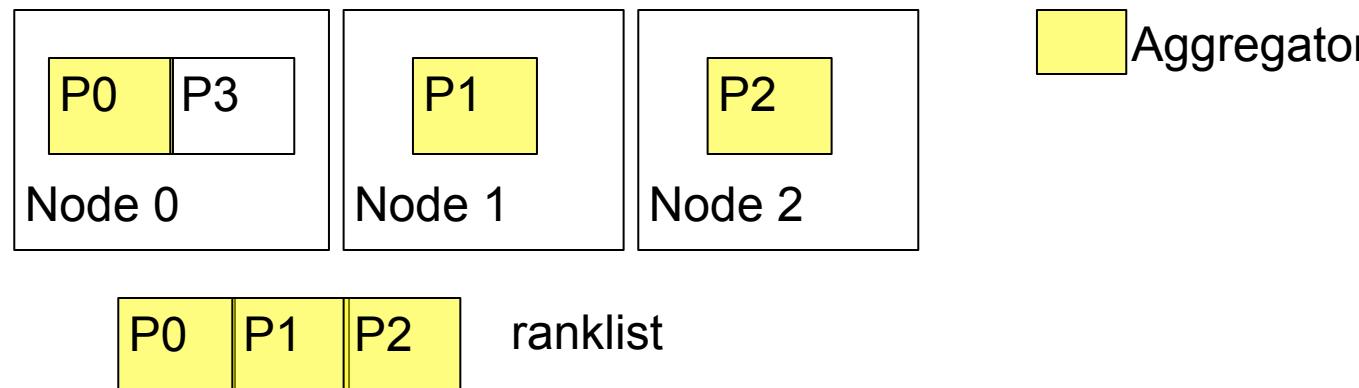
Example of Two-Phase I/O

37



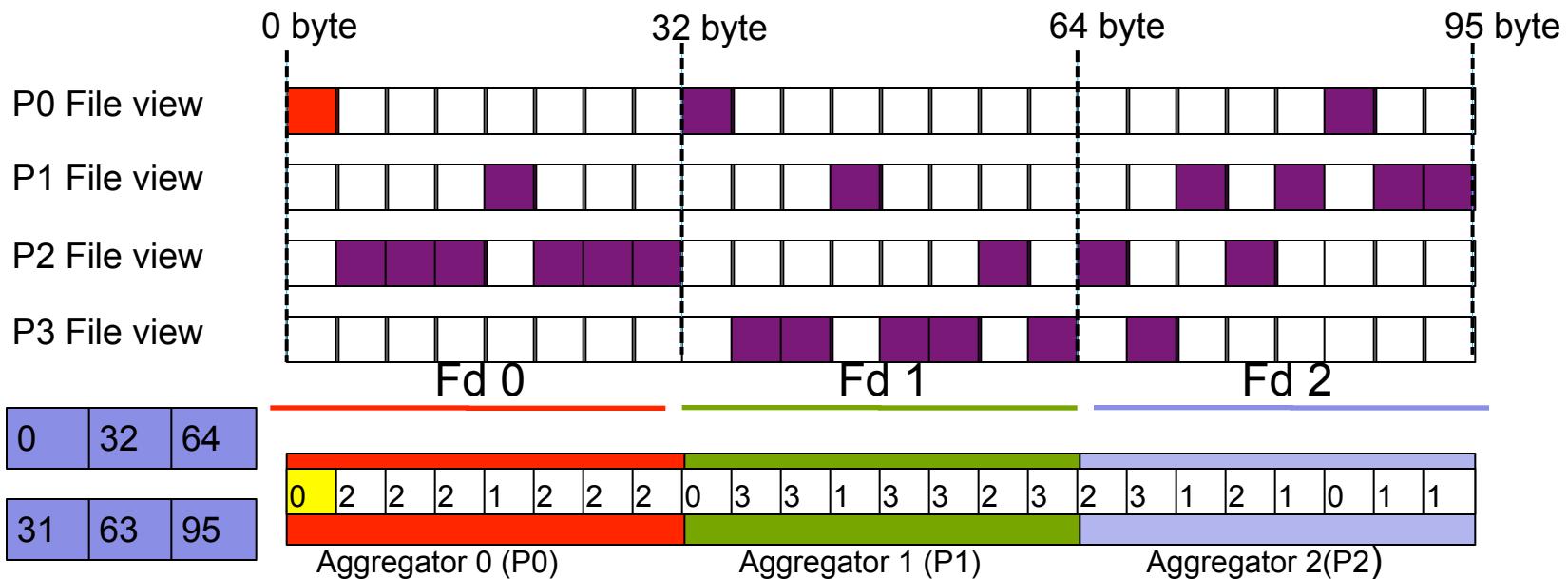
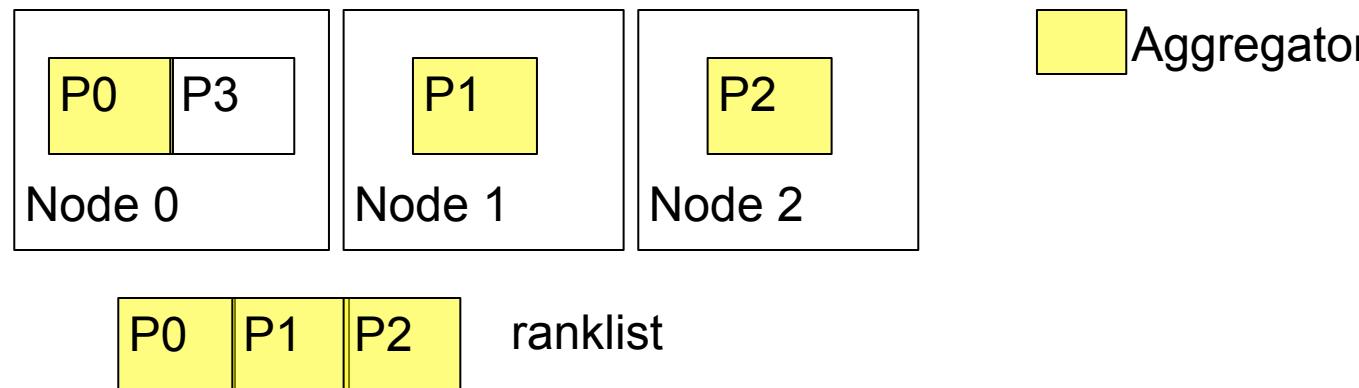
Example of Two-Phase I/O

38



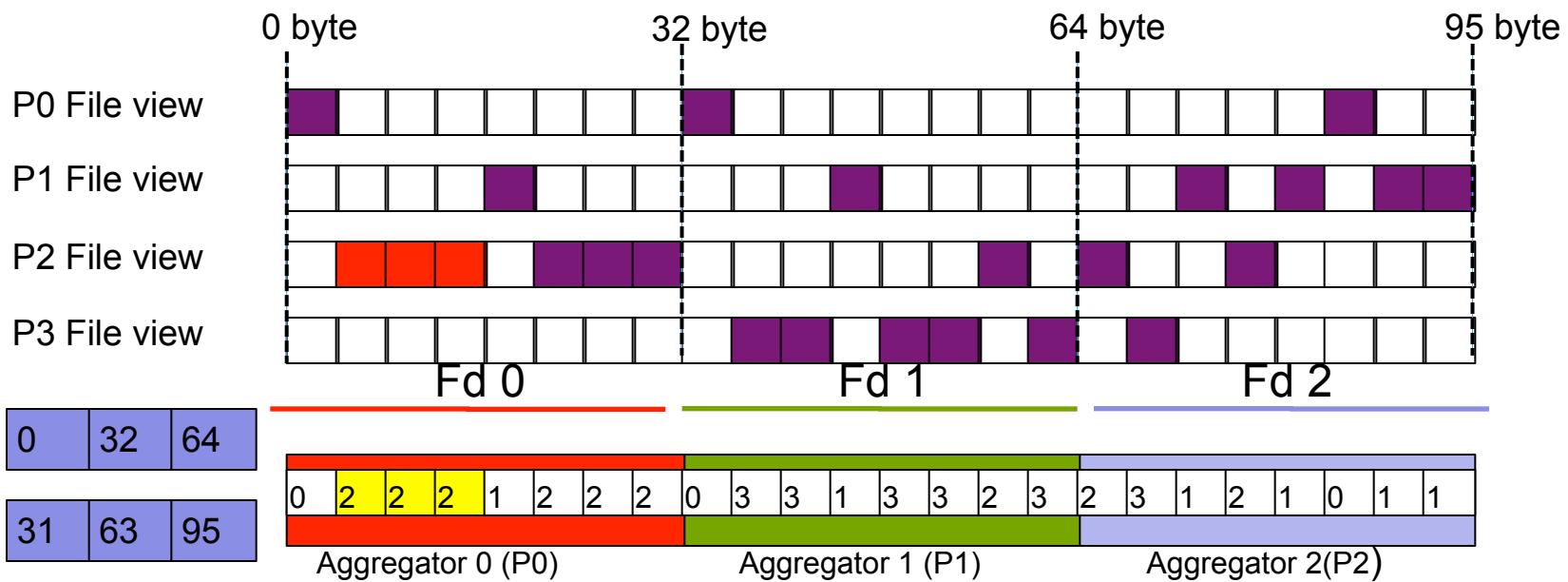
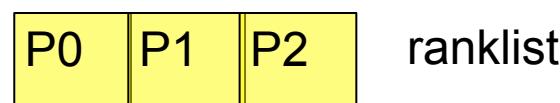
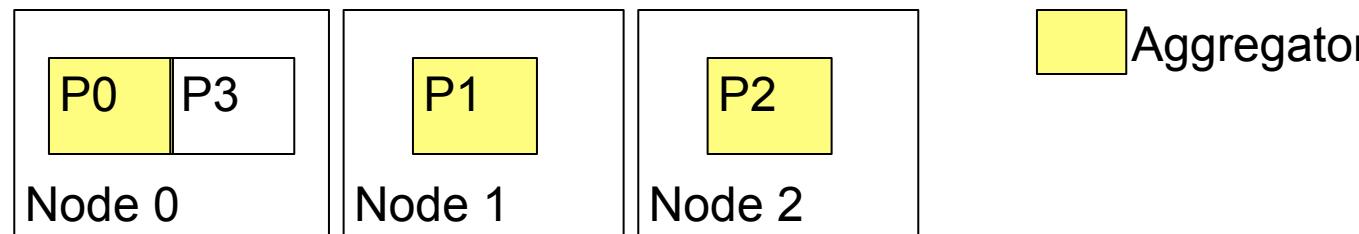
Example of Two-Phase I/O

39



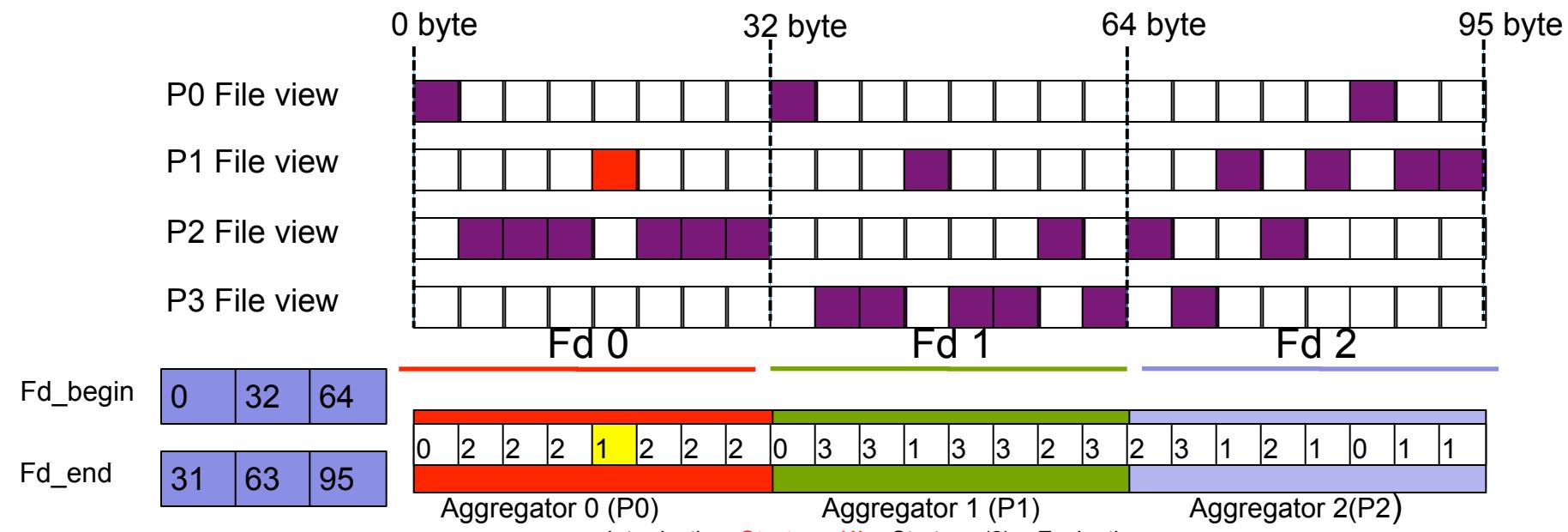
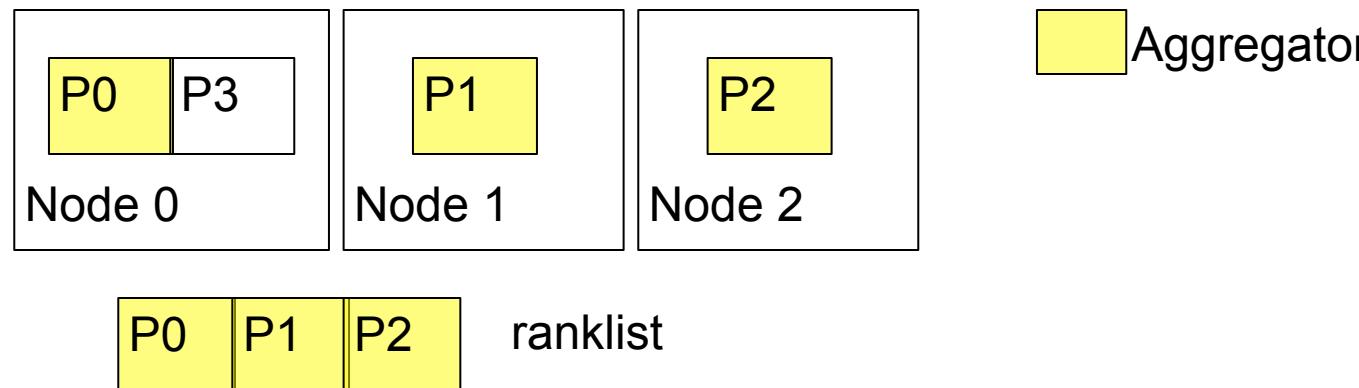
Example of Two-Phase I/O

40



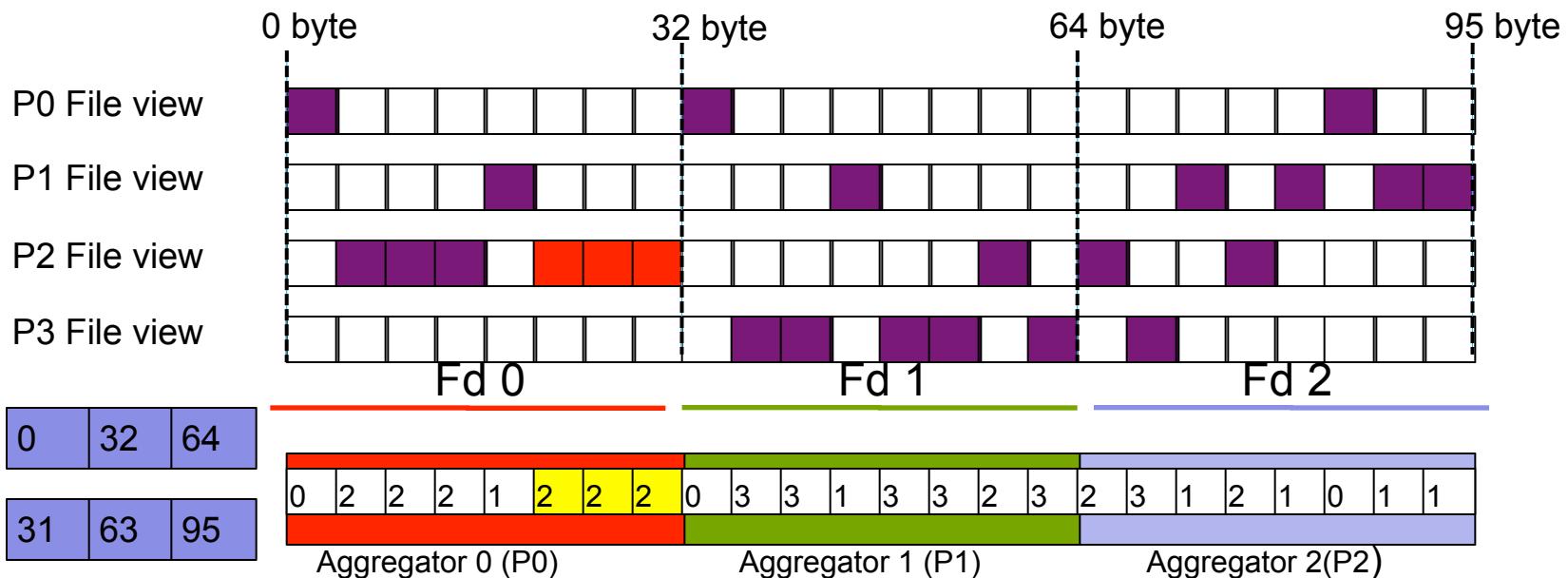
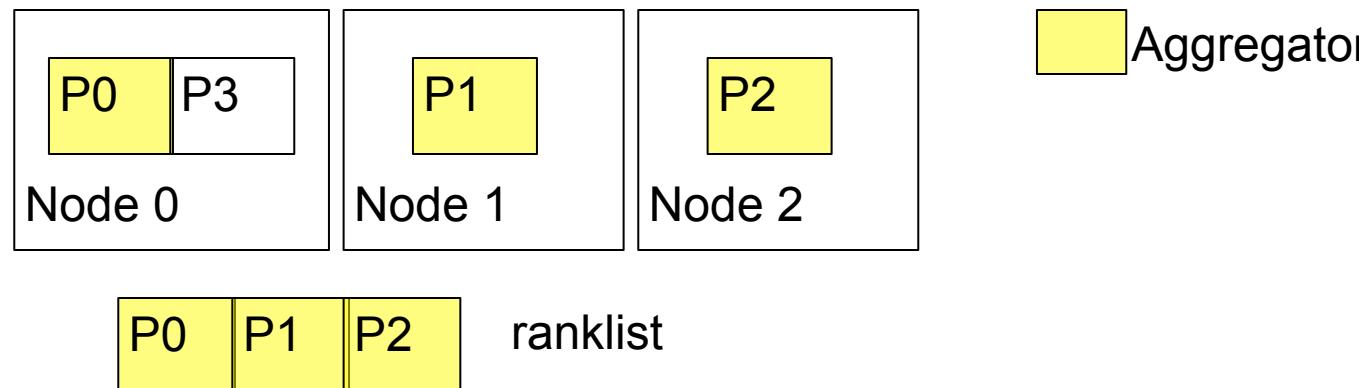
Example of Two-Phase I/O

41



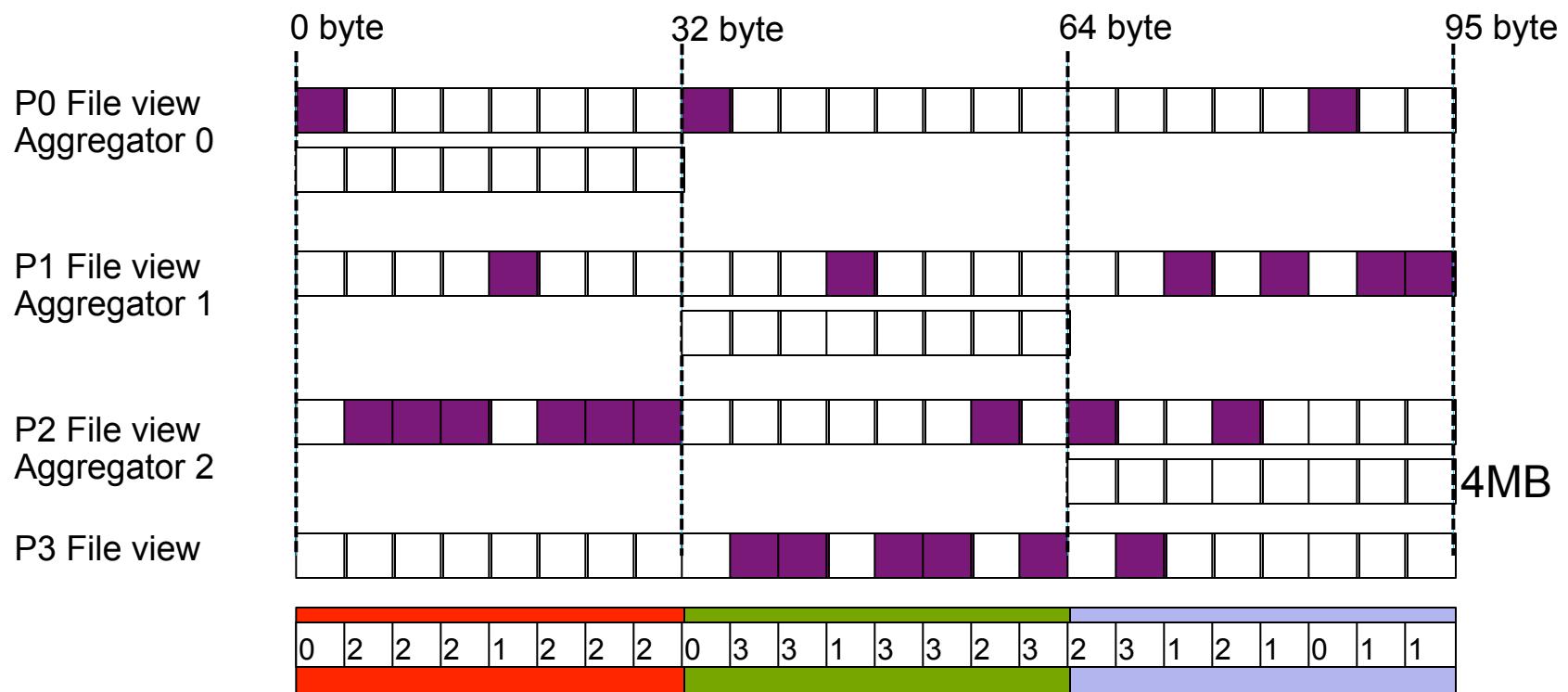
Example of Two-Phase I/O

42



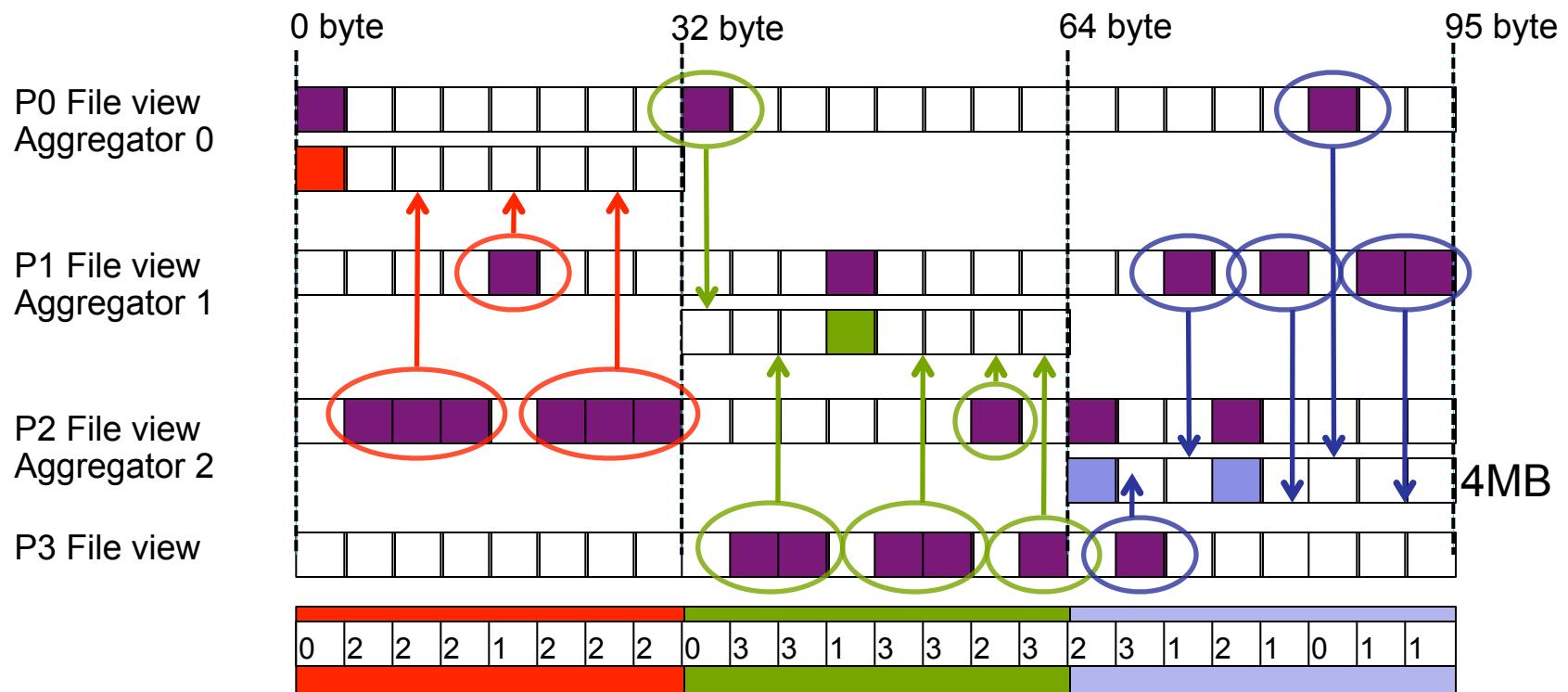
Example of Two-Phase I/O

43



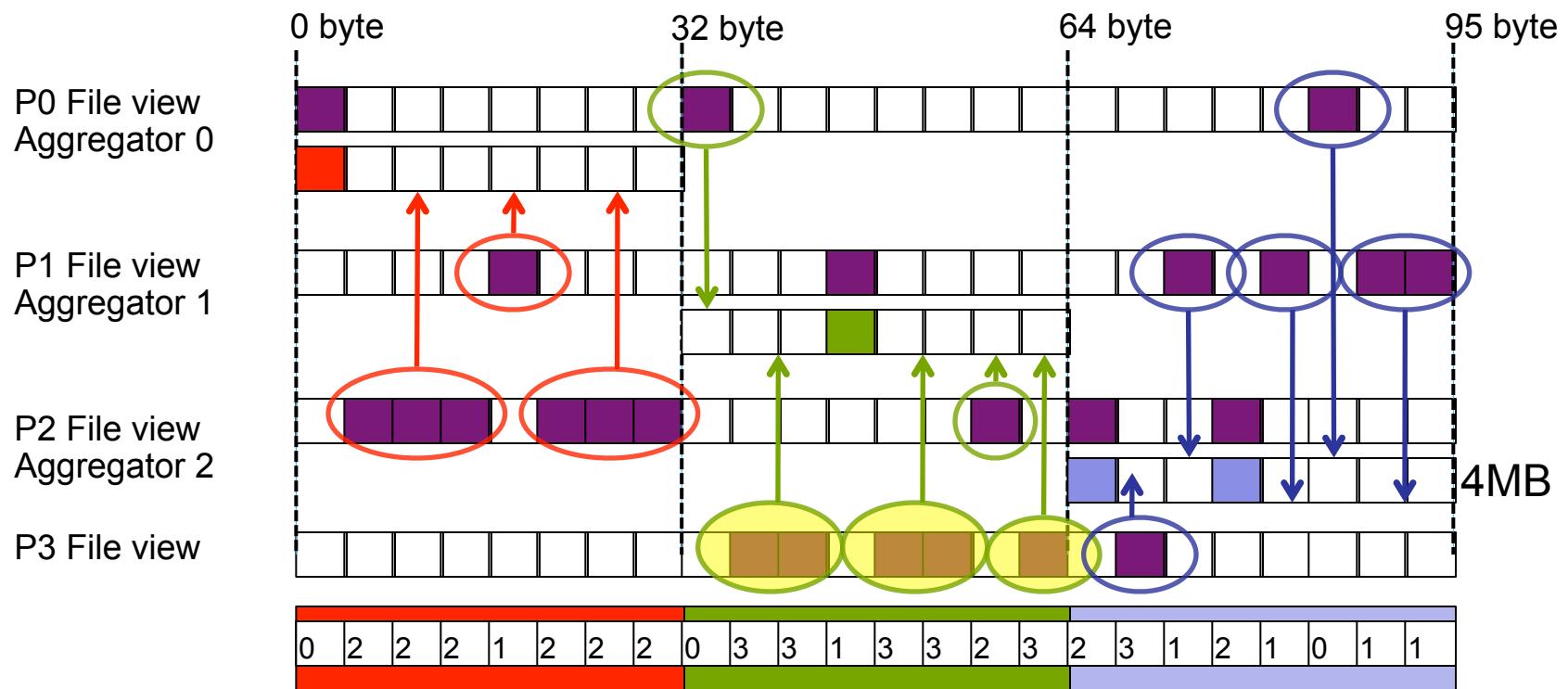
Example of Two-Phase I/O

44



Example of Two-Phase I/O

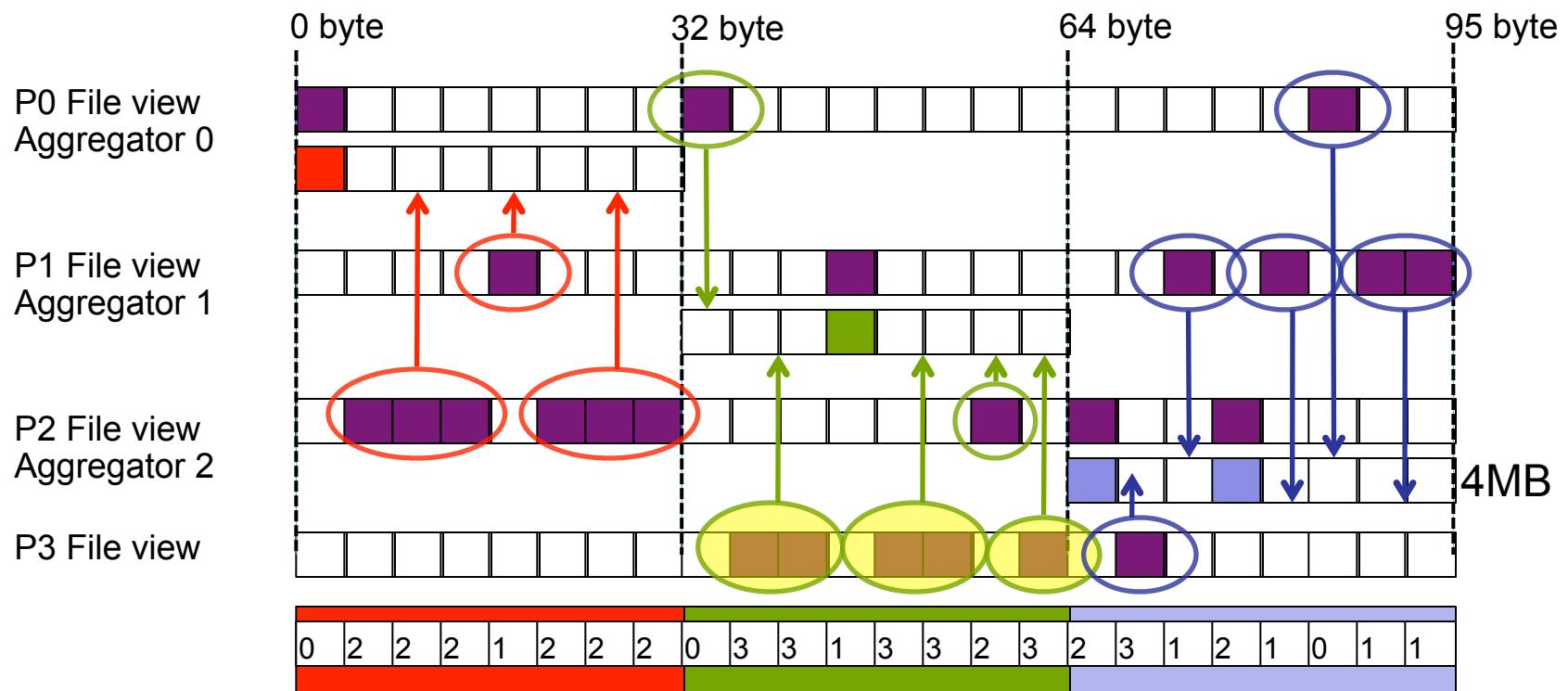
45



Communication: 13 message
76 bytes of 96 bytes are transferred among the processes.

Example of Two-Phase I/O

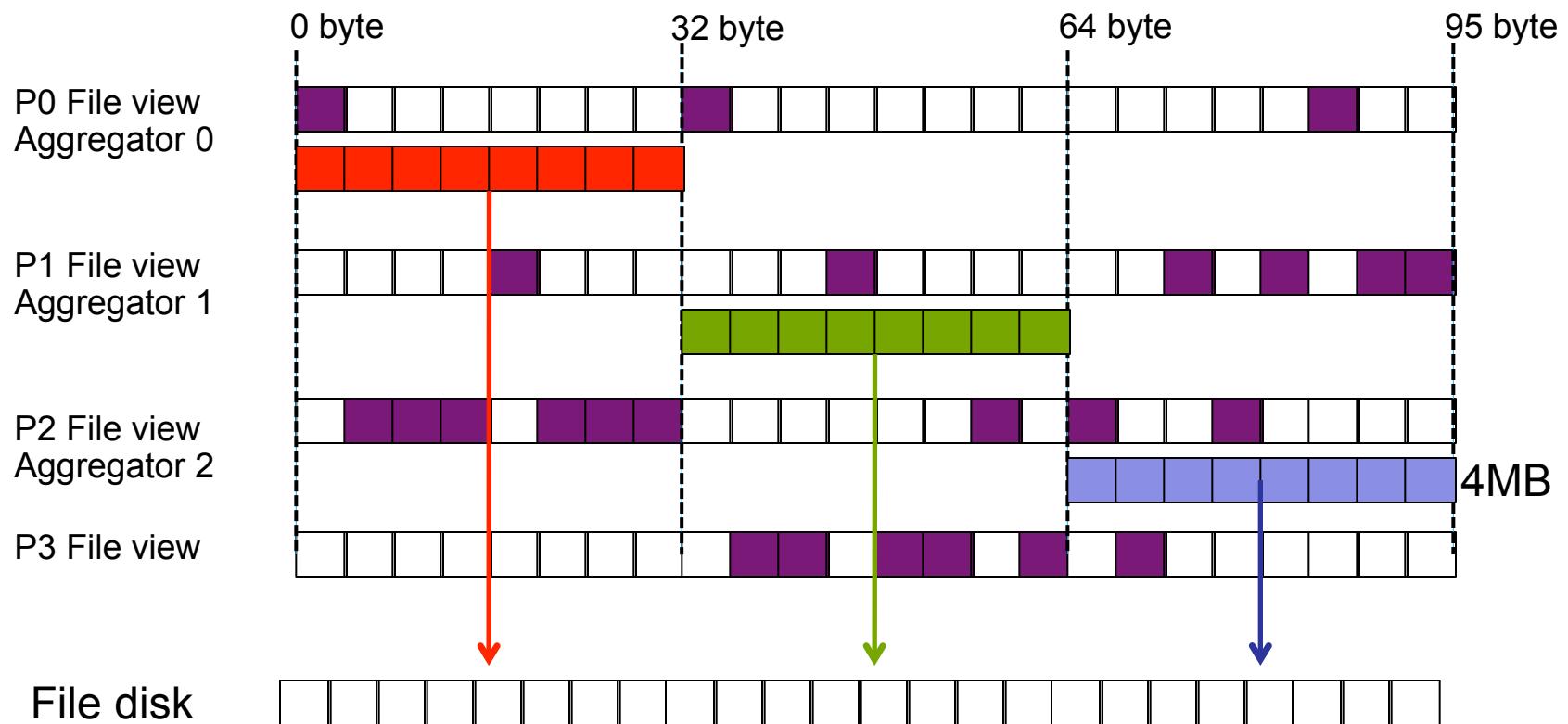
46



The number of communications could be reduced if each aggregator is assigned to the process more adequately.

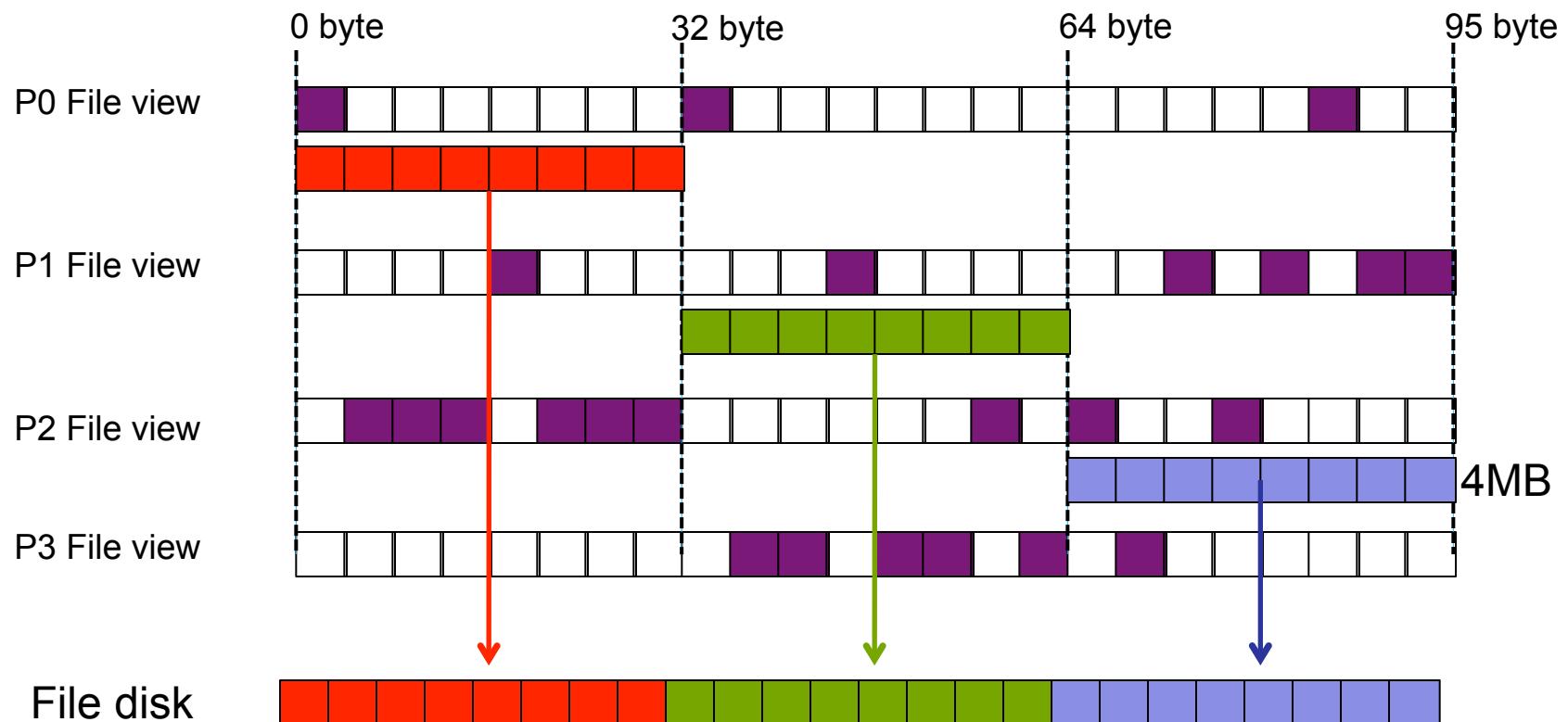
Example of Two-Phase I/O

47



Example of Two-Phase I/O

48



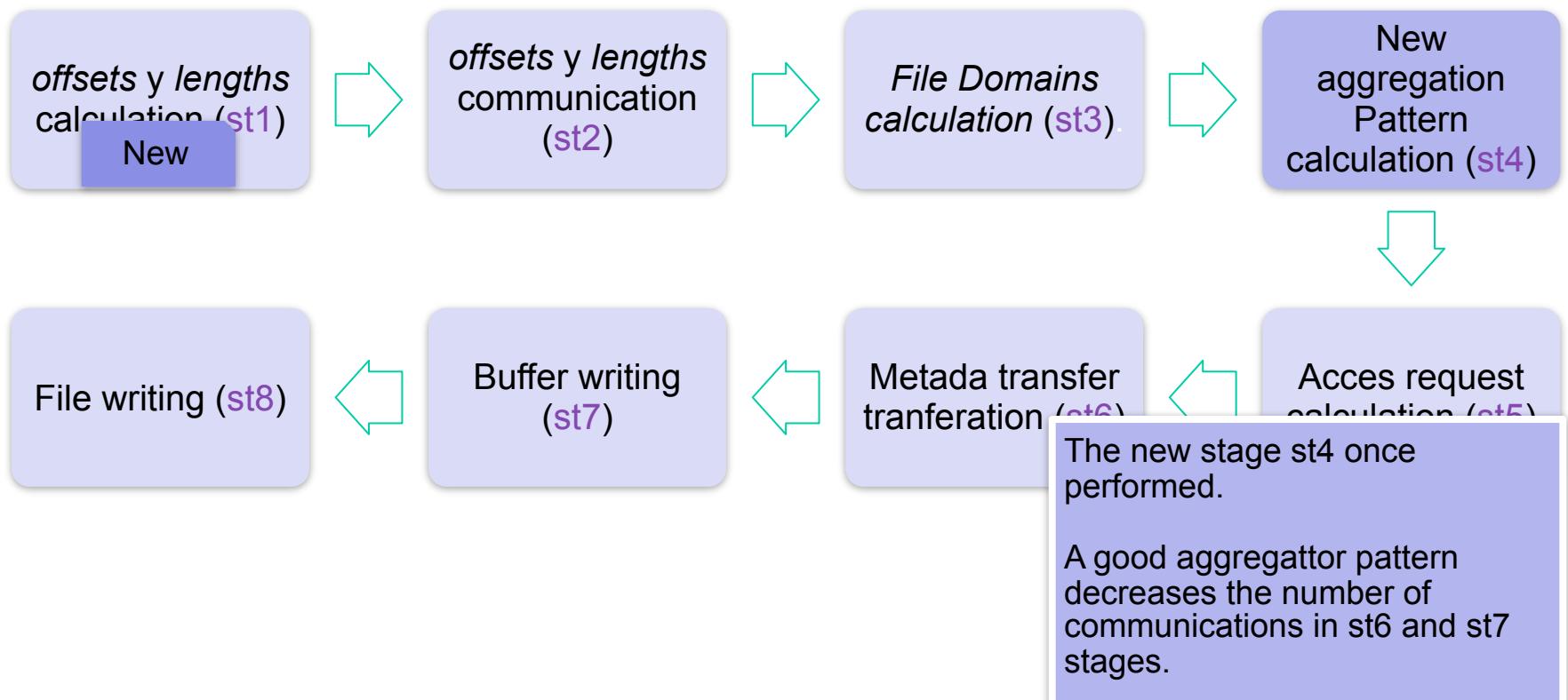
Phd. Thesis Proposal

49

- Replace the rigid assignment of aggregator by new one based on the next aggregation-criteria:
 - Reduce the number of communications:
 - Each aggregator is assigned to the process who has the highest number of contiguous data blocks.
 - Uses the initial distribution of data over the processes.
 - Employs the Linear Assignment Problem (LAP)
- The implementation of the **new dynamic** aggregator pattern → new version of Two_Phase I/O:
 - Locality_Aware_Two_Phase I/O (LA_TwoPhase I/O)

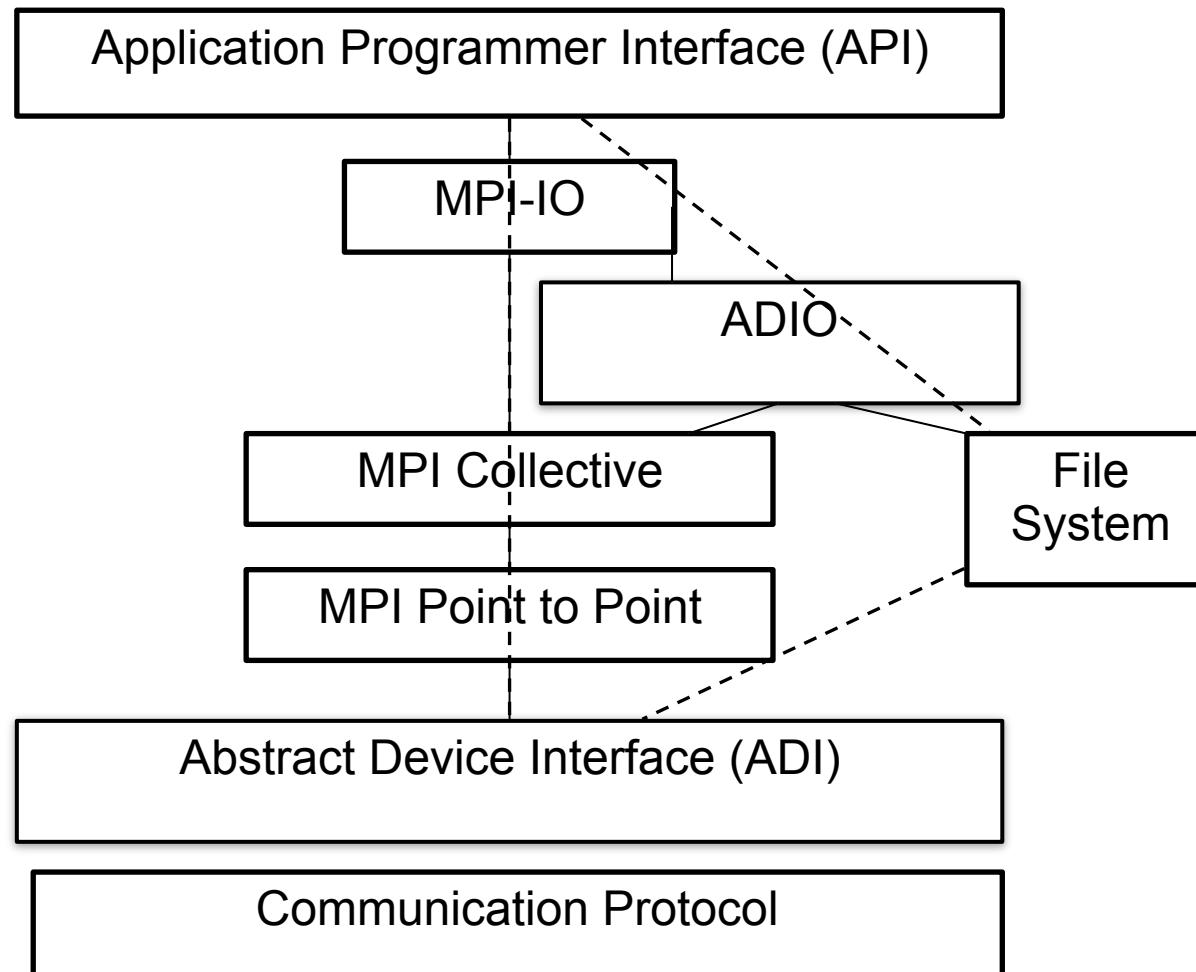
LA_Two_Phase I/O

50



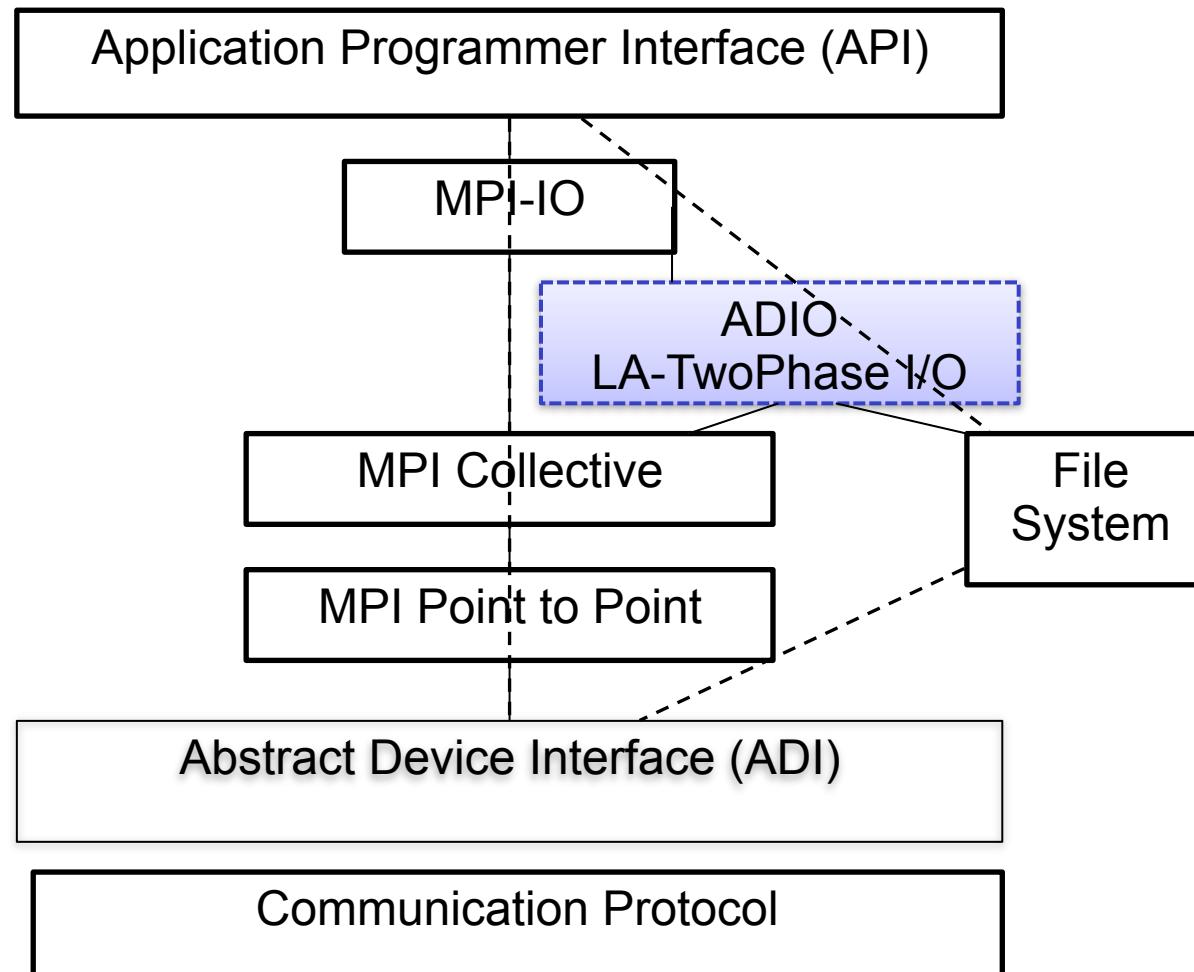
MPICH1.2 Architecture

51

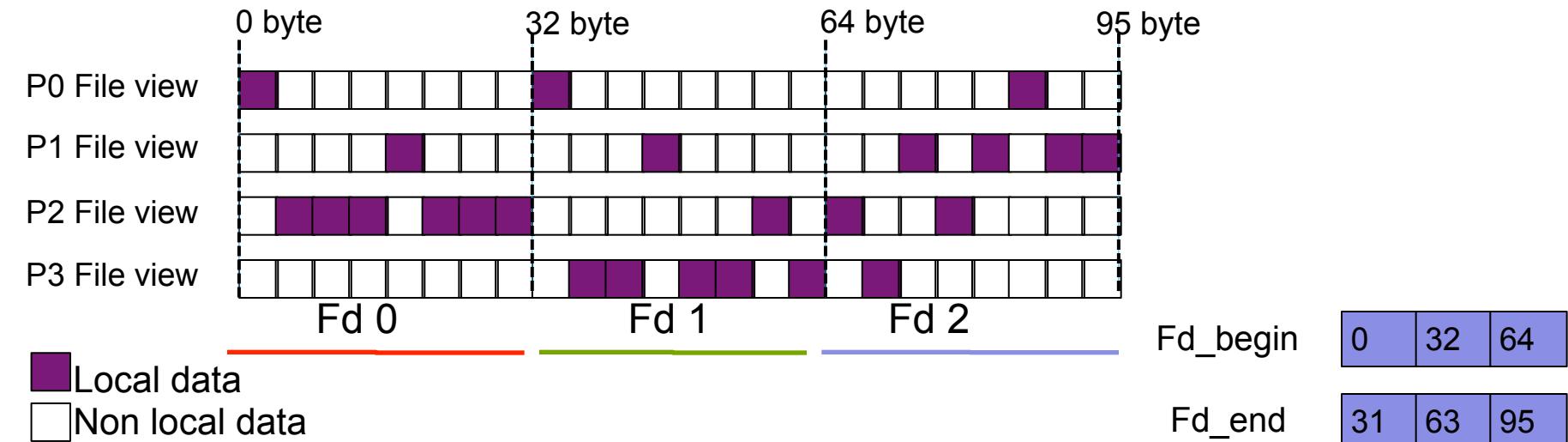


MPICH1.2 Modification

52

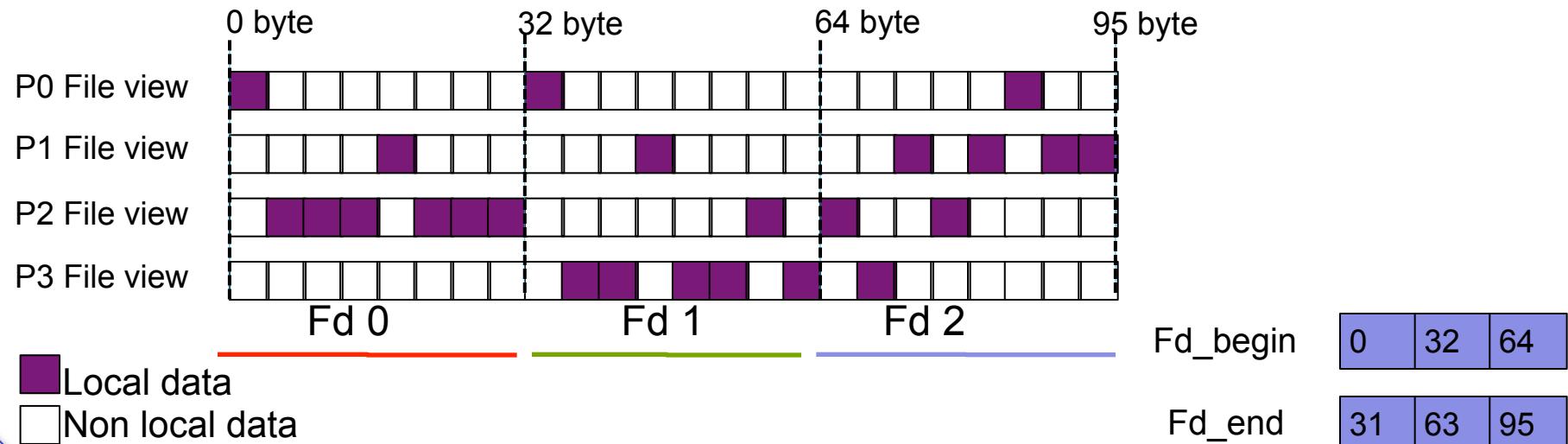


Example of LA_Two_Phase I/O

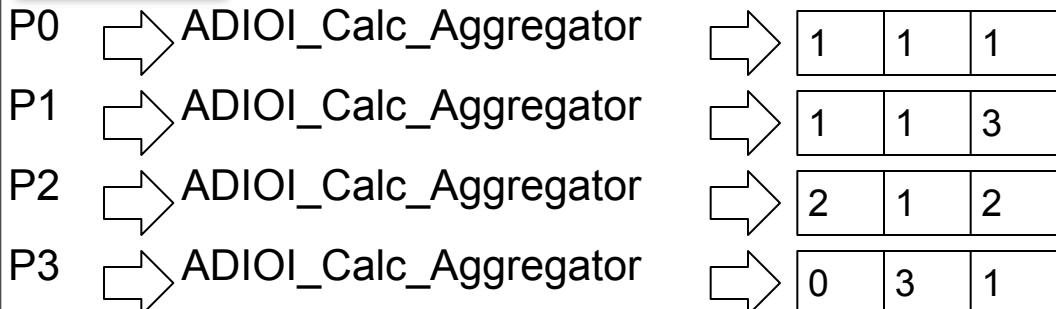


Calculation and recollection of distribution data

54

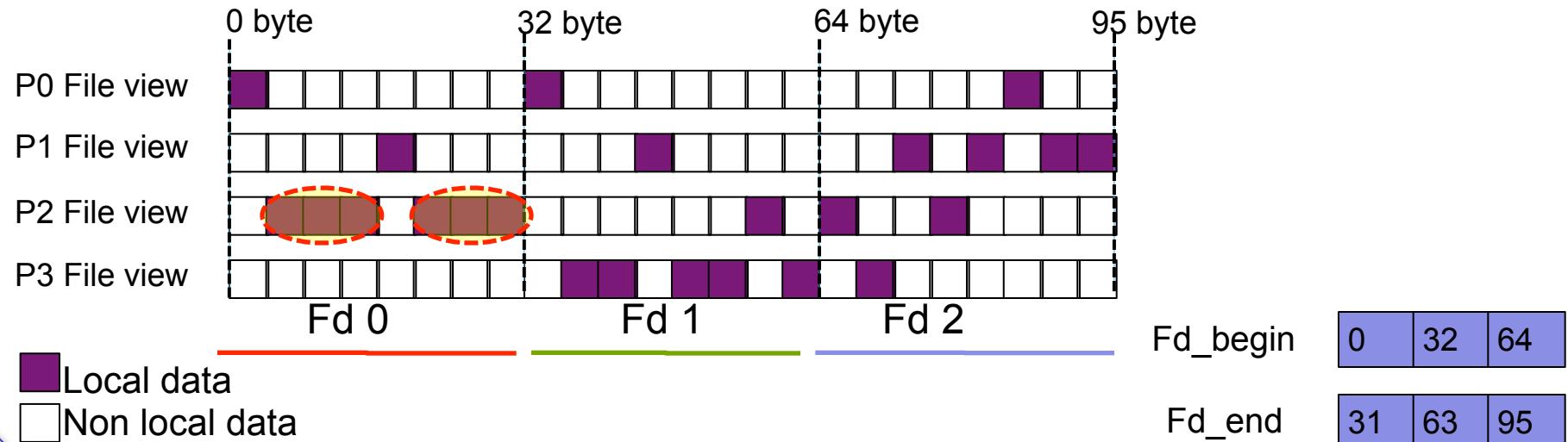


Stage:
St4.1



Calculation and recollection of distribution data

55



Stage:
St4.1

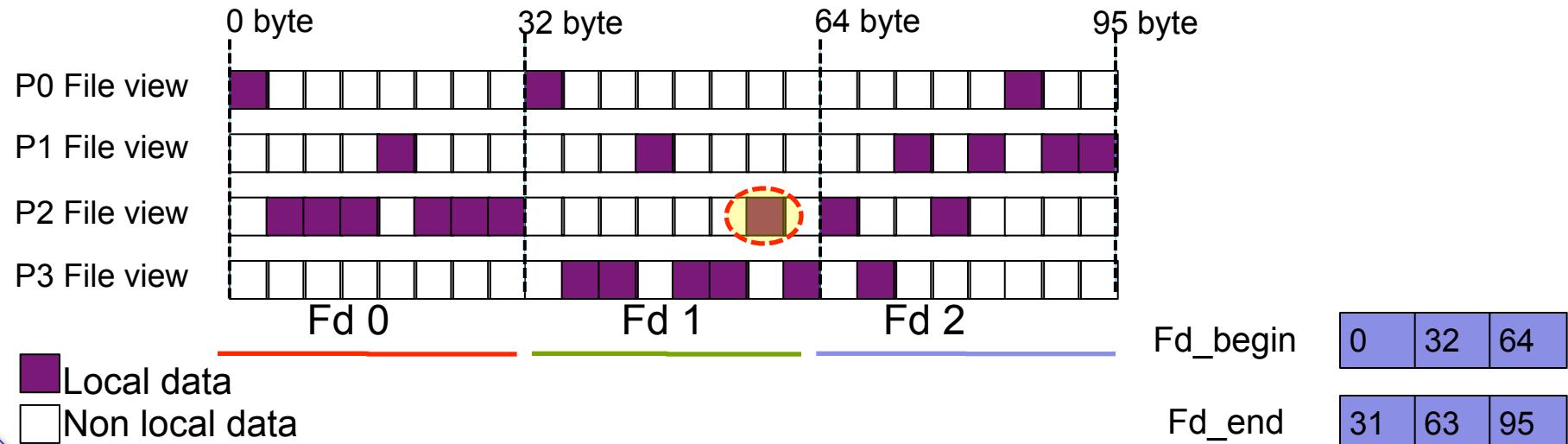
- P0 → ADIOI_Calc_Aggregator
- P1 → ADIOI_Calc_Aggregator
- P2 → ADIOI_Calc_Aggregator
- P3 → ADIOI_Calc_Aggregator

Array_data

→	1	1	1
→	1	1	3
→	2	1	2
→	0	3	1

Calculation and recollection of distribution data

56



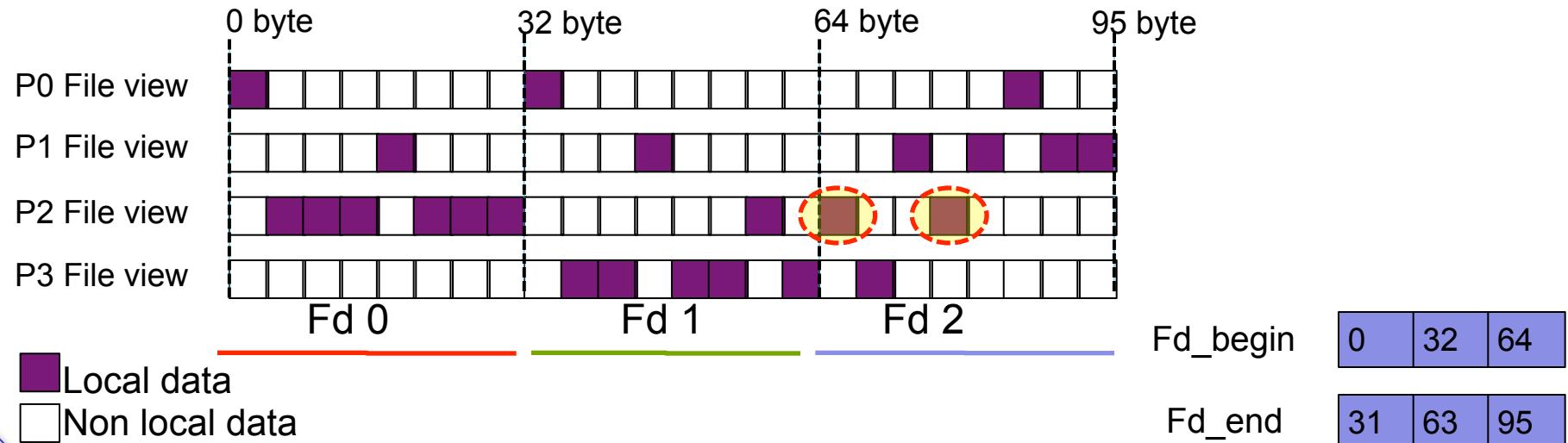
Stage:
St4.1

Array_data

P0	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>1</td><td>1</td><td>1</td></tr></table>	1	1	1
1	1	1			
P1	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>1</td><td>1</td><td>3</td></tr></table>	1	1	3
1	1	3			
P2	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>2</td><td>1</td><td>2</td></tr></table>	2	1	2
2	1	2			
P3	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>0</td><td>3</td><td>1</td></tr></table>	0	3	1
0	3	1			

Calculation and recollection of distribution data

57



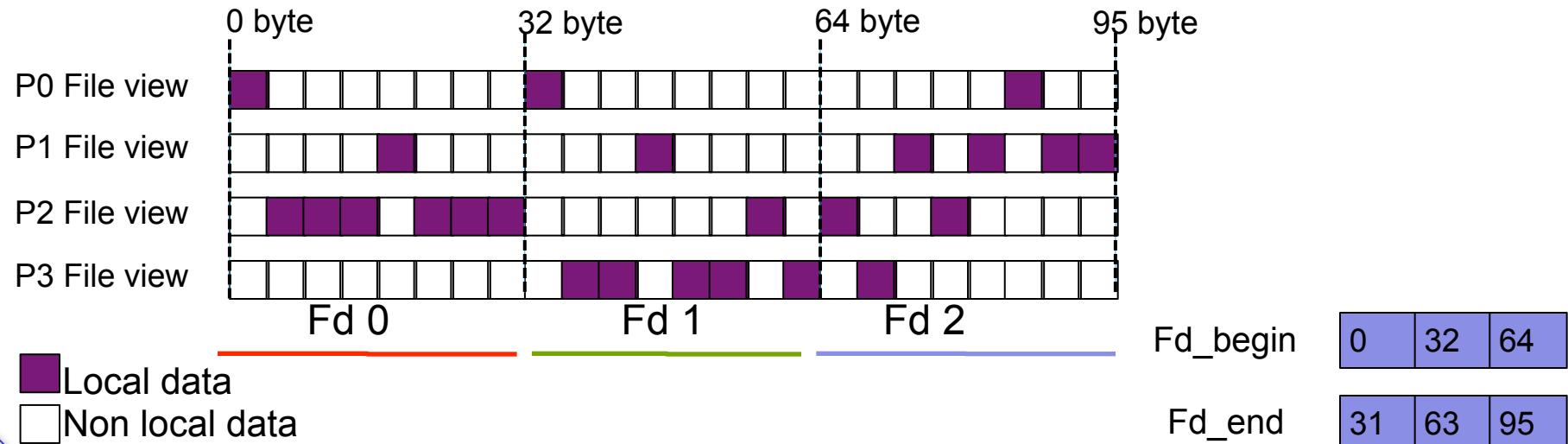
Stage:
St4.1

Array_data

P0	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>1</td><td>1</td><td>1</td></tr></table>	1	1	1
1	1	1			
P1	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>1</td><td>1</td><td>3</td></tr></table>	1	1	3
1	1	3			
P2	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>2</td><td>1</td><td>2</td></tr></table>	2	1	2
2	1	2			
P3	→ ADIOI_Calc_Aggregator	→ <table border="1"><tr><td>0</td><td>3</td><td>1</td></tr></table>	0	3	1
0	3	1			

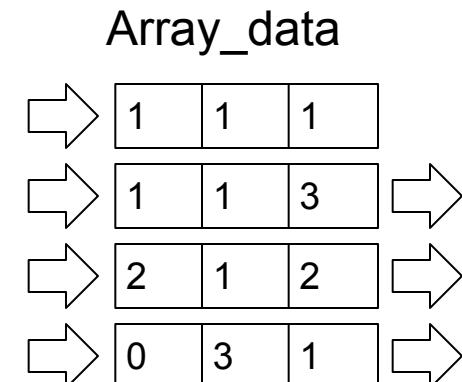
Calculation and recollection of distribution data

58



Stage :
St4.2

- P0 → ADIOI_Calc_Aggregator
- P1 → ADIOI_Calc_Aggregator
- P2 → ADIOI_Calc_Aggregator
- P3 → ADIOI_Calc_Aggregator



P0 Assignment_matrix

	FD0	FD1	FD2
P0	1	1	1
P1	1	1	3
P2	2	1	2
P3	0	3	1

Calculation of new aggregator pattern

59

Stage:
St4.3

LAP →

P0	Assignment_matrix		
	FD0	FD1	FD2
P0	1	1	1
P1	1	1	3
P2	2	1	2
P3	0	3	1

$A_j = \max(P_i, FD)$

Calculation of new aggregator pattern

60

Stage:
St4.3

P0

LAP →

New ranklist

		Assignment_matrix		
		FD0	FD1	FD2
P0	1	1	1	
P1	1	1	3	
P2	2	1	2	
P3	0	3	1	

$A_j = \max(P_i, FD)$

P2	P3	P1
----	----	----

Distribution of new aggregator pattern

61

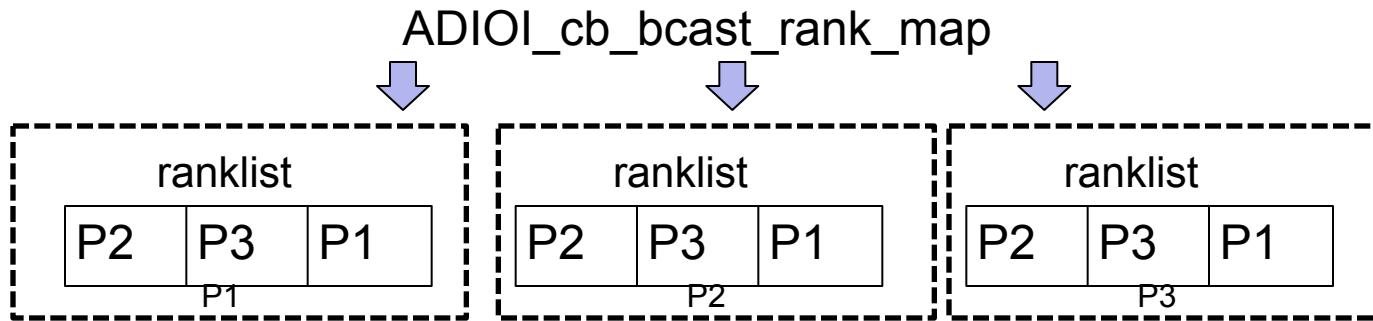
Stage:
St4.3

LAP →

P0	Assignment_matrix		
	FD0	FD1	FD2
P0	1	1	1
P1	1	1	3
P2	2	1	2
P3	0	3	1

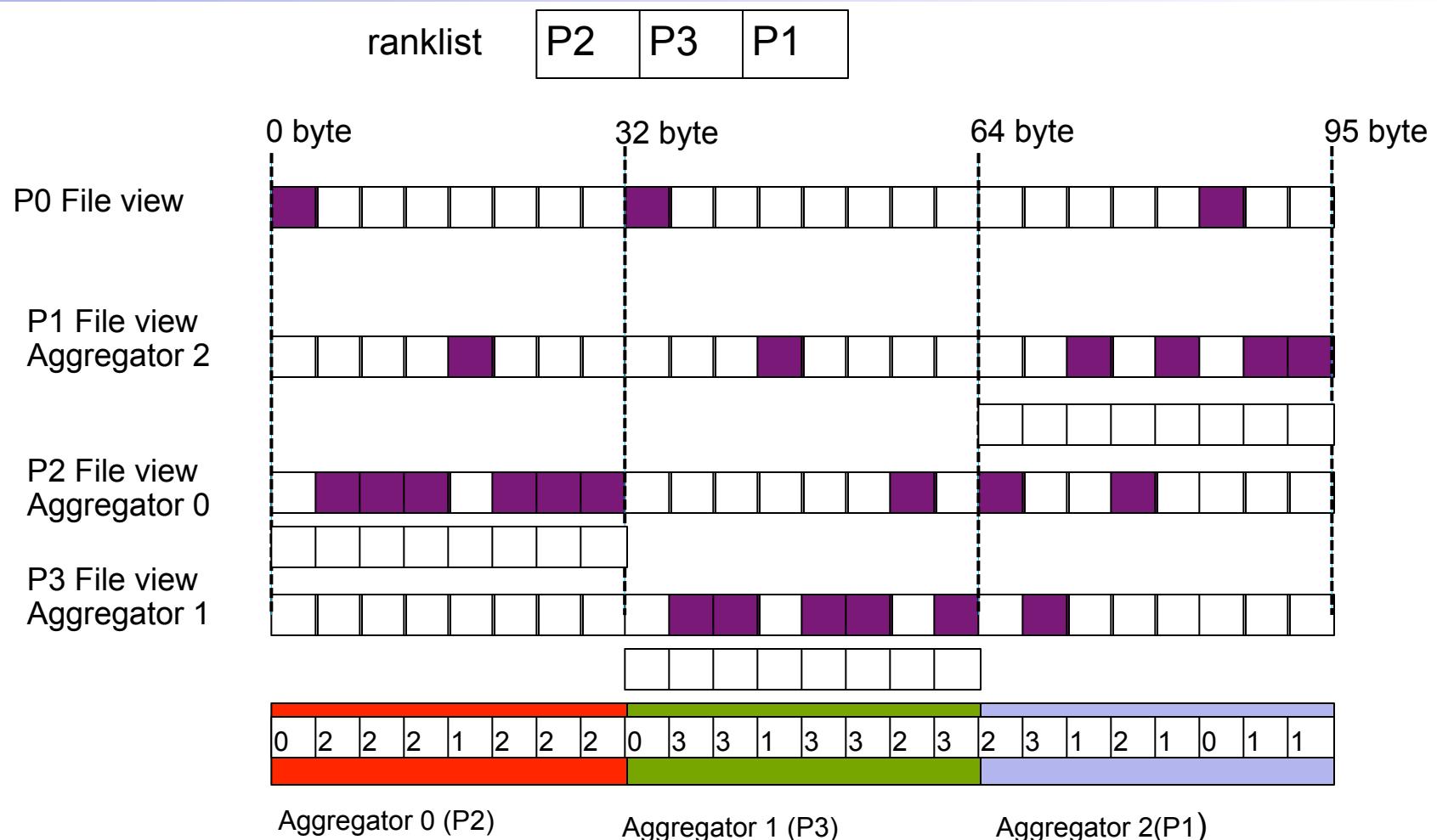
$Aj = \text{Max}(Pi, FD)$

New ranklist: P2 | P3 | P1



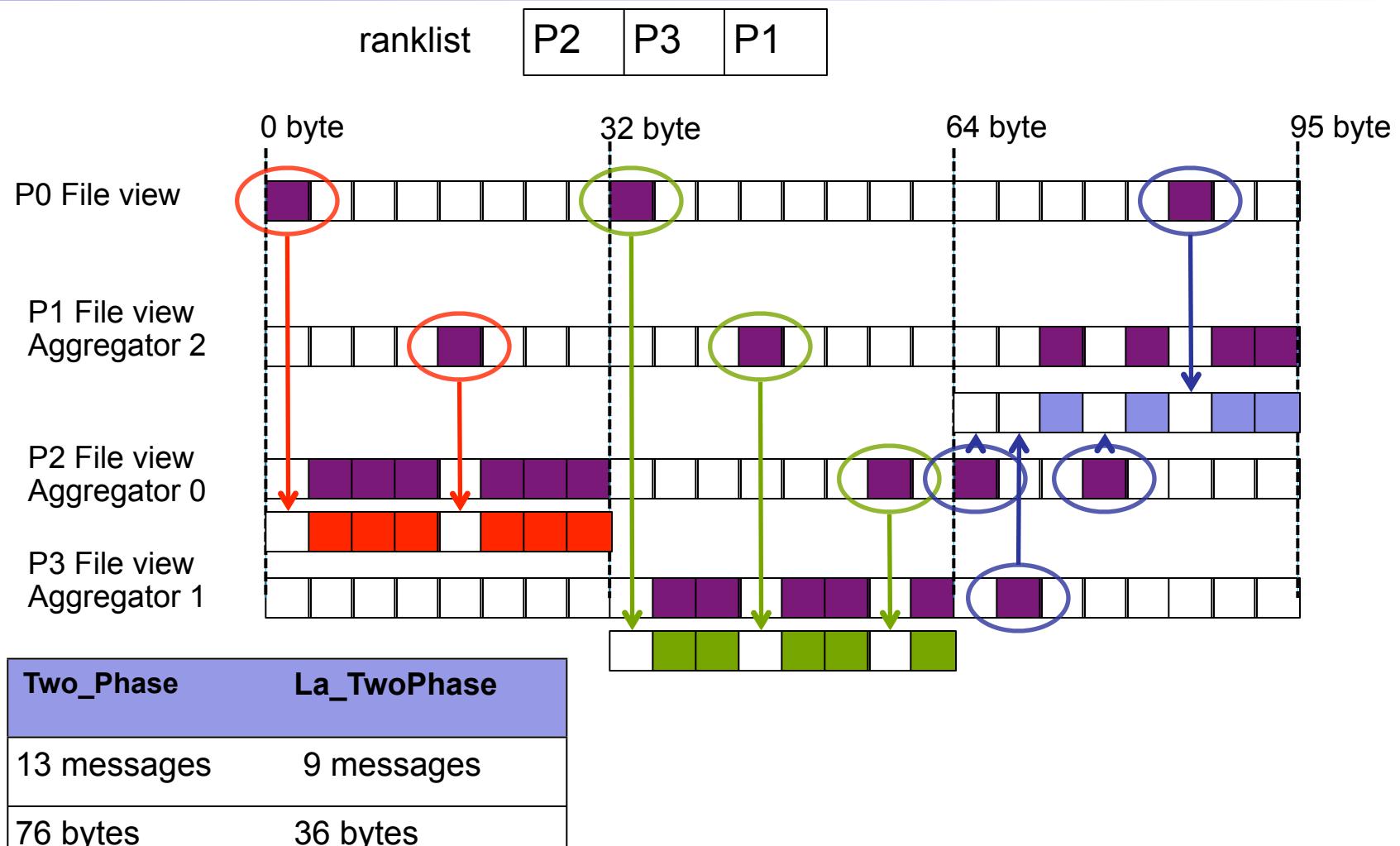
Reduction in the number of communication with the new pattern.

62

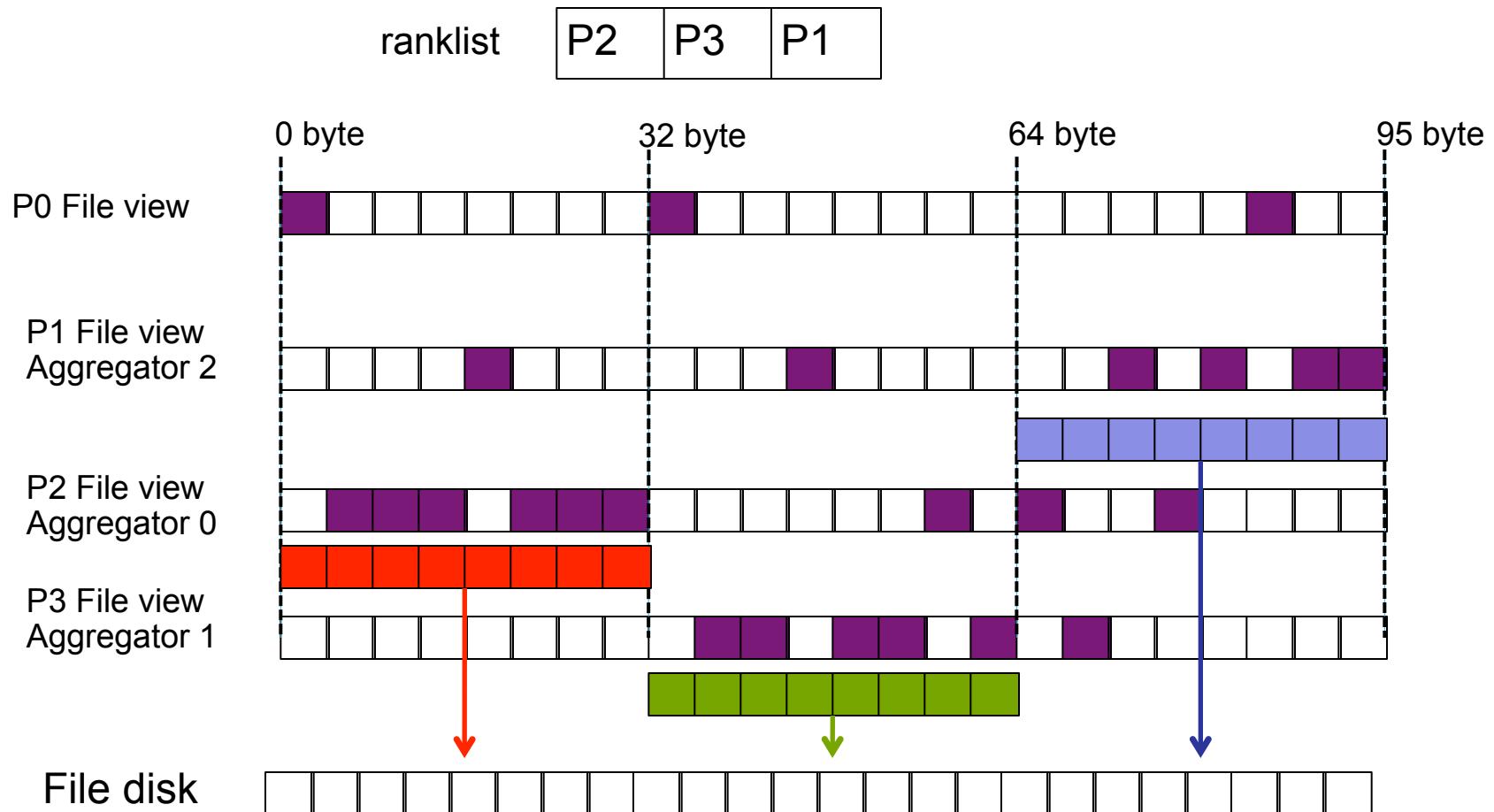


Reduction in the number of communication with the new pattern.

63

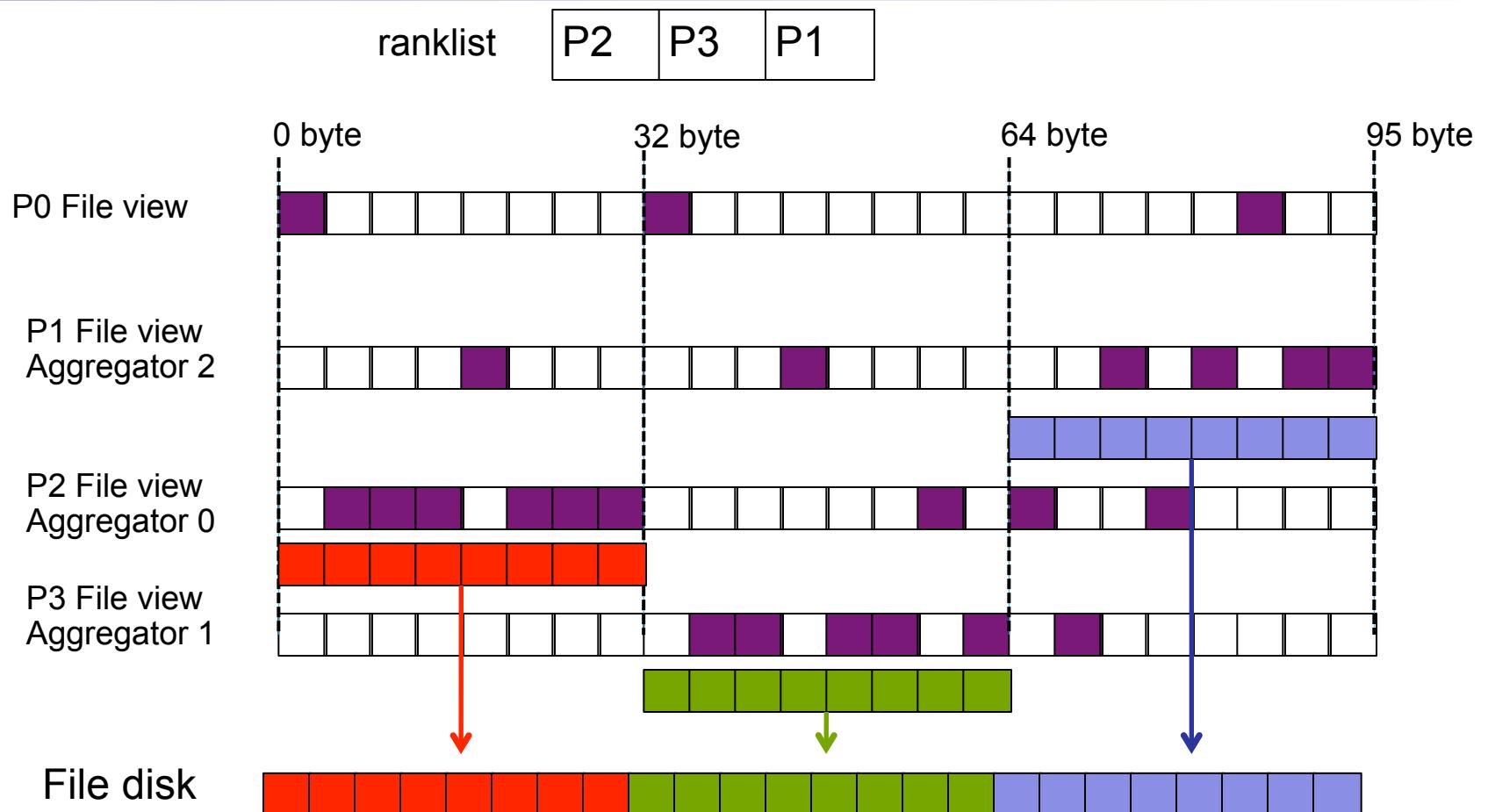


Collective write



Collective write

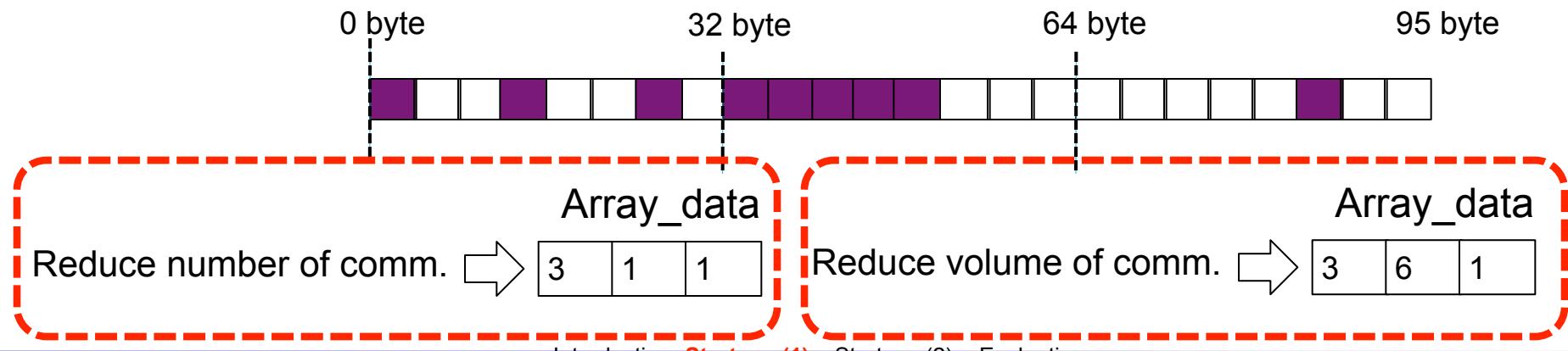
65



New Proposal: Dynamic and adaptive I/O aggregator pattern

66

- We propose to use two aggregation-criteria:
 - Reduce the number of communications:
 - Each aggregator to the **node** who has more highest number of contiguous data blocks.
 - Reduce the volume the communications:
 - Each aggregator to the **node** who has more data of file domain associated with the aggregator.



New Proposal: Dynamic and adaptive I/O aggregator pattern

67

- Is dynamic because is calculated at runtime
- Is adaptive because could select the aggregation-criterion that the reduce more the communication phase.
- Before: each aggregator was assigned according local data that each **cores** stores.
- Now: MPICH2 → according local data of each **node** stores.

Summary

68

1. Problem description
2. Main Objectives
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance communication operations.
5. Evaluation

Strategies for improve the performance of I/O and communication operations.

69

Scientific applications need:

- Large number of computer nodes
- Huge volume of data transferred among the processes



Program model using in cluster → MPI

- Communication among the processes by messages



Communication system has become one of the major limiting factors

- Reduction the performance



Strategy to reduce the cost of communication in MPI by using lossless compression techniques.

Phd. Thesis Proposal: Communication Compression

70

- Reduce the cost of communications:
 - By MPI messages compression in **run-time**
- Lossless compressions algorithms
- Compress **all** MPI primitives.
- We have developed three different strategies:
 - Runtime Compression (RC)
 - Runtime Adaptive Compression (RAS)
 - Guided Strategy (GS)

Phd. Thesis Proposal: Communication Compression

71

- Reduce the cost of communications:
 - By MPI messages compression in **run-time**
 - Lossless compressions algorithms
 - Compress **all** MPI primitives.
 - We have developed three different strategies:
 - Runtime Compression (RC)
 - Runtime Adaptive Compression (RAS)
 - Guided Strategy (GS)
- 

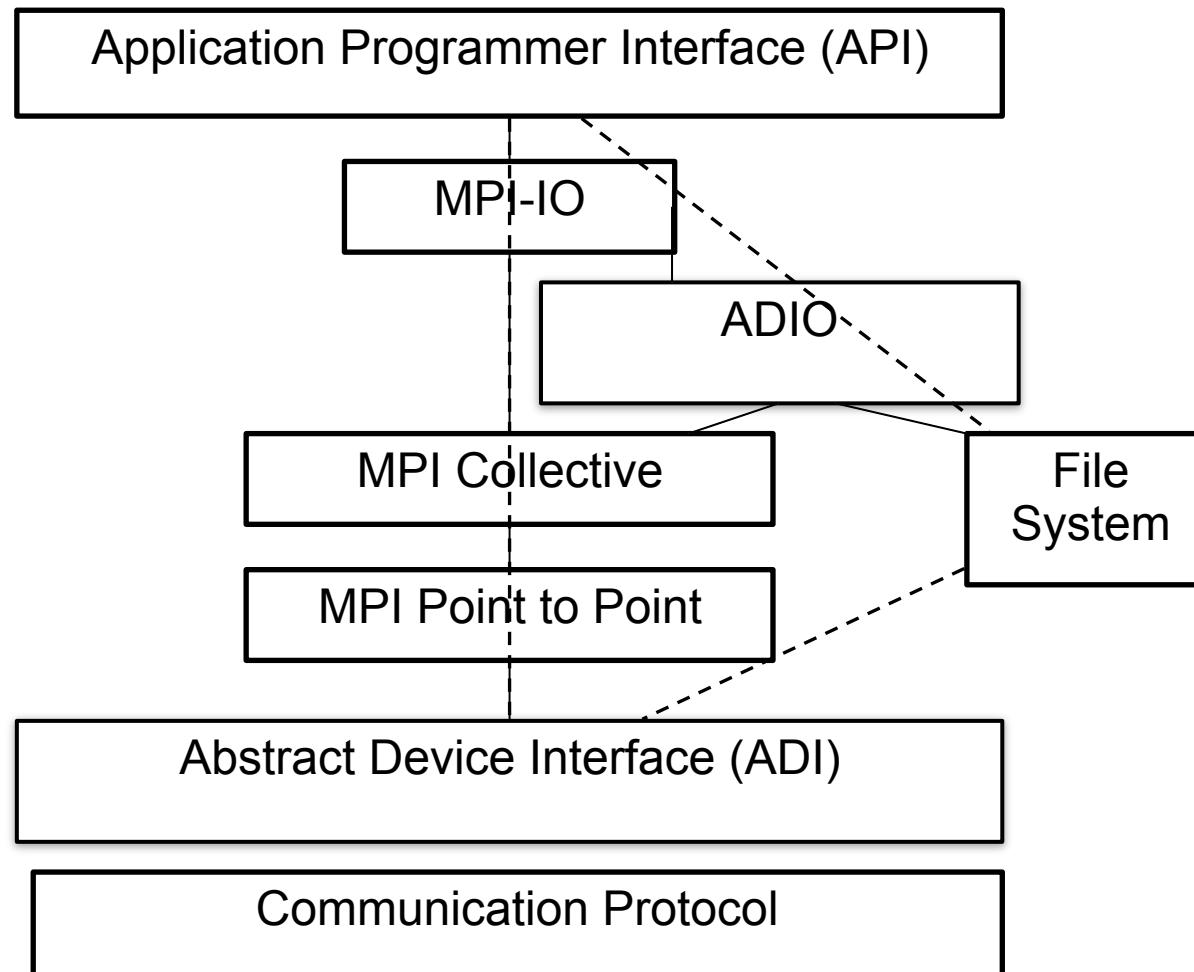
1 Compression Strategy: Runtime Compression

72

- Runtime Compression (RC) uses different compression algorithms:
 - LZO, RLE, HUFFMAN, RICE, FPC.
- Compress all the transferred message (> 2KB) with the algorithm indicated by the user (`MPI_hint`)
- Compression not can be disable
- Compress all MPI primitives:
 - Messages: point-to-point and collectives.
 - Communications: blocking and no blocking
 - Datatypes: contiguous y non-contiguous

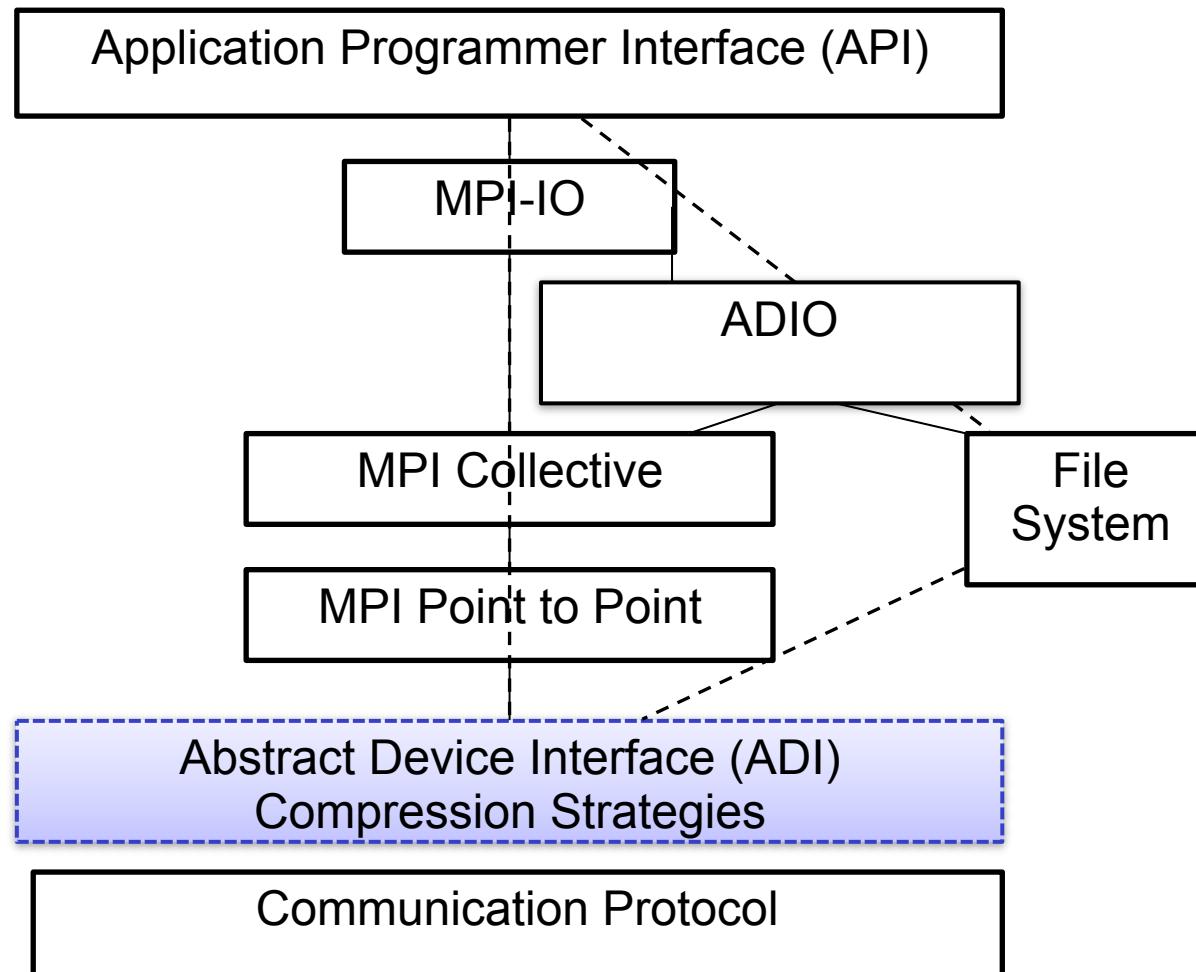
MPICH1.2 Architecture

73



MPICH1.2 Architecture

74





Modification over ADI layer

75

- Header in the exchanged message to inform:
 - Compression used or not, algorithm and length
 - All compression algorithm are included in a single **Compression Library**:
 - We can add new compression algorithms easily.

Compressi
on stages

- Messag
e size



Decompre
ssion stages

- Header



2 Compression Strategy: Runtime Adaptive Compression

76

- Runtime Adaptive Compression Strategy (RAS), per message transferred takes two decision:
 - Turn on and off the compression.
 - Select itself the best compression algorithm.
- Learn in run-time from previous messages
- Decision depending on:
 - Message feature:
 - Datatype and length
 - Network performance:
 - Latency and bandwidth
 - Compression algorithms

Decisions → Speedup (1/2)

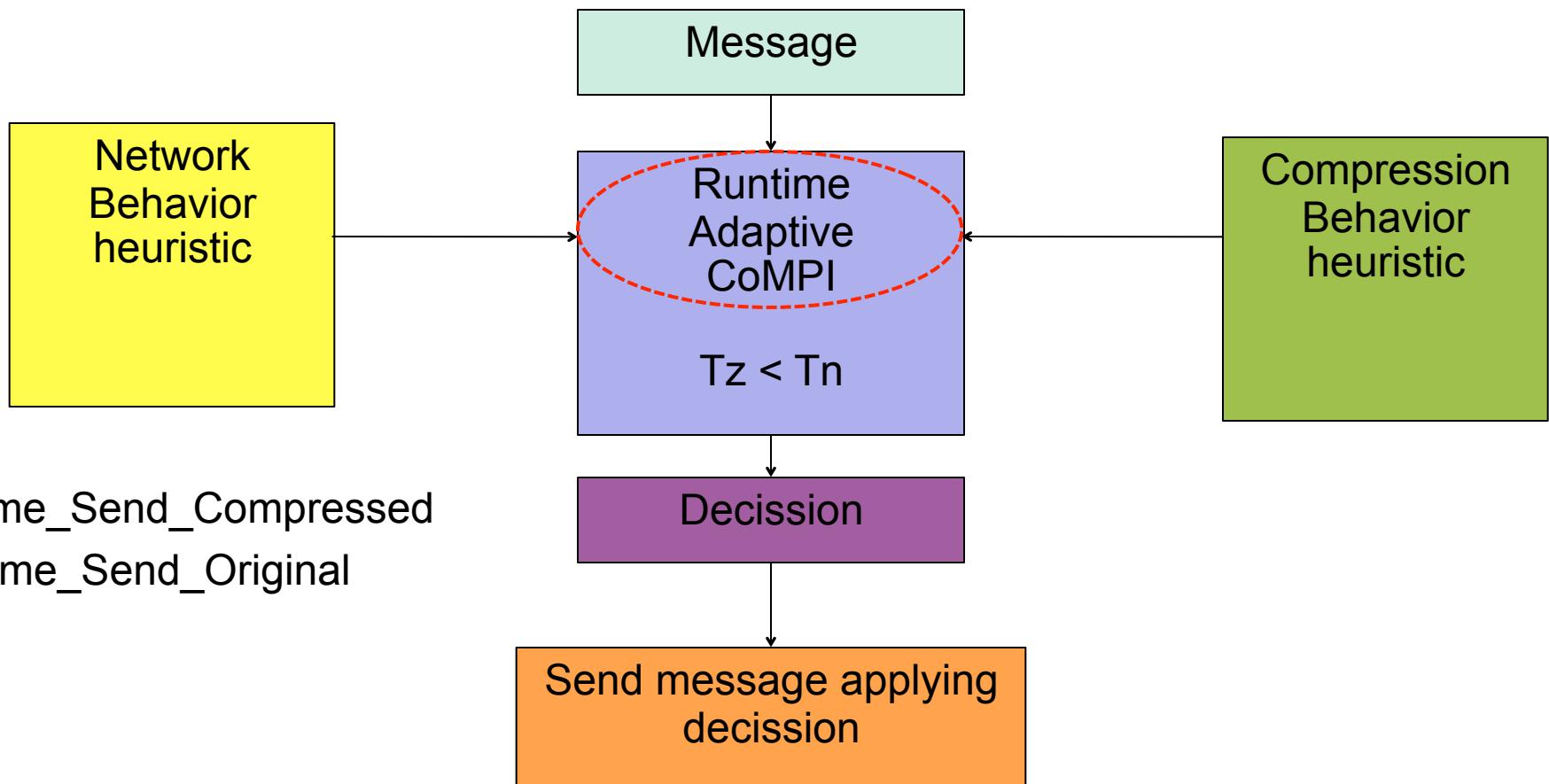
77

- Speedup to decide if send the message with/without compression.
- So the decision depends:
 - Original message transmission time
 - Compressed message transmission time
 - Compression and decompression time

$$Speedup = \frac{Time_Sent_Orig.}{(Time_Sent_Compr.+ time_compress.+ time_decompr.)}$$

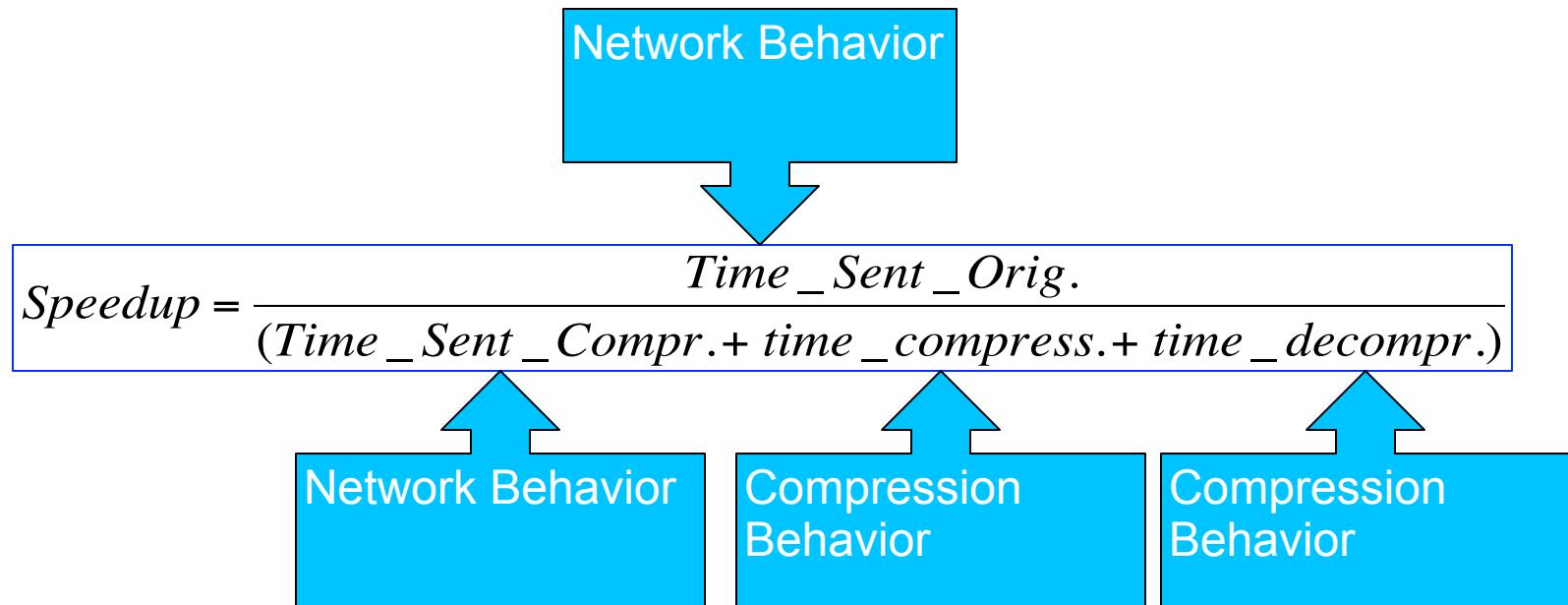
Decisions → Speedup

78



Compression behavior

79



Decision Methodology

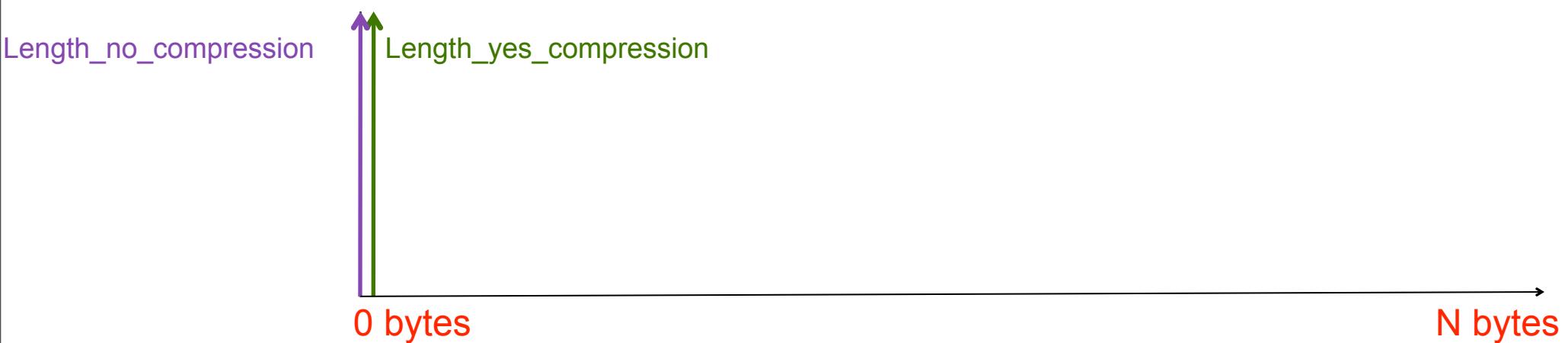
80

- Calculate the speedup per message? No → **high overhead computation time**
- According to Compression Behavior and Network data Behavior, RAS decides:
 - Datatype:
 - Integer y Float → LZO
 - Double → LZO or FPC
 - Others → LZO, RLE, RICE or HUFFMAN
 - Message size → Decision Threshold:
 - Each datatype has its thresholds
 - Length_yes_compression
 - Length_no_compression

Decision Methodology

81

Process 0

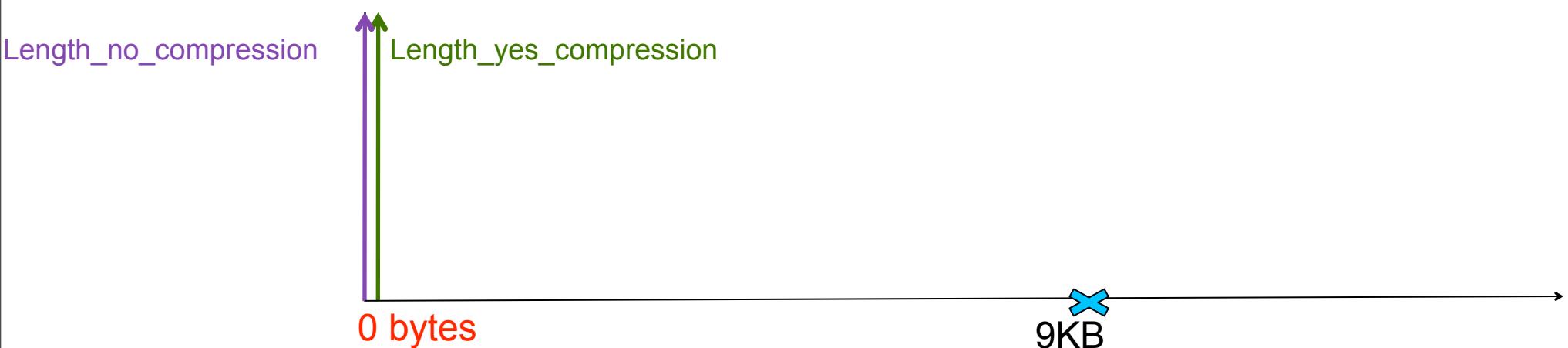


Decision Methodology

82

Process 0

Message Integer
Size 9KB



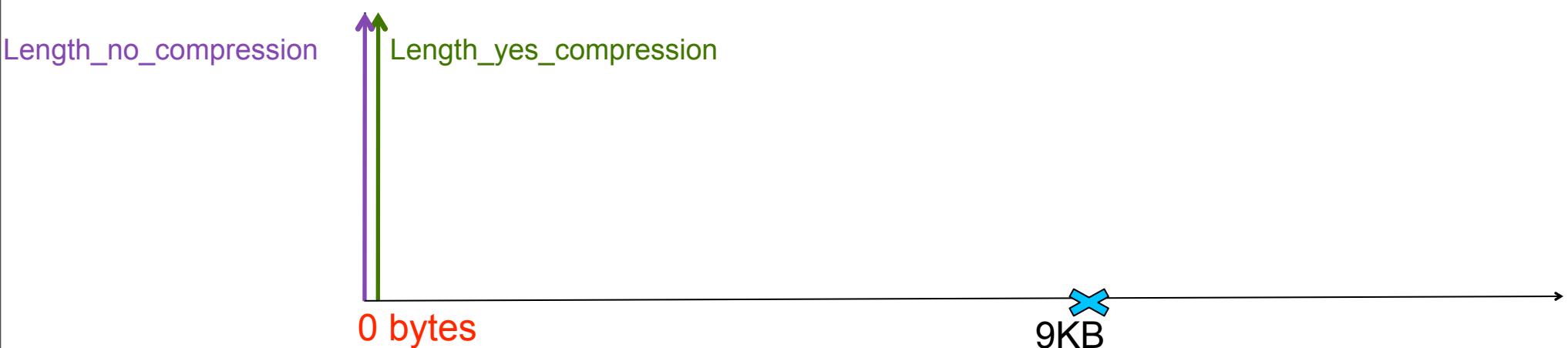
Decision Methodology

83

Process 0

Message Integer
Size 9KB

Study size message
And thresholds

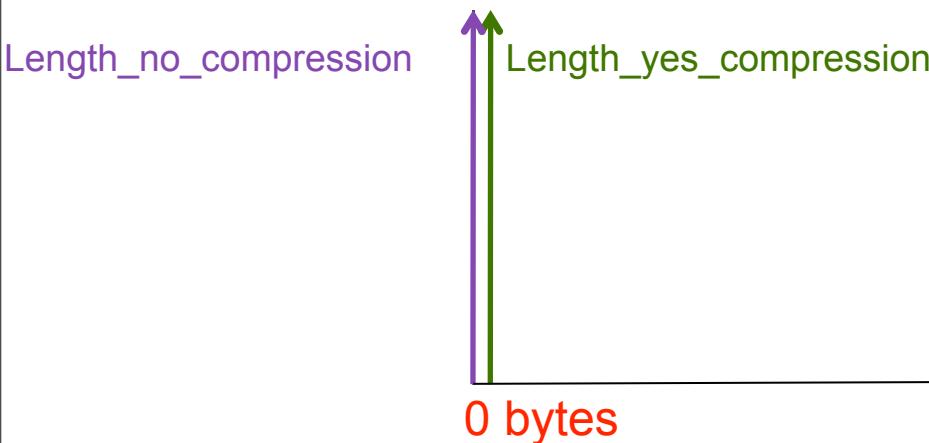
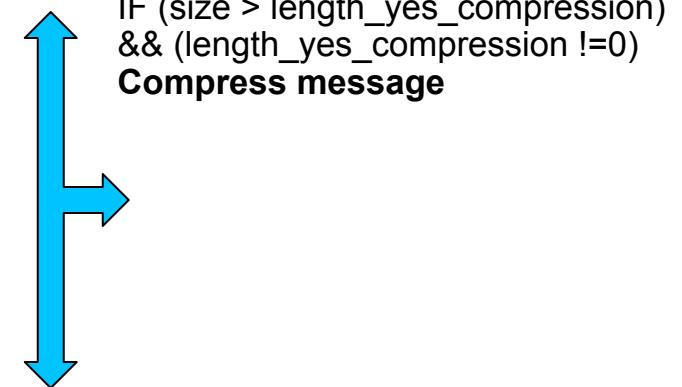


Decision Methodology

Process 0

Message Integer
Size 9KB

Study size message
And thresholds



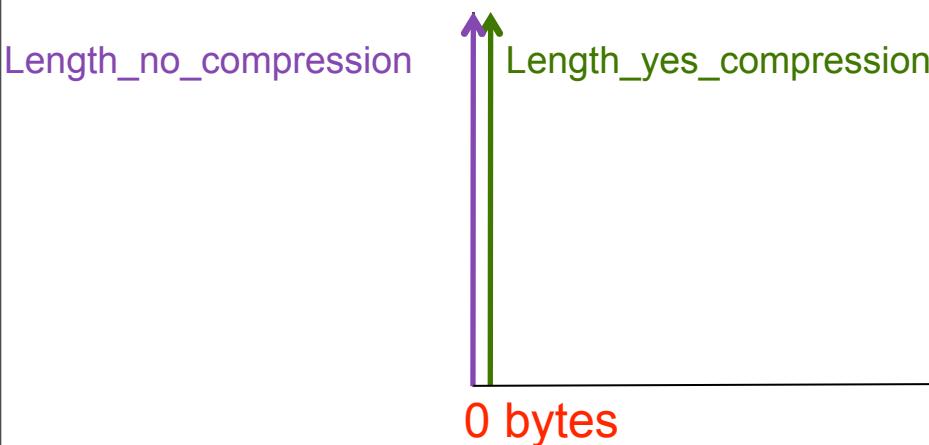
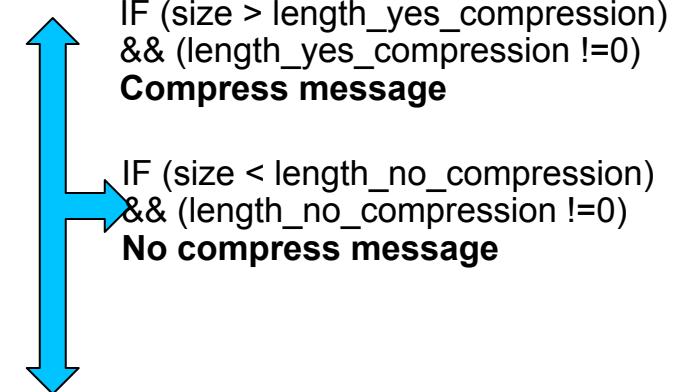
Decision Methodology

85

Process 0

Message Integer
Size 9KB

Study size message
And thresholds



Decision Methodology

Process 0

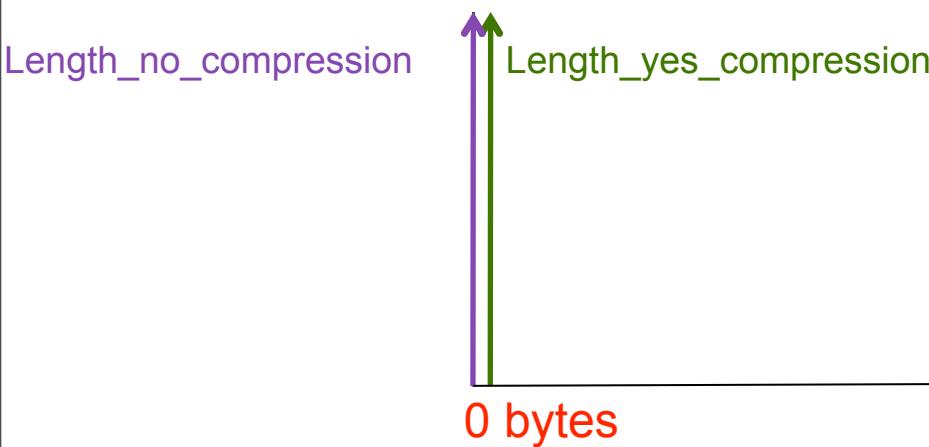
Message Integer
Size 9KB

Study size message
And thresholds

IF (size > length_yes_compression)
 $\&\&$ (length_yes_compression !=0)
Compress message

IF (size < length_no_compression)
 $\&\&$ (length_no_compression !=0)
No compress message

IF (Between thresholds)
 $\|$ (thresholds == 0)
Speedups calculation



Decision Methodology

87

Process 0

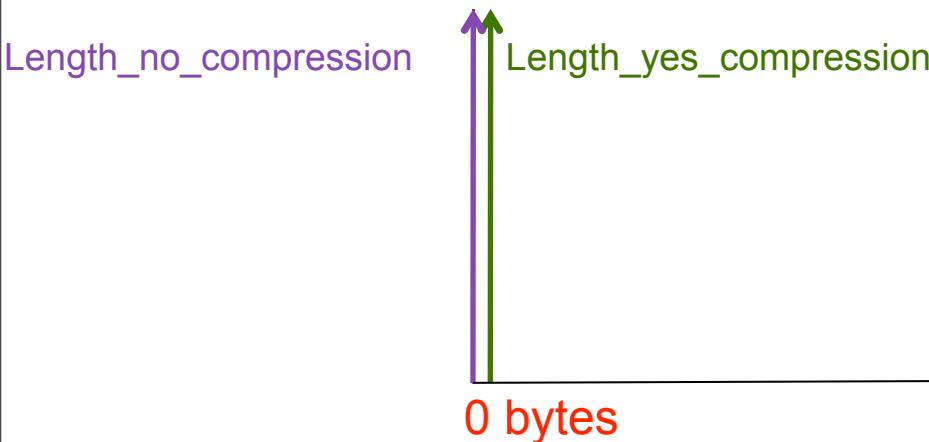
Message Integer
Size 9KB

Study size message
And thresholds

IF (size > length_yes_compression)
&& (length_yes_compression !=0)
Compress message

IF (size < length_no_compression)
&& (length_no_compression !=0)
No compress message

IF (Between thresholds)
|| (thresholds == 0)
Speedups calculation



Decision Methodology

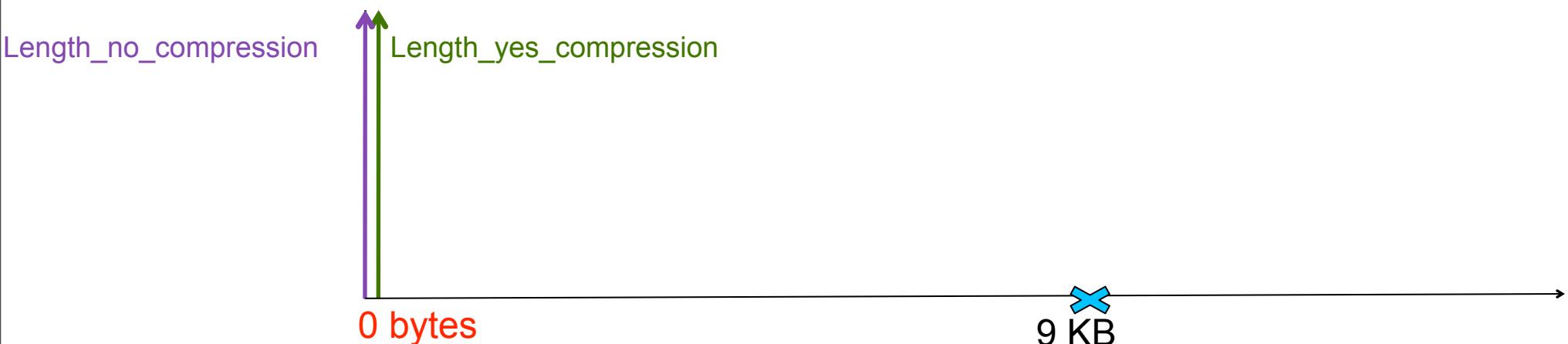
88

Process 0

Message Integer
Size 9KB

Speedups
Calculation

- I. Select compressors
- II. Speedups calculation
- III. Select compressor with greater speedup
- IV. Or (if all speedup less or equal 0)
No compress



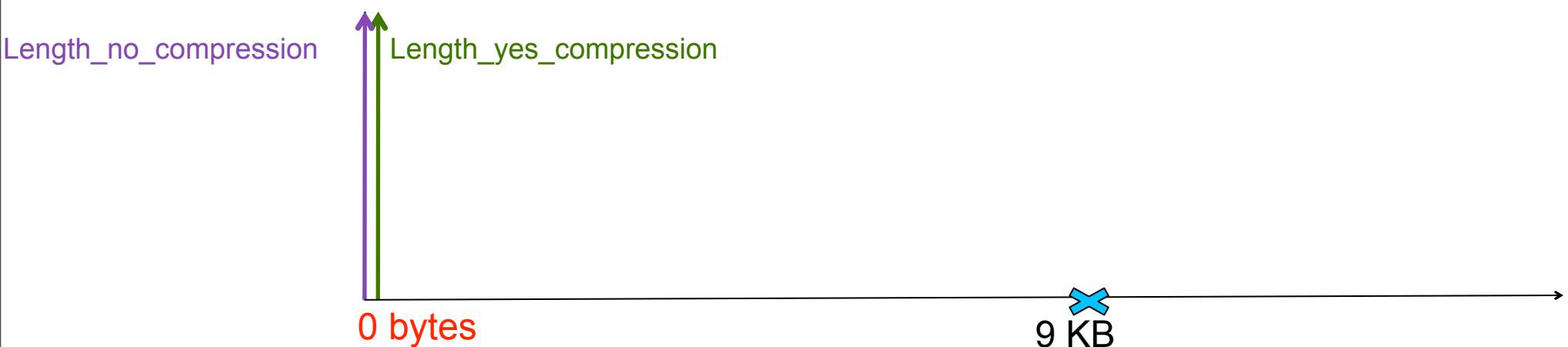
Decision Methodology

89

Process 0

Message Integer
Size 9KB

Speedups
Calculation → Compress with LZO



Decision Methodology

90

Process 0

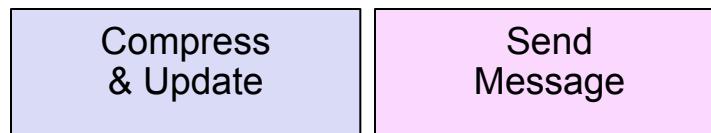
Compress
& Update



Decision Methodology

91

Process 0



Decision Methodology

92

Process 0

Message Integer
Size 16KB



Decision Methodology

93

Process 0

Message Integer
Size 16KB

Study size message
And thresholds

IF (size > length_yes_compression)
&& (length_yes_compression !=0)
Compress message

IF (size < length_no_compression)
&& (length_no_compression !=0)
No compress message

IF (Between thresholds)
|| (thresholds == 0)
Speedups Calculation

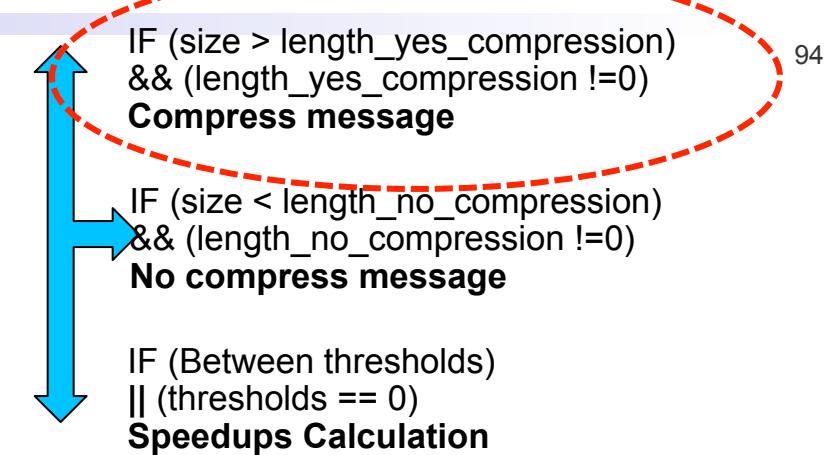


Decision Methodology

Process 0

Message Integer
Size 16KB

Study size message
And thresholds



Length_no_compression



0 bytes

Length_yes_compression

9 KB

16 KB

94

Decision Methodology

95

Process 0

Compression
With LZO



Decision Methodology

96

Process 0

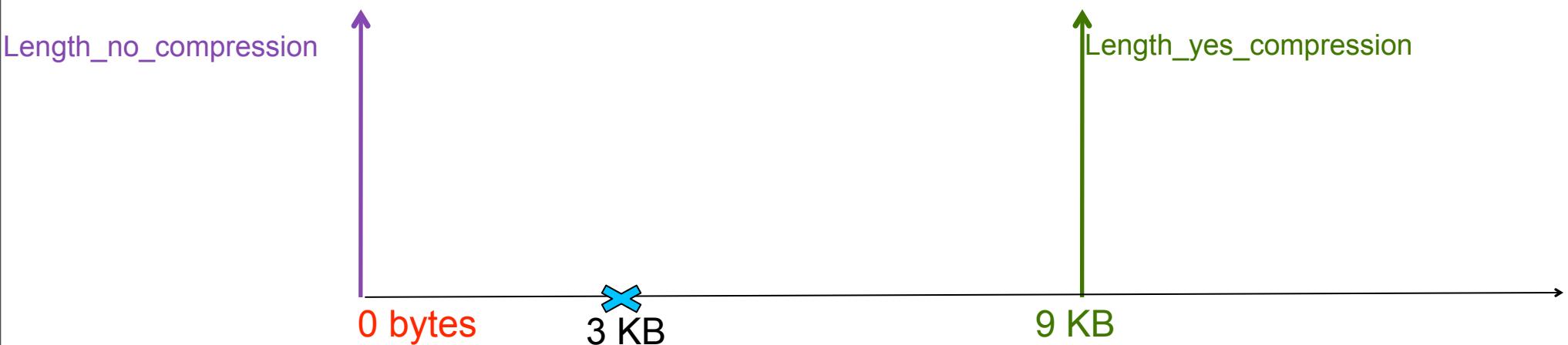


Decision Methodology

97

Process 0

Message Integer
Size 3KB



Decision Methodology

98

Process 0

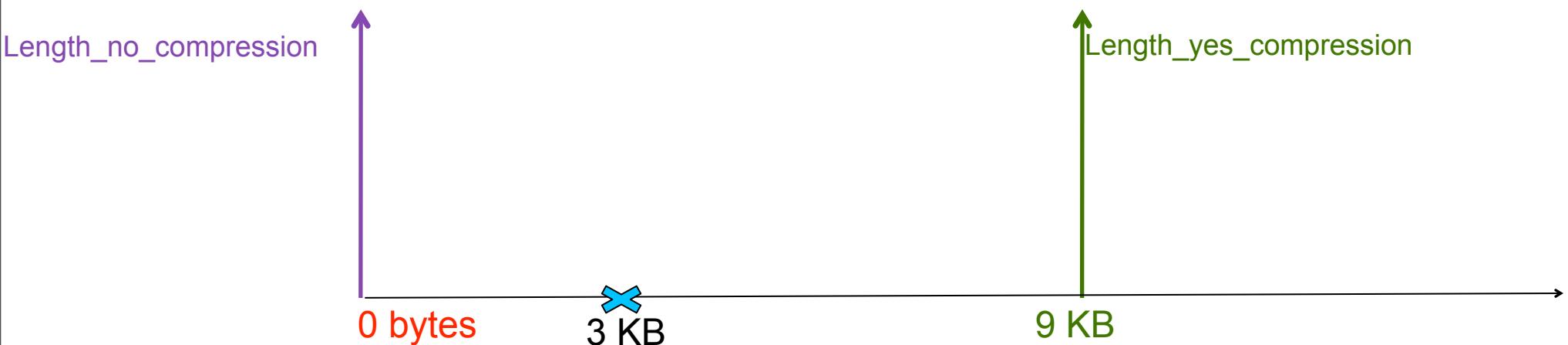
Message Integer
Size 3KB

Study size message
And thresholds

IF (size > length_yes_compression)
&& (length_yes_compression !=0)
Compress message

IF (size < length_no_compression)
&& (length_no_compression !=0)
No compress message

IF (Between thresholds)
|| (thresholds == 0)
Speedups Calculation



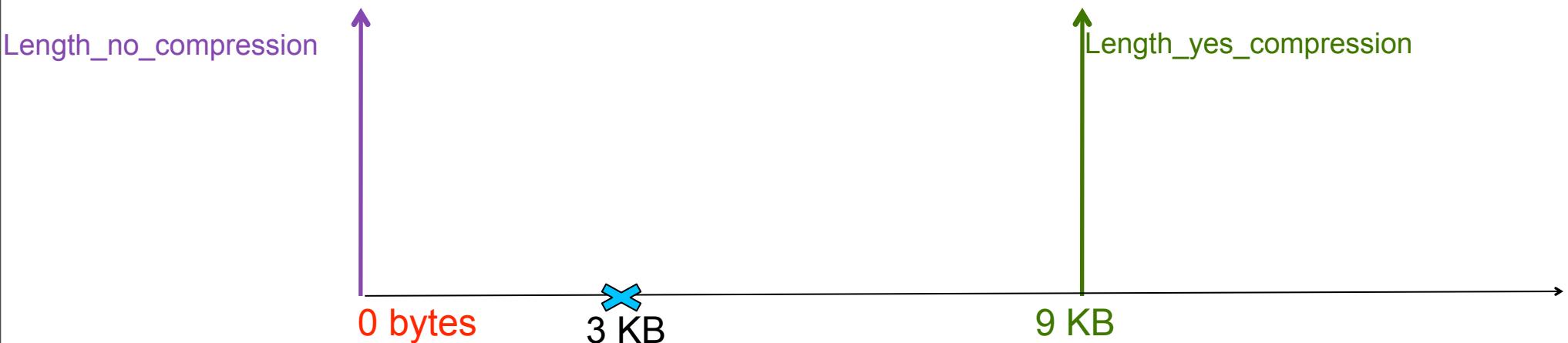
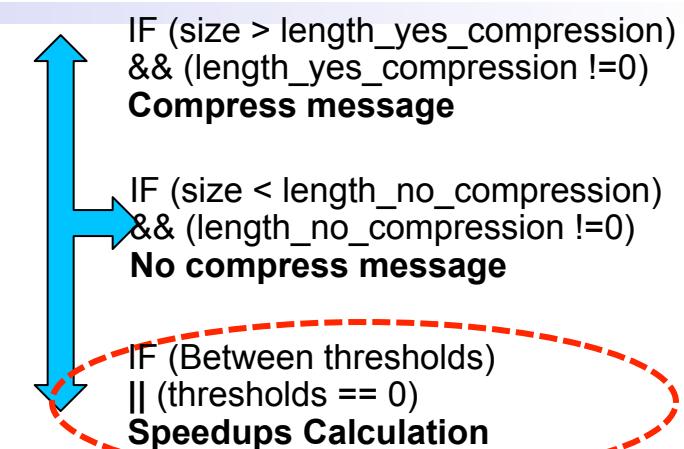
Decision Methodology

99

Process 0

Message Integer
Size 3KB

Study size message
And thresholds



Decision Methodology

100

Process 0

Message Integer
Size 3KB

Speedups
Calculation → No Compress



Decision Methodology

101

Process 0

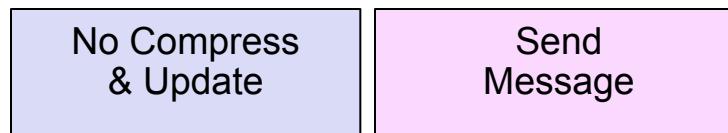
No Compress
& Update



Decision Methodology

102

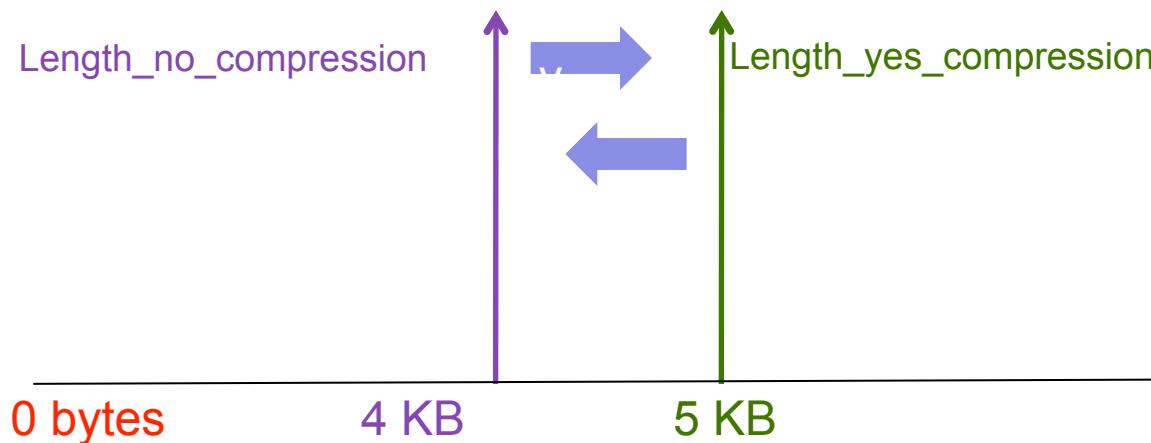
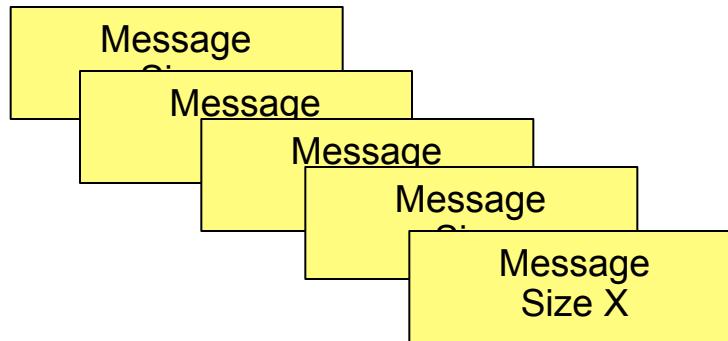
Process 0



Decision Methodology

103

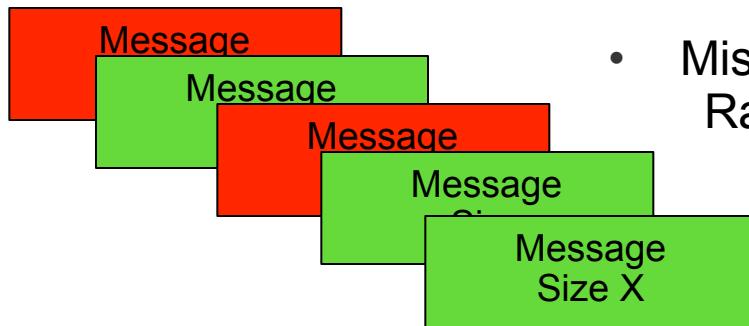
Process 0



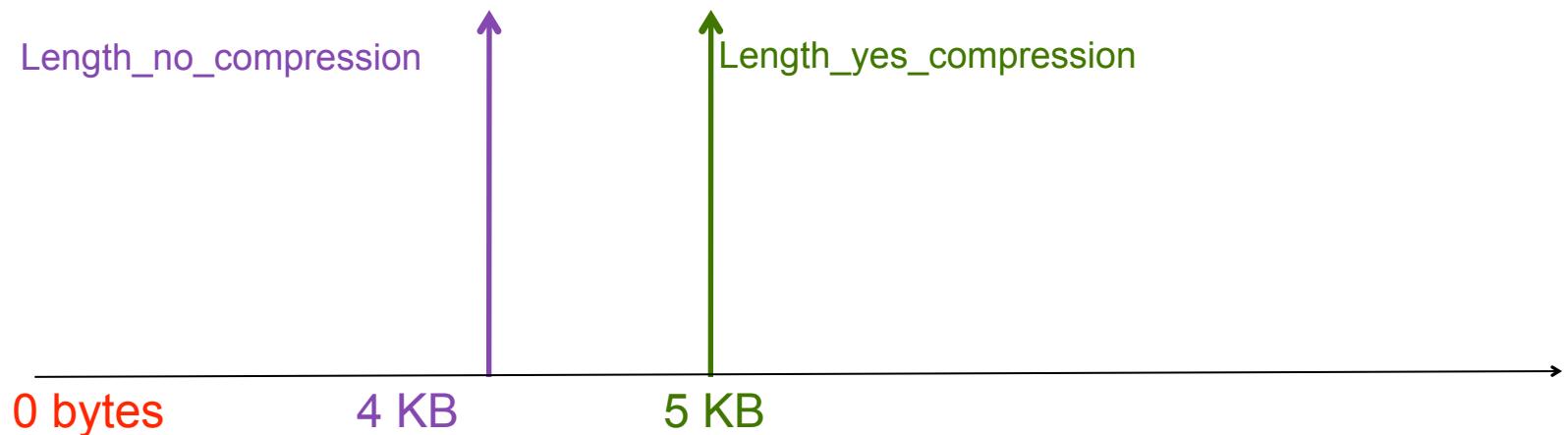
Re-evaluation Methodology

104

Process 0



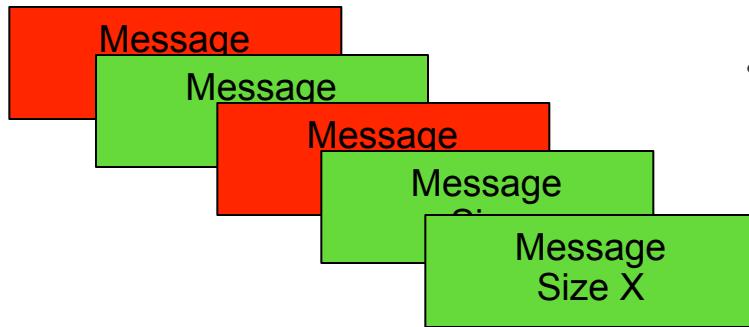
- Mistake:
Ratio compress message < MIN_RATIO



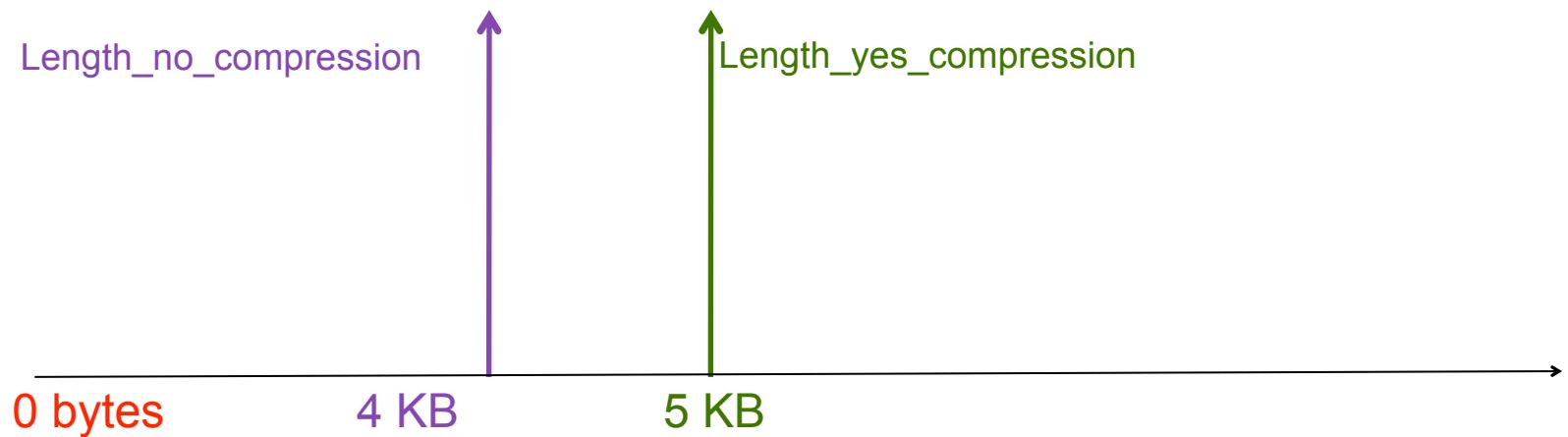
Re-evaluation Methodology

105

Process 0



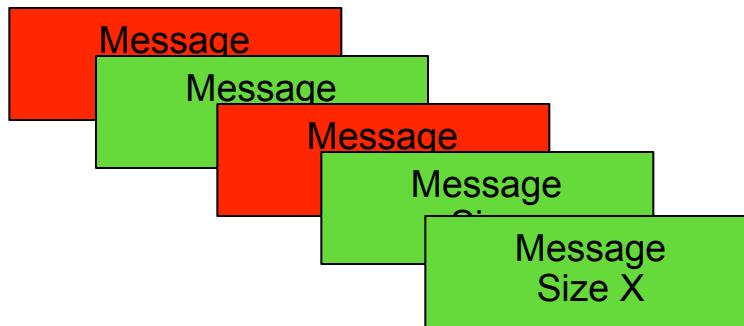
- Number of message without compress > MAX_MESSAGE



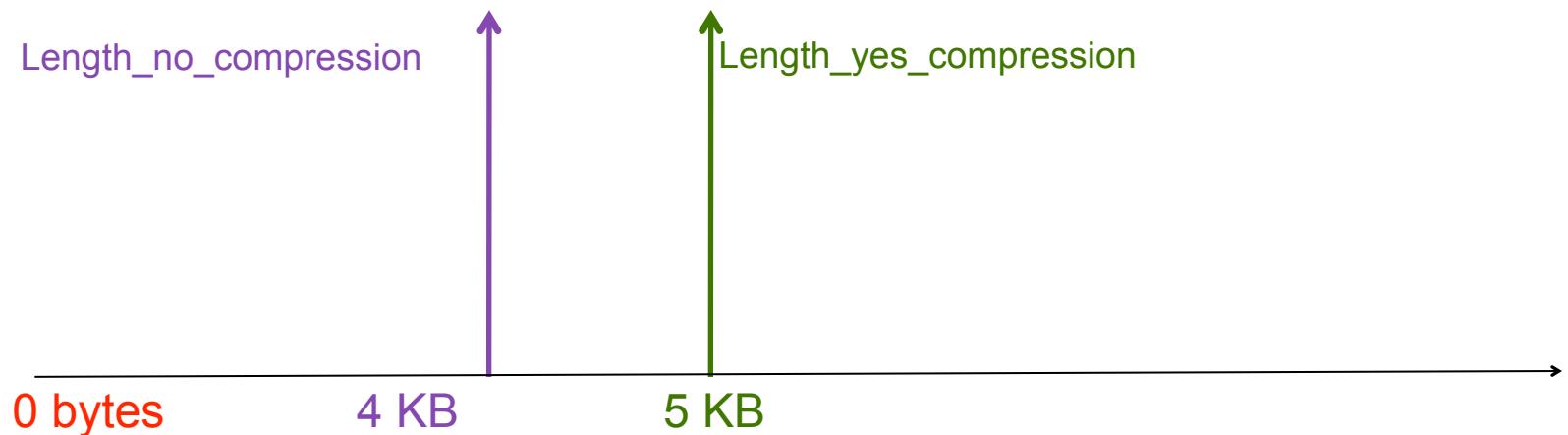
Re-evaluation Methodology

106

Process 0



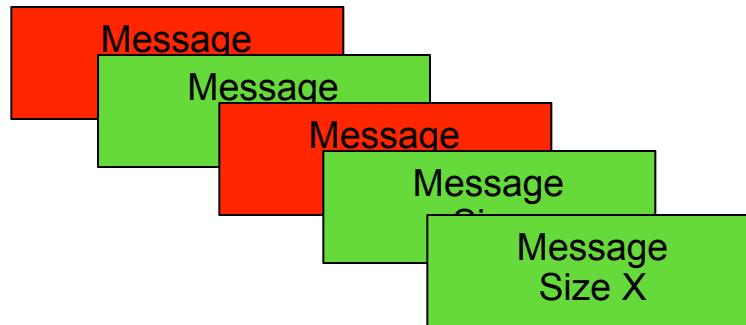
- Number of message without compress > MAX_MESSAGE
- Number Mistakes > MAX_MISTAKES
&& Frequency Mistakes < MIN_FRECUENCY



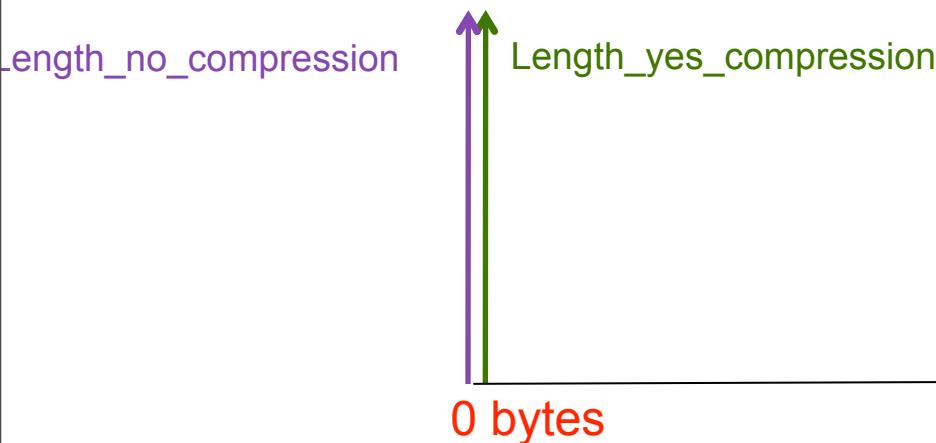
Re-evaluation Methodology

107

Process 0



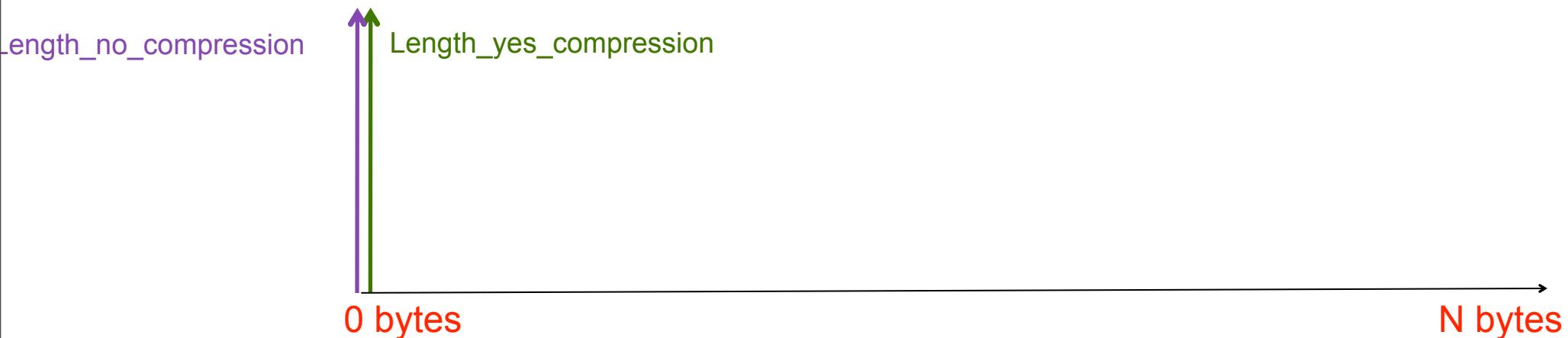
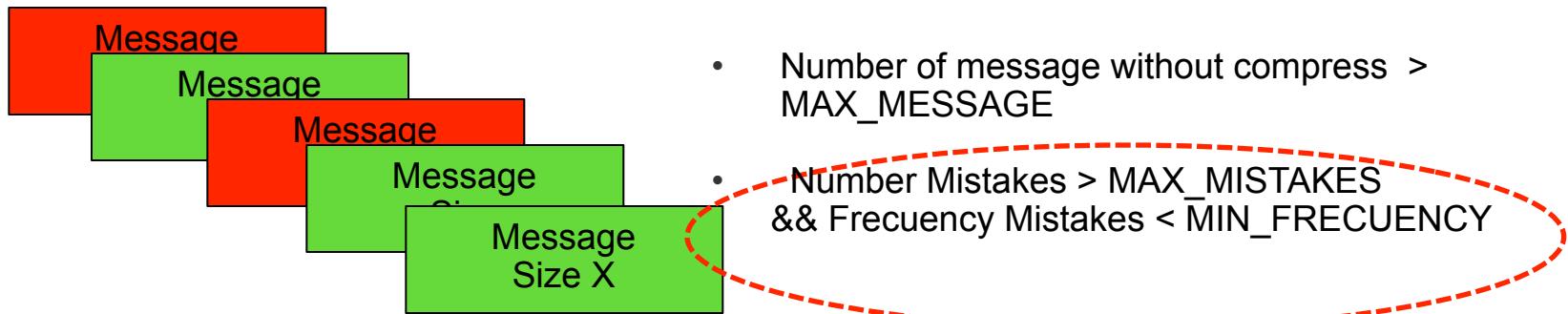
- Number of message without compress > MAX_MESSAGE
- Number Mistakes > MAX_MISTAKES
&& Frequency Mistakes < MIN_FRECUENCY



Re-evaluation Methodology

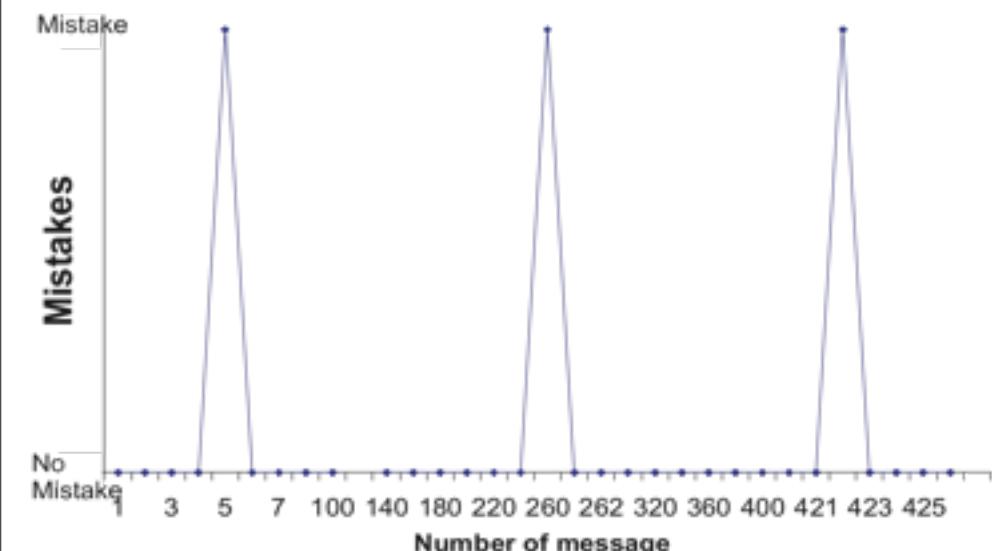
108

Process 0



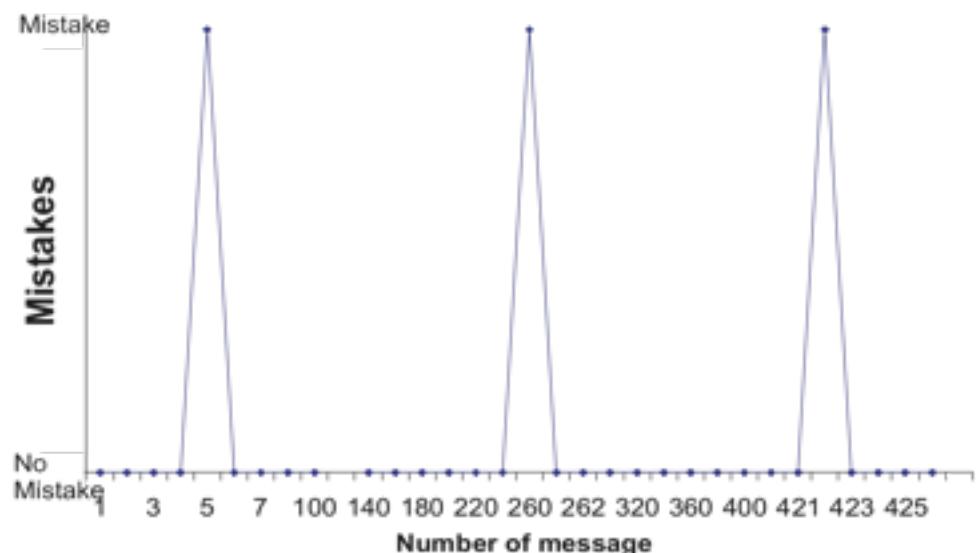
Different cases of re-evaluation

109

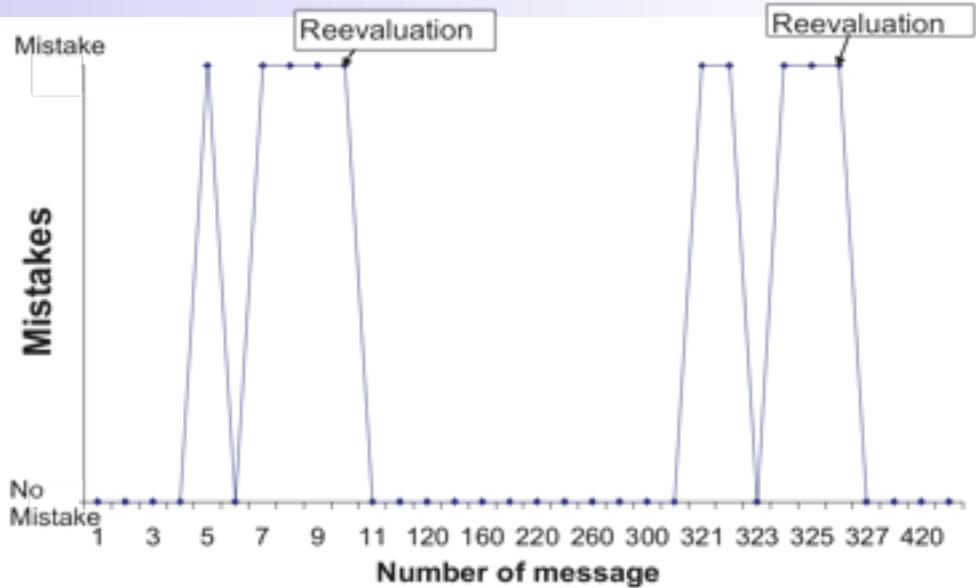


High frequency mistakes

Different cases of re-evaluation

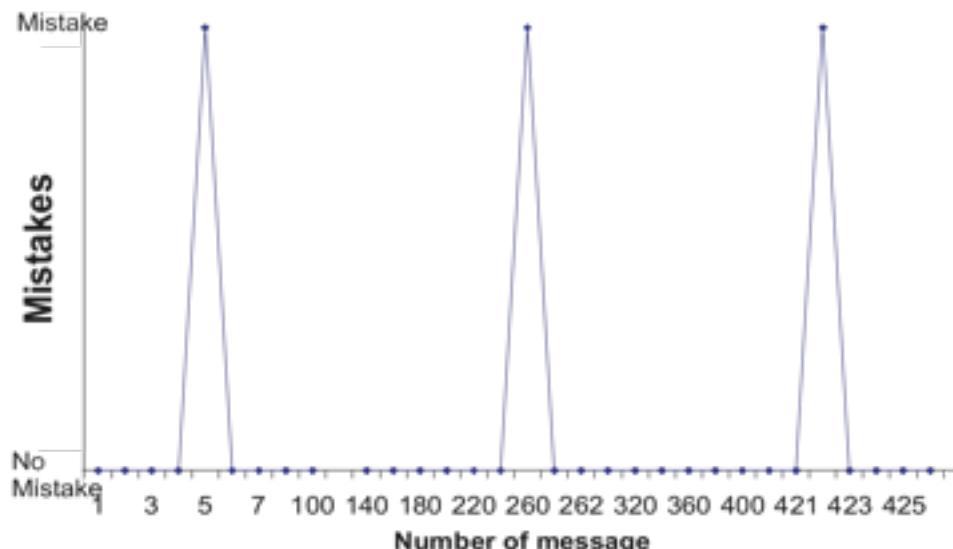


High frequency mistakes

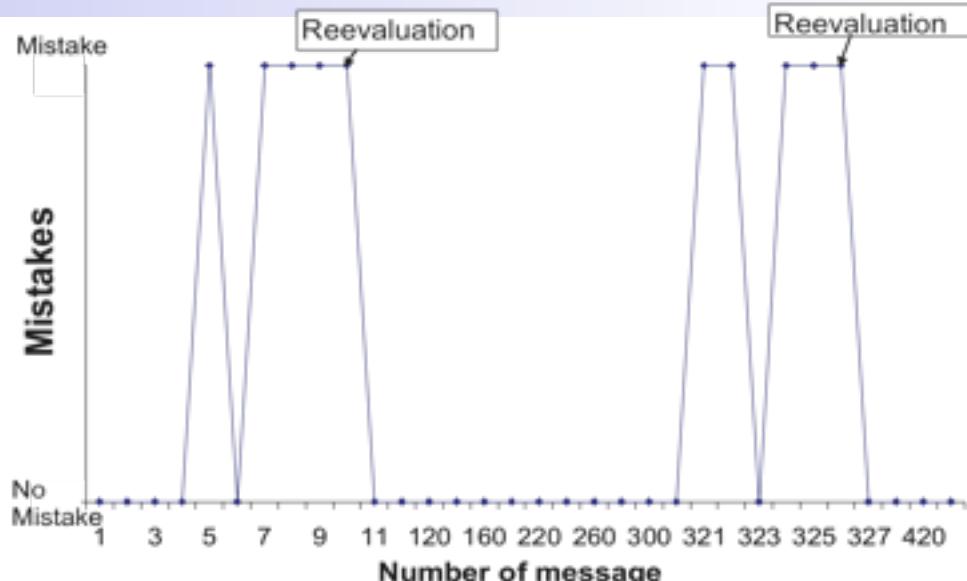


Low frequency mistakes

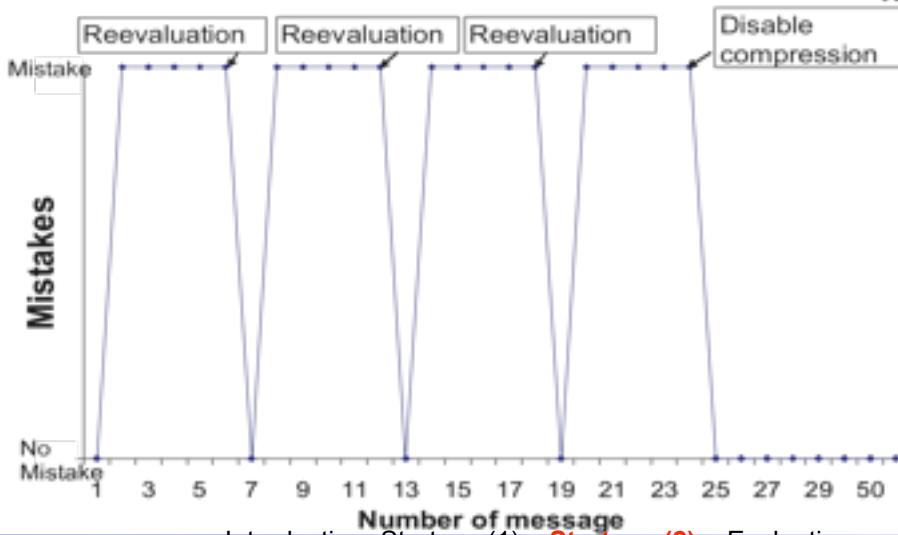
Different cases of re-evaluation



High frequency mistakes



Low frequency mistakes



Many reevaluations
in a short time

3 Compression Strategy: Guided Strategy

112

- Scientific Application:
 - Are executed many times in the same cluster
 - With same input parameters.
- Guided strategy (GS) → takes the decisions off-line by storing all message from the first execution.
- Decisions (compress or not, and which algorithm) are stored in a file called Decision Rules
- Subsequent executions of the same application and same input parameters, GS applies rules of Decision Rules.

And now??

113

- Use MPICH2 instead MPICH1.2 and adapt the different compression strategies.
- No apply compression when the processes are in the same node
- New compression algorithms
- New ideas:
 - Take account the features of communications
 - MPI_Broadcast → Apply compression once per operation.

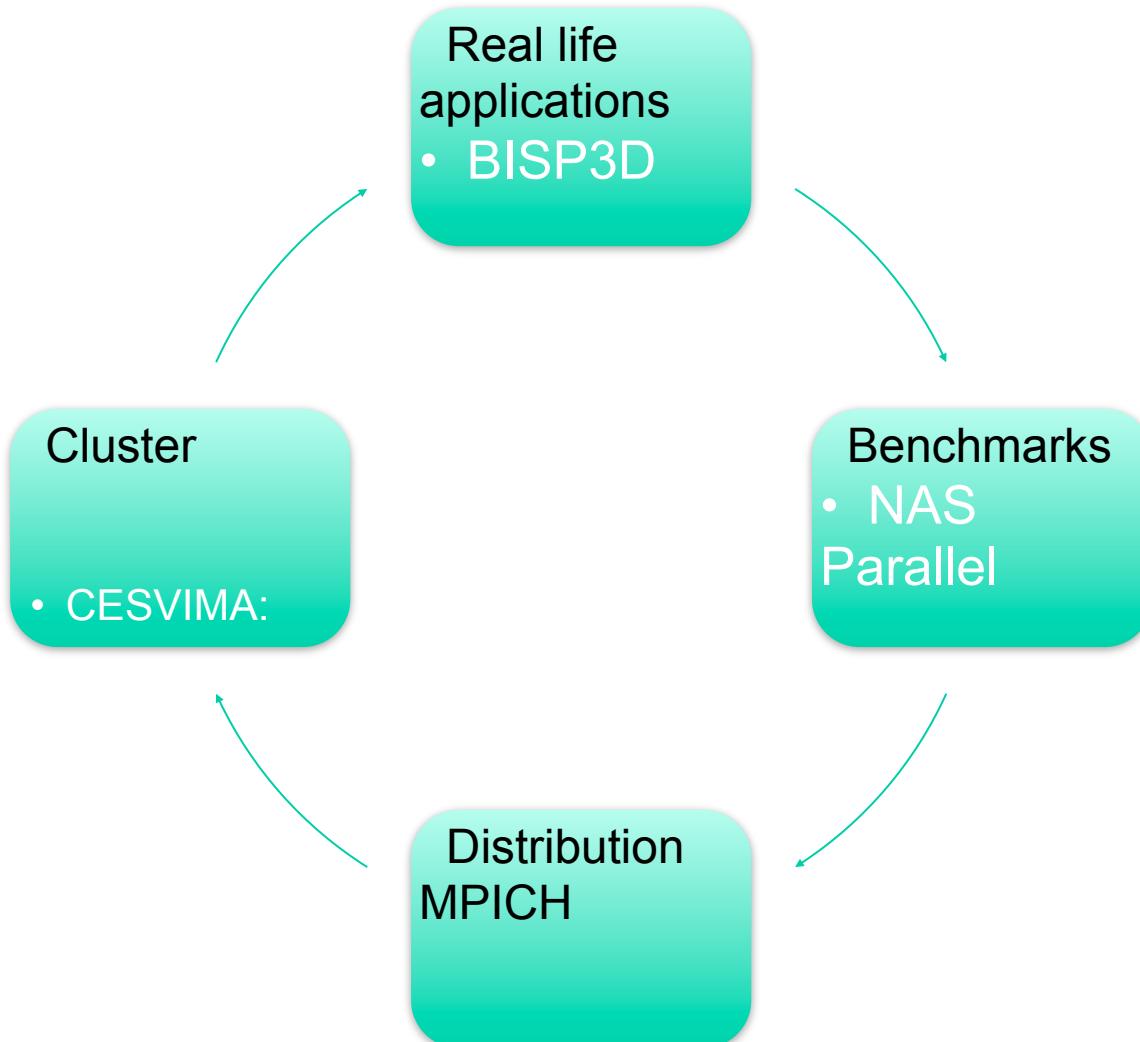
Summary

114

1. Problem description
2. Main Objectives
3. Strategies for improve the performance of collective I/O operation
4. Strategies for improve the performance of I/O and communication operations.
5. Evaluation tools

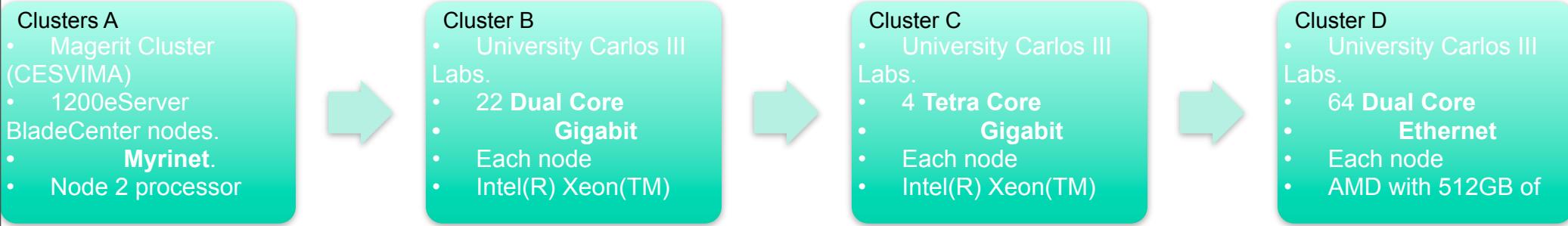
Phd. Thesis evaluation tools

115



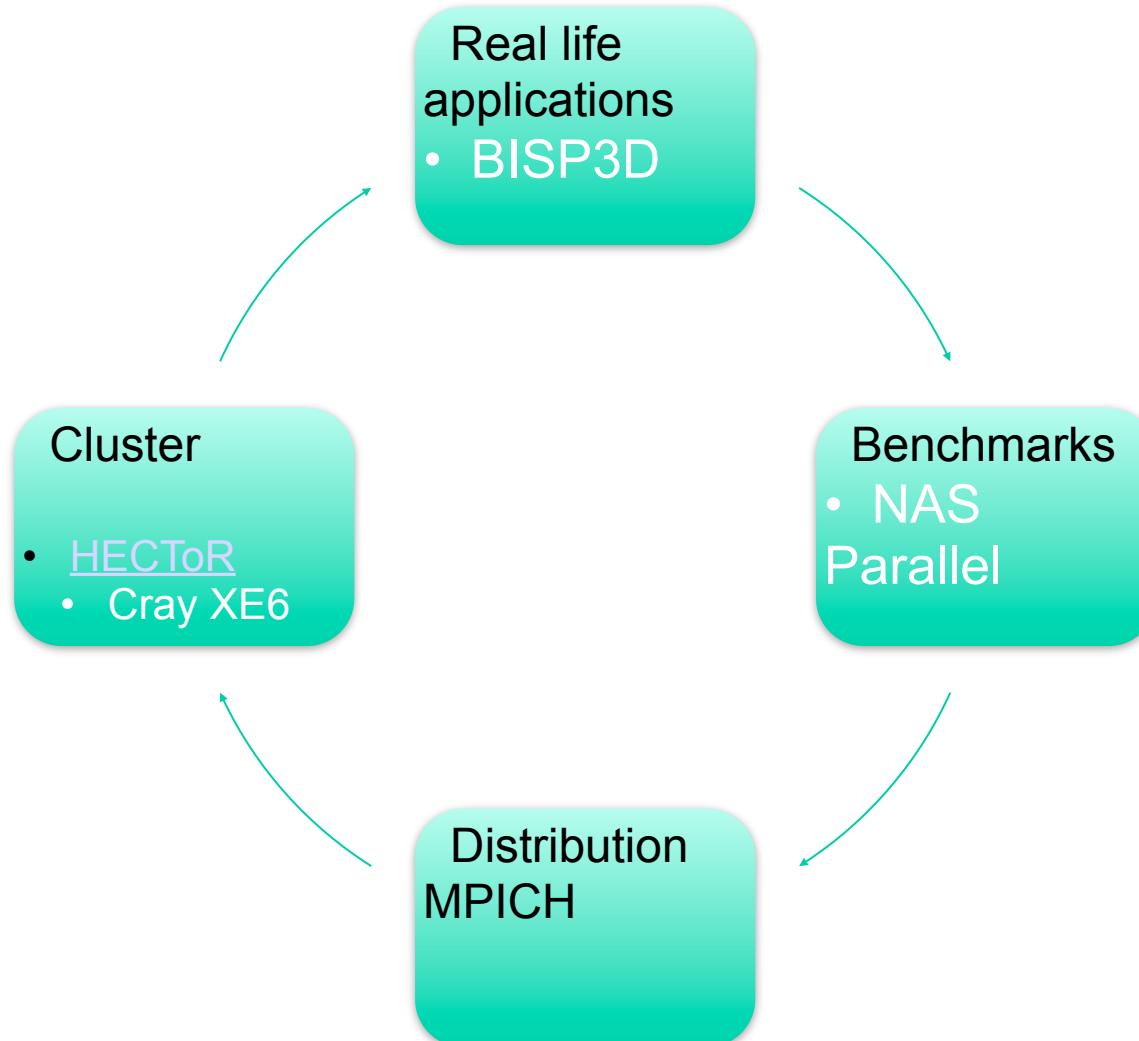
Clusters used in our experiments

116



New evaluation tool: HECToR

117



HECToR Hardware

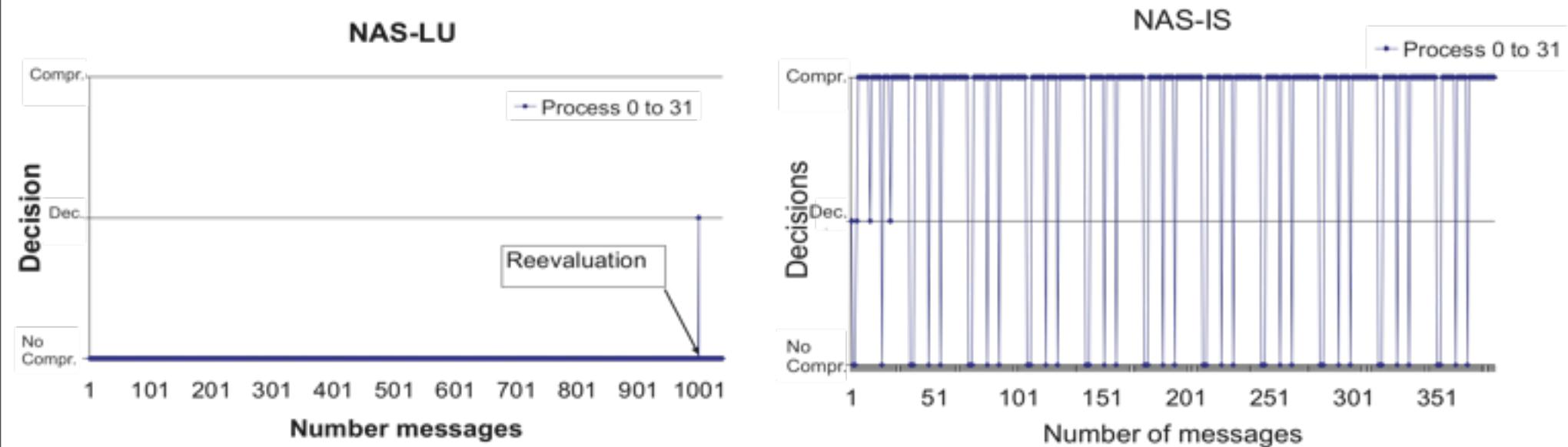
118

Clusters HECToR

- Cray XE6
- 1800 compute nodes
- **Gemini interconect.**
- Each Node 2 AMD 2.1 GHz 12 core Opteron processors → 44,544 cores
- Lustre

Runtime Adaptive Compression (RAS) Compression decisions.

119

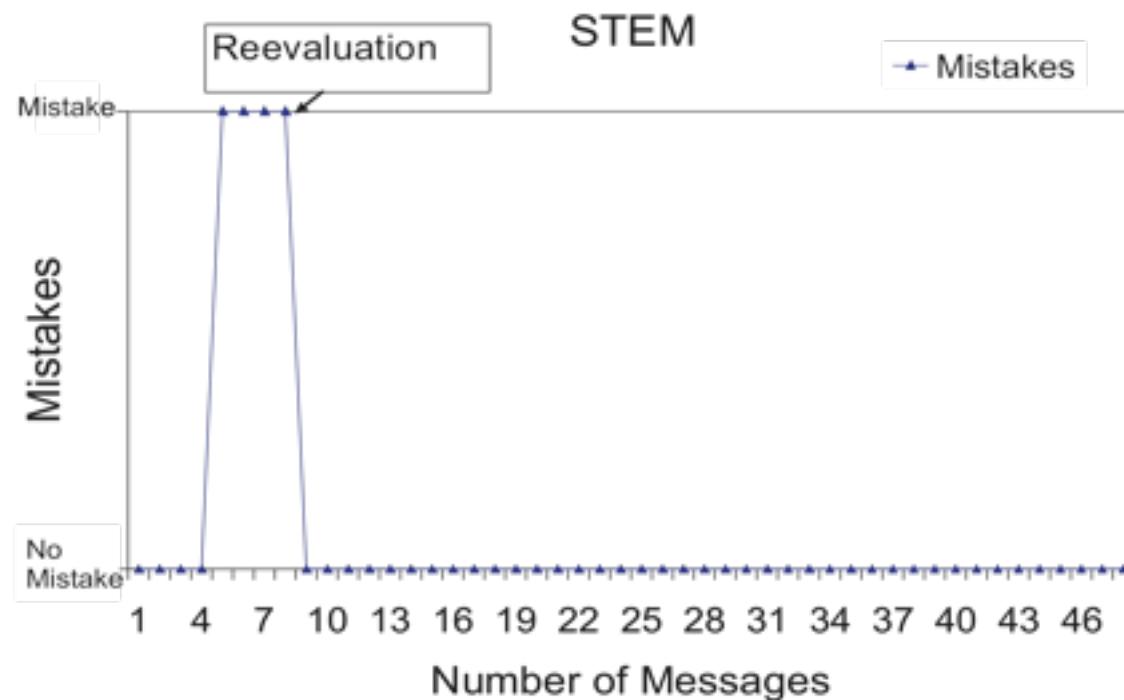


Benchmarks, which different kind of data, datatype and message size.

RAS adapts itself at run-time to the application behavior
and takes different compression decisions.

Reevaluation case with STEM

120



- Several mistake in a short period of time:
STEM. Initially selects RLE, and before the re-evaluation selects LZO.
- Number of messages sent without compression is higher than defined ratio. (NAS-LU)

Questions??

121

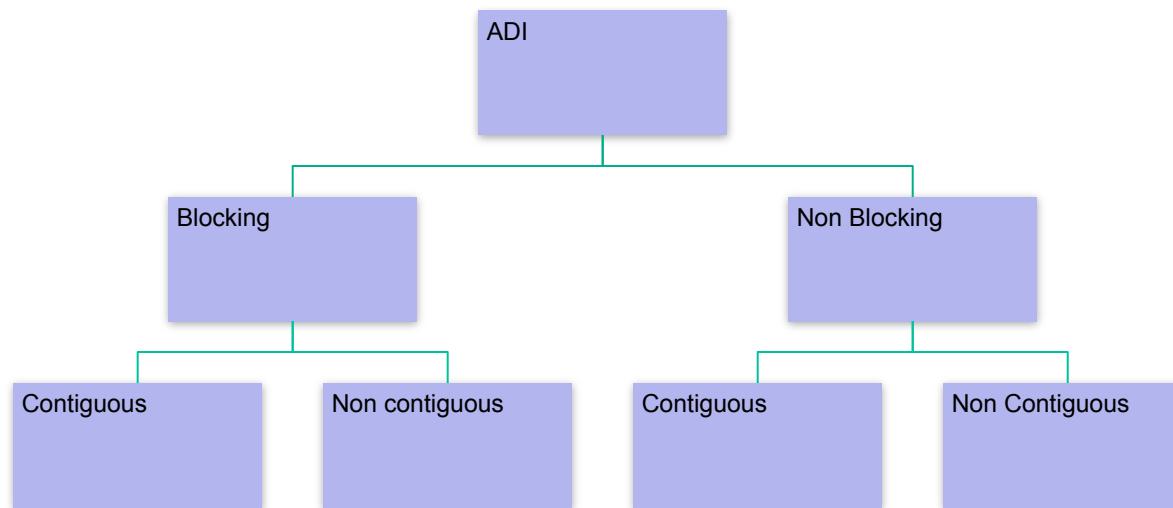
- Thanks!!
- <http://www.arcos.inf.uc3m.es/~rosaf/>

Dynamic optimization techniques to enhance scalability and performance of MPI-based applications on Multi-core cluster.

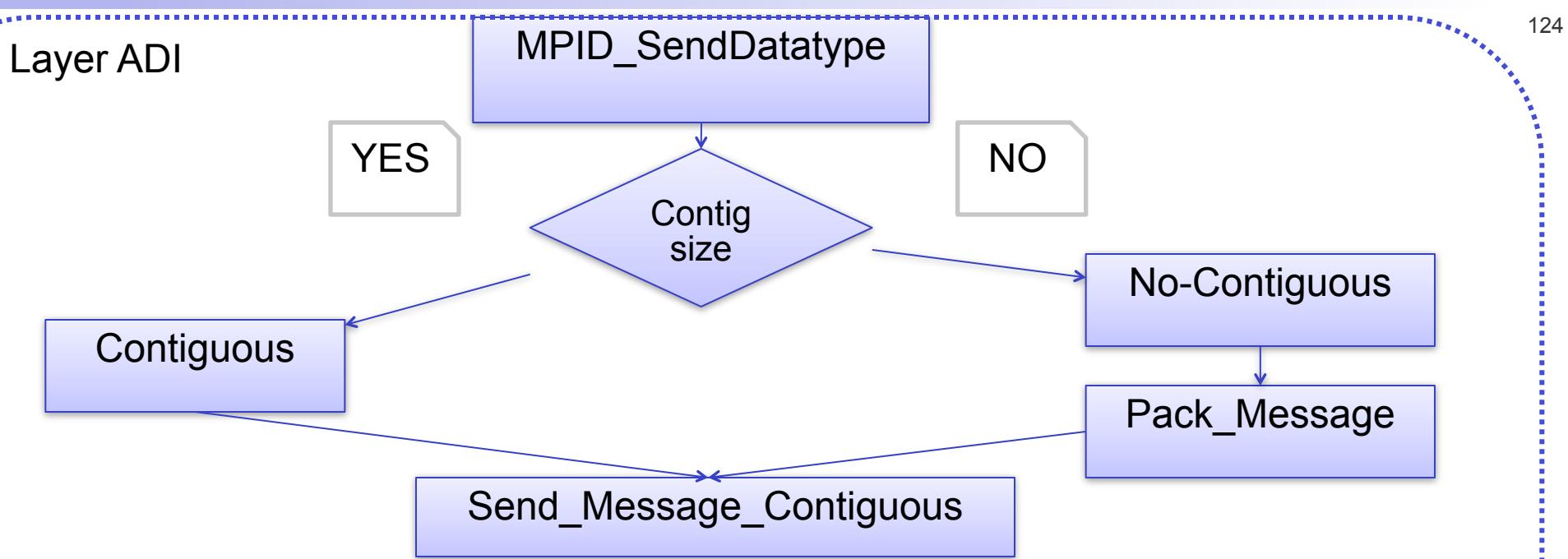
Author: Rosa Filgueira Vicente

ADI Layer structure

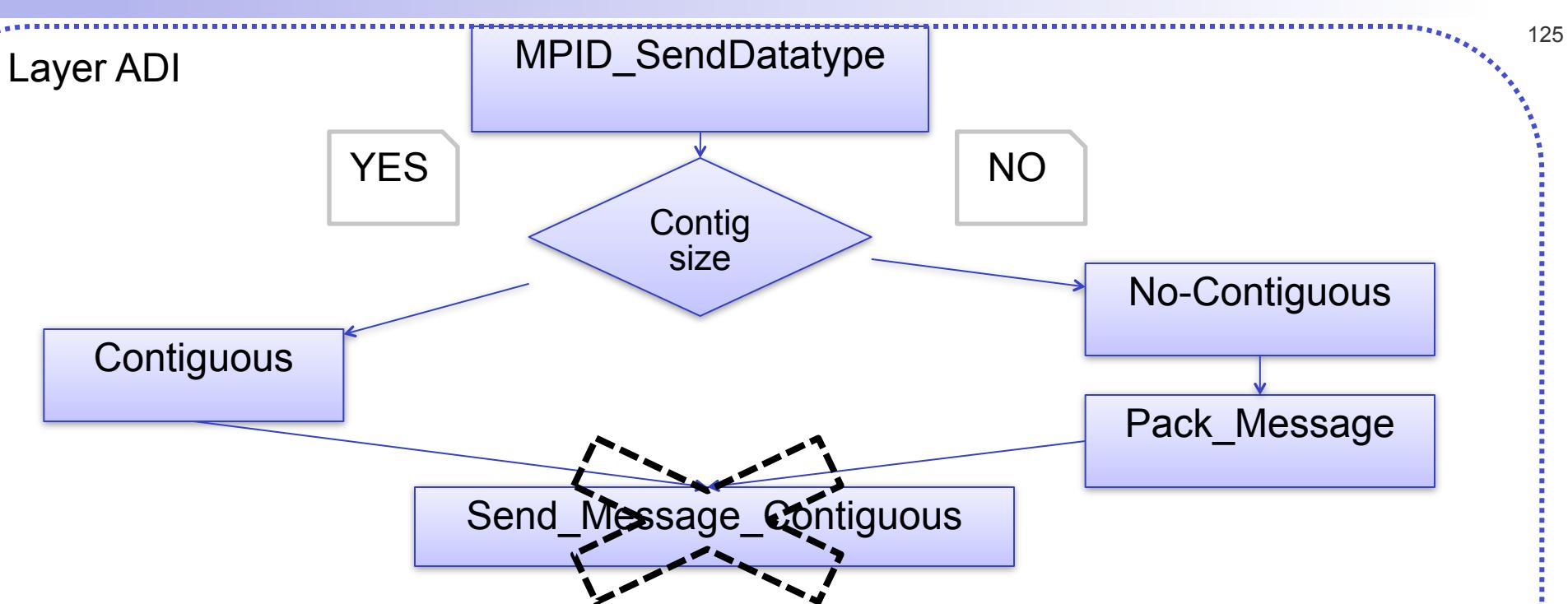
123



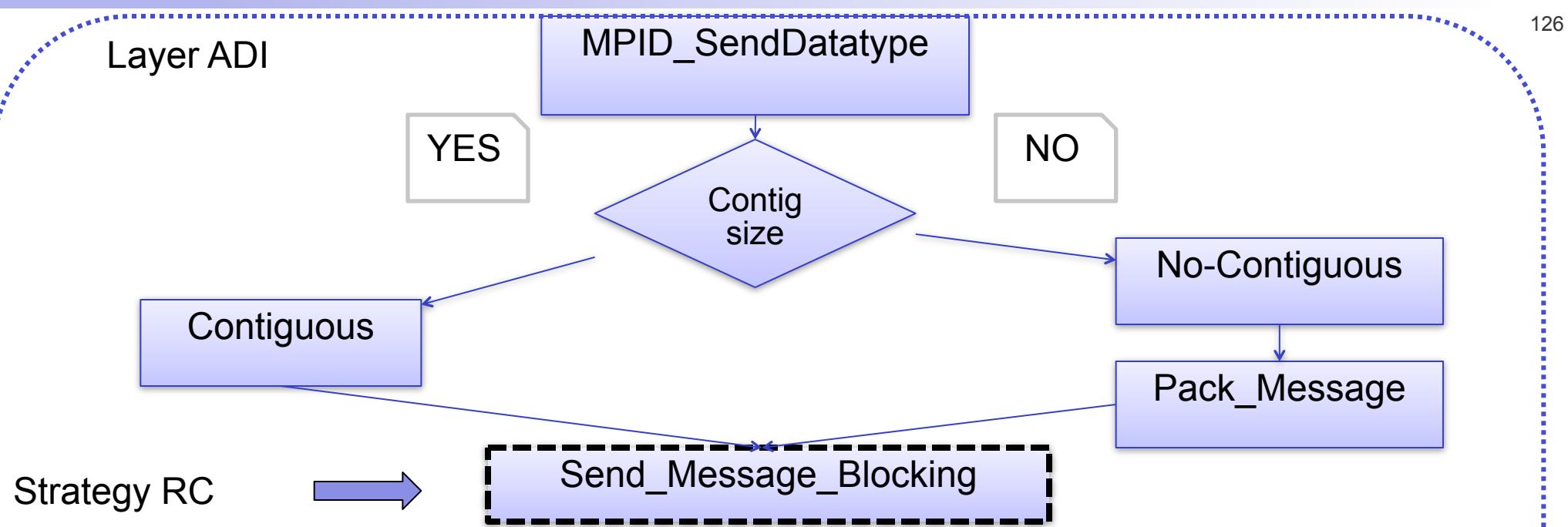
Send Message Blocking (MPICH)



Send Message Blocking (RC)



Send Message Blocking (RC)



Send Message Blocking (RC)

Layer ADI

MPID_SendDatatype

YES

NO

Contig
size

Contiguous

Send_Message_Blocking

No-Contiguous

Pack_Message

Strategy RC

YES

Size >
Threshold

NO

Compression

Add Header

Send_Message_Contiguous

Send Message Blocking (RC)

Layer ADI

MPID_SendDatatype

YES

NO

Contig
size

Contiguous

Send_Message_Blocking

No-Contiguous

Pack_Message

Strategy RC

YES

Size >
Threshold

NO

Compression

Add Header

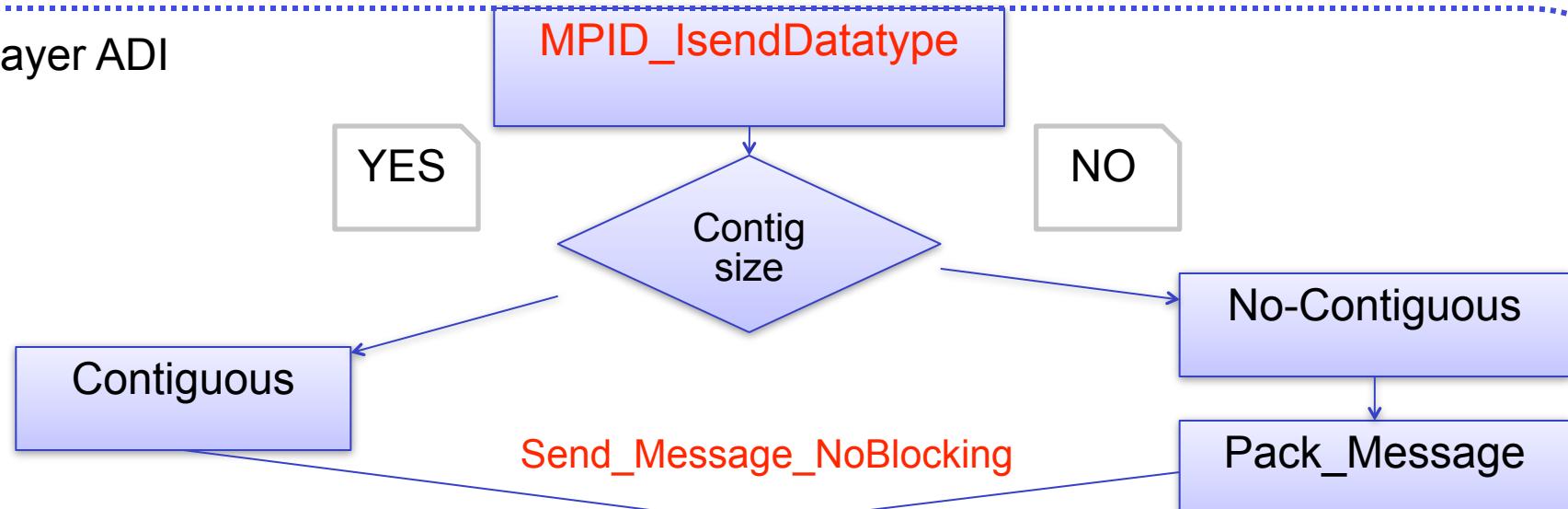
Send_Message_Contiguous

1. Select Compressor
2. Compress
3. Add Header

128

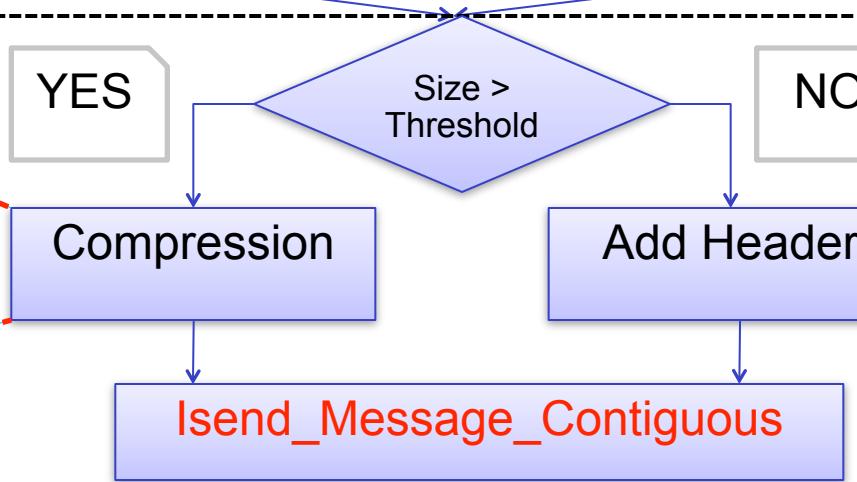
Send Message No Blocking (RC)

Layer ADI



Strategy RC

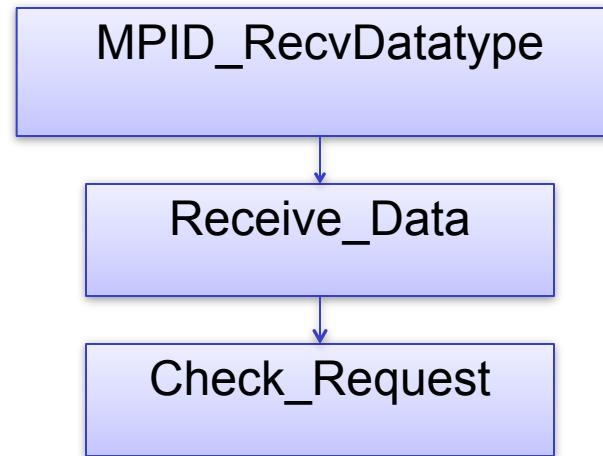
1. Select Compressor
2. Compress
3. Add Header



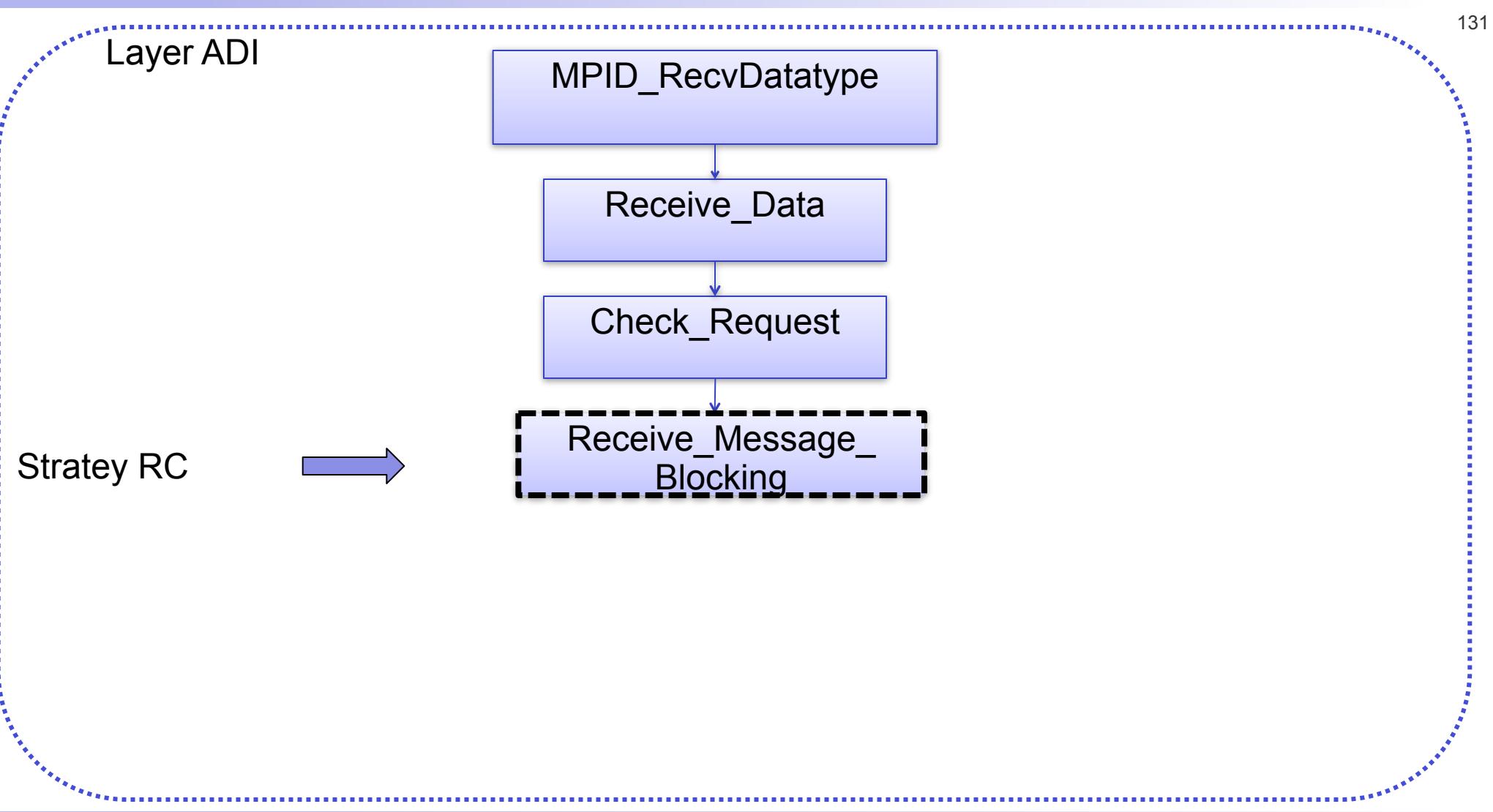
Receive message blocking (MPICH)

Layer ADI

130



Receive message blocking (RC)



Receive message blocking (MPICH)

Layer ADI

132

MPID_RecvDatatype

Receive_Data

Check_Request

Receive_Message_Blocking

Study Header

Decompressed? ?

NO

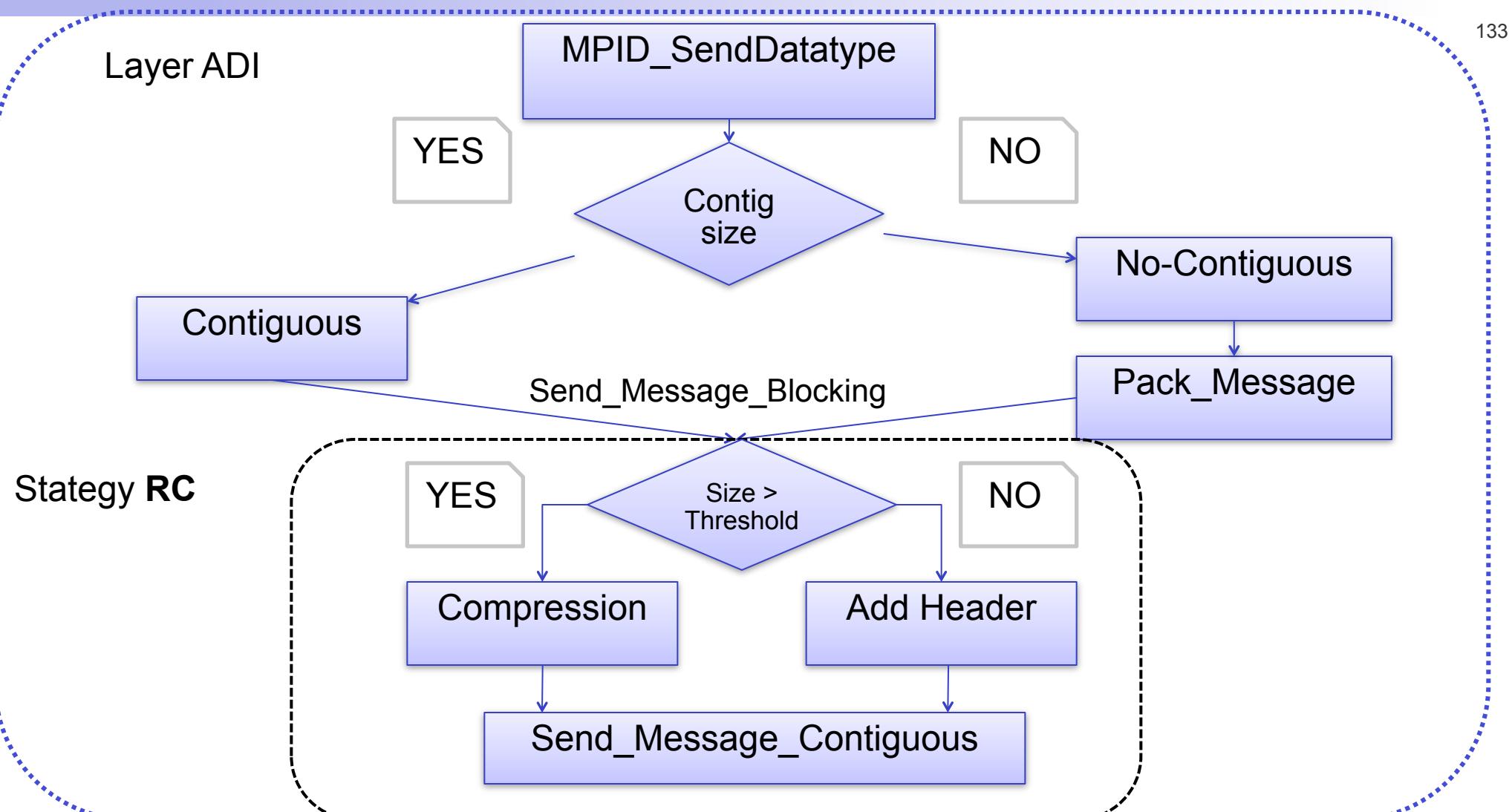
YES

Decompress_message

Strategy RC

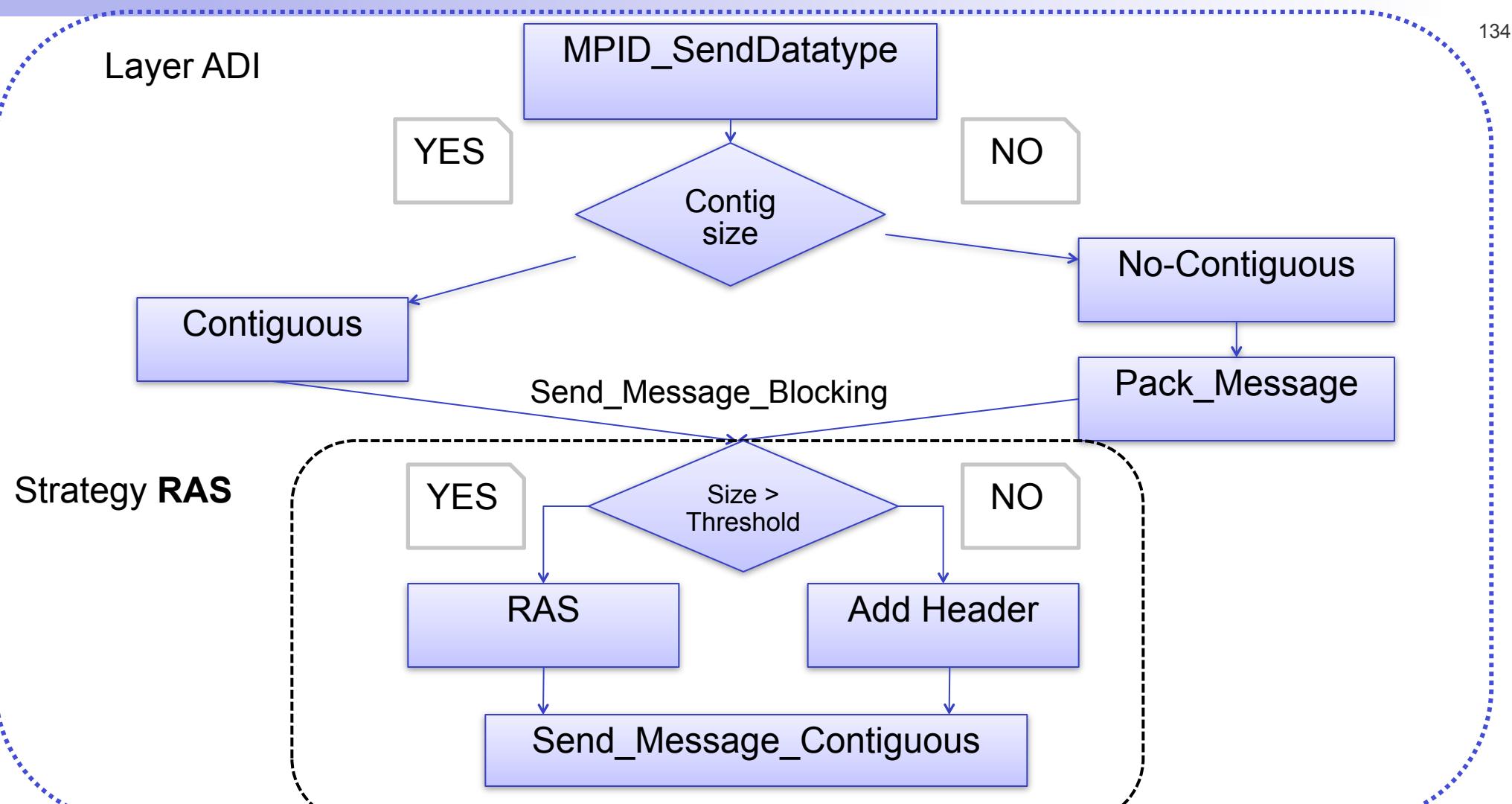
Send_Message_Blocking

Runtime Compression



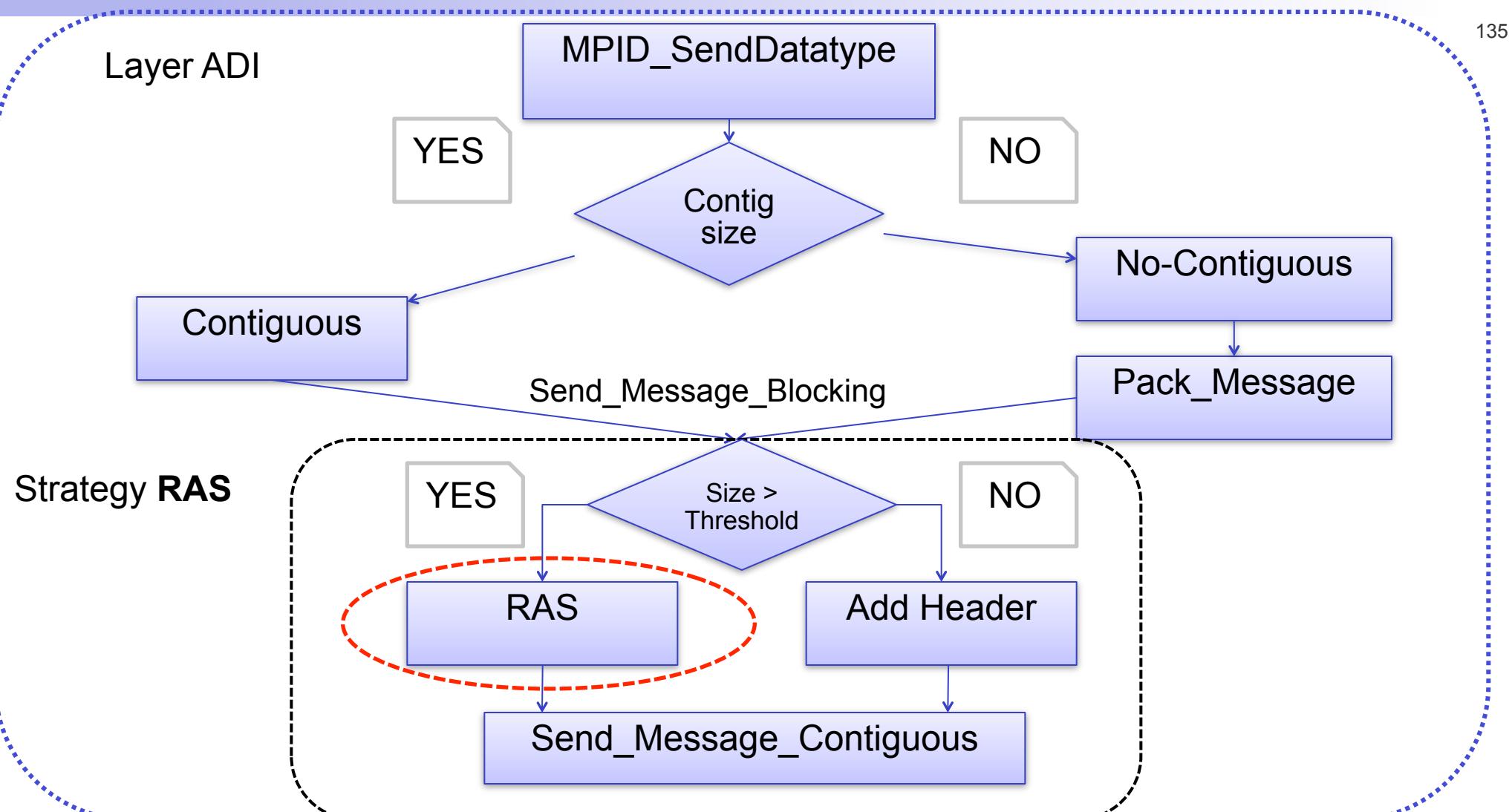
Send_Message_Blocking

Runtime Adaptive Compression



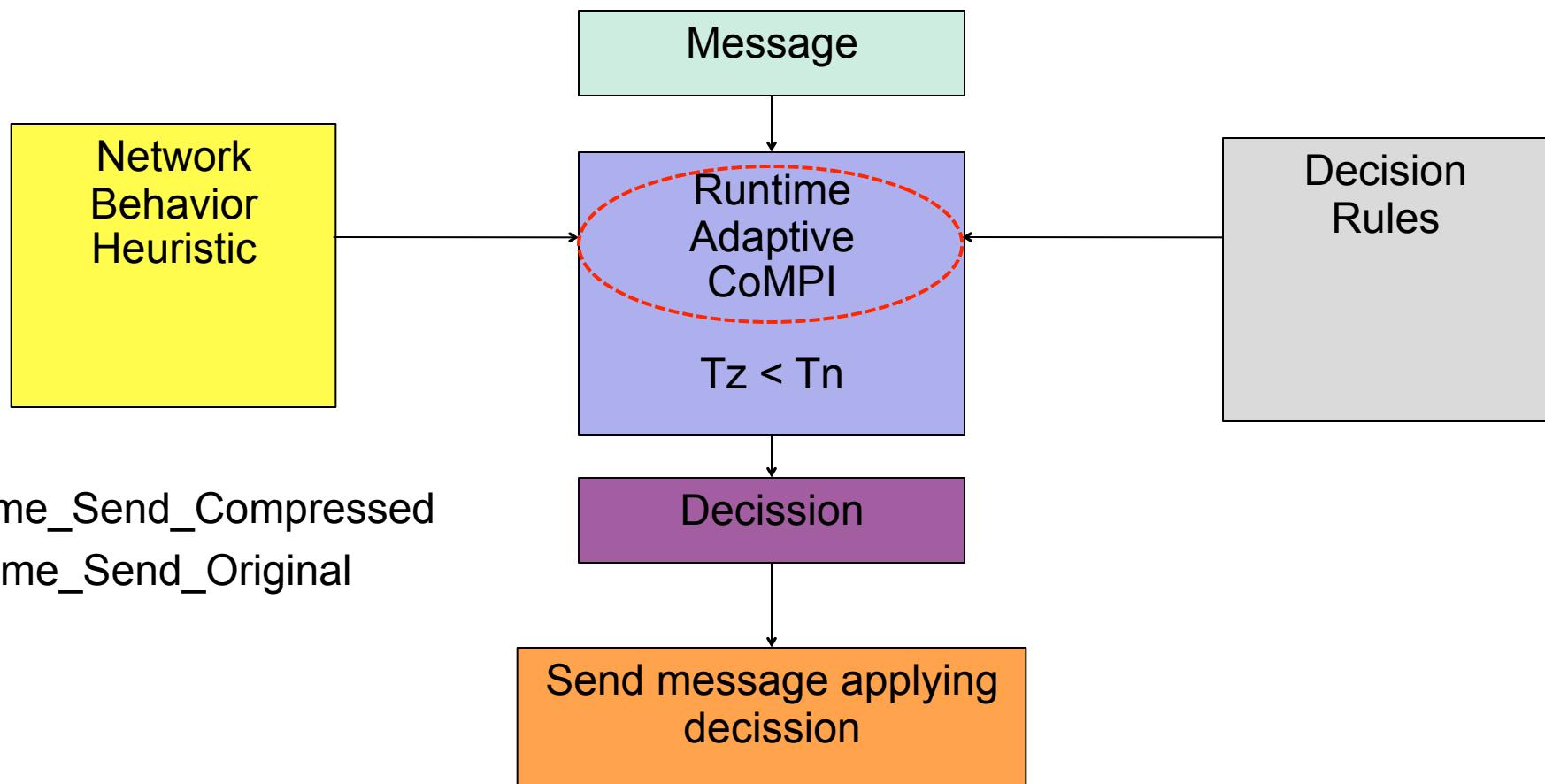
Send_Message_Blocking

Runtime Adaptive Compression



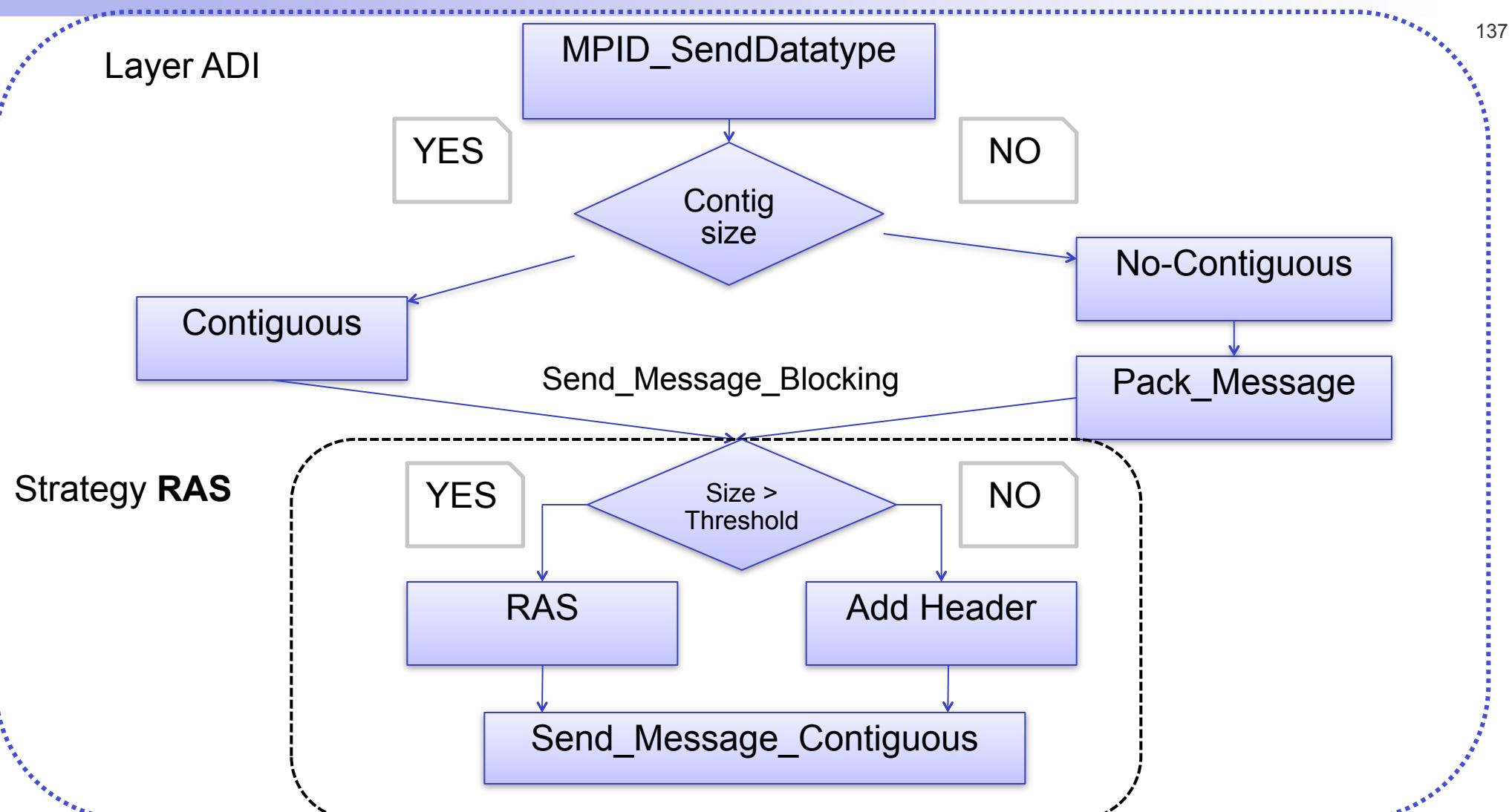
Guided Strategy

136



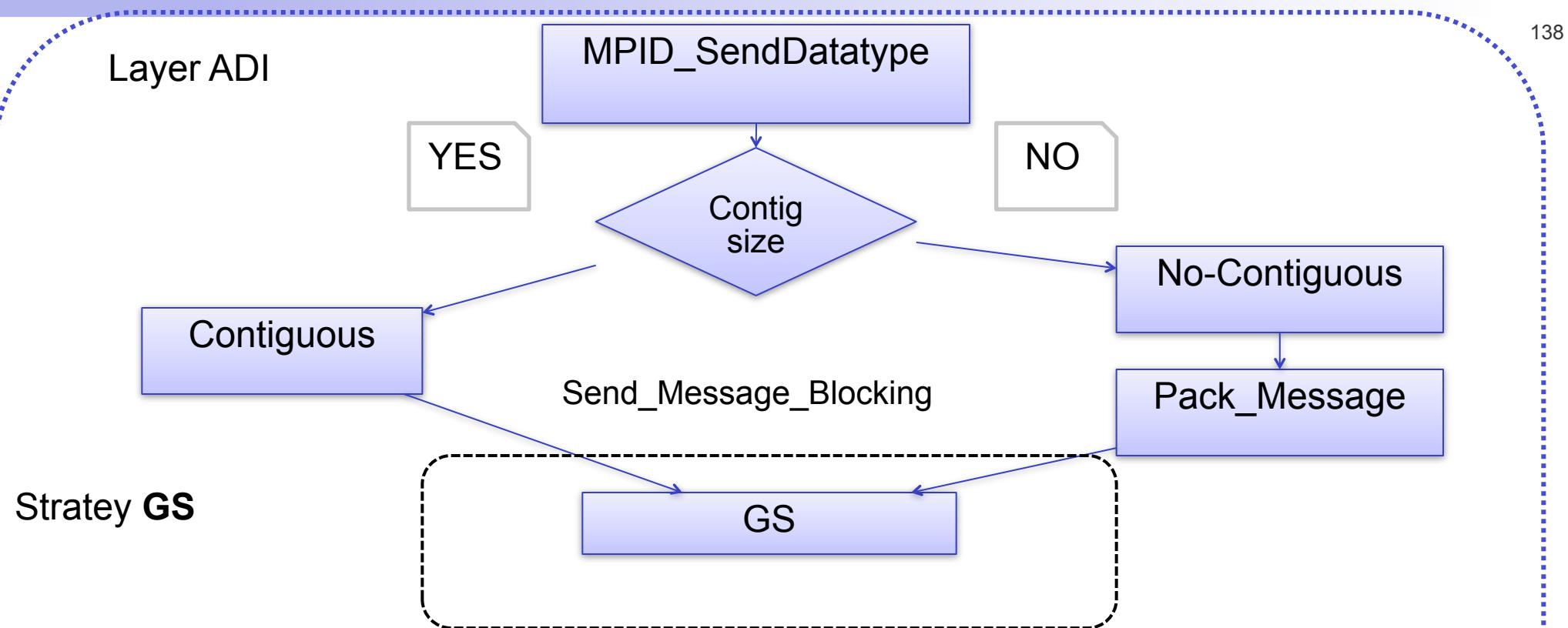
Send_Message_Blocking

Runtime Adaptive Compression



Send_Message_Blocking

Guided Strategy

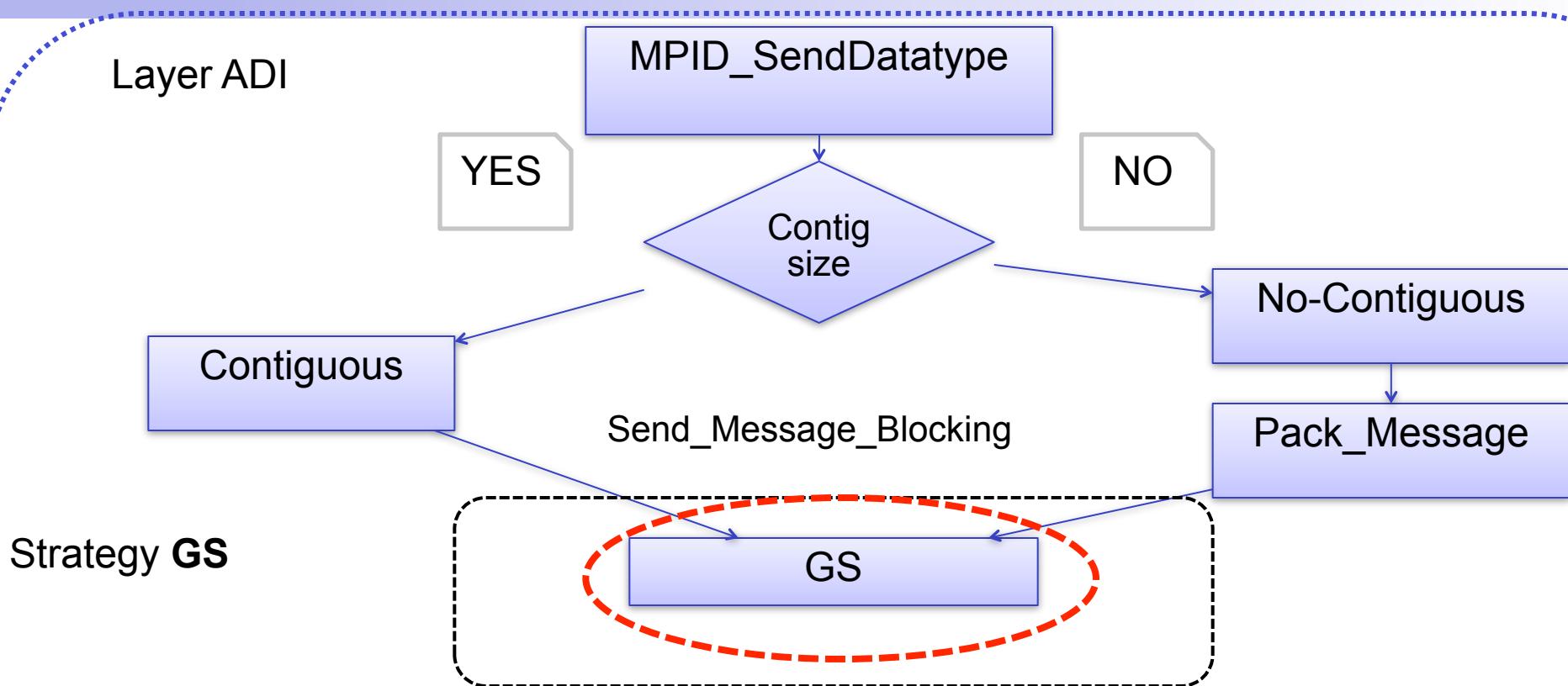




Send_Message_Blocking

Guided Strategy

139



Guided Strategy

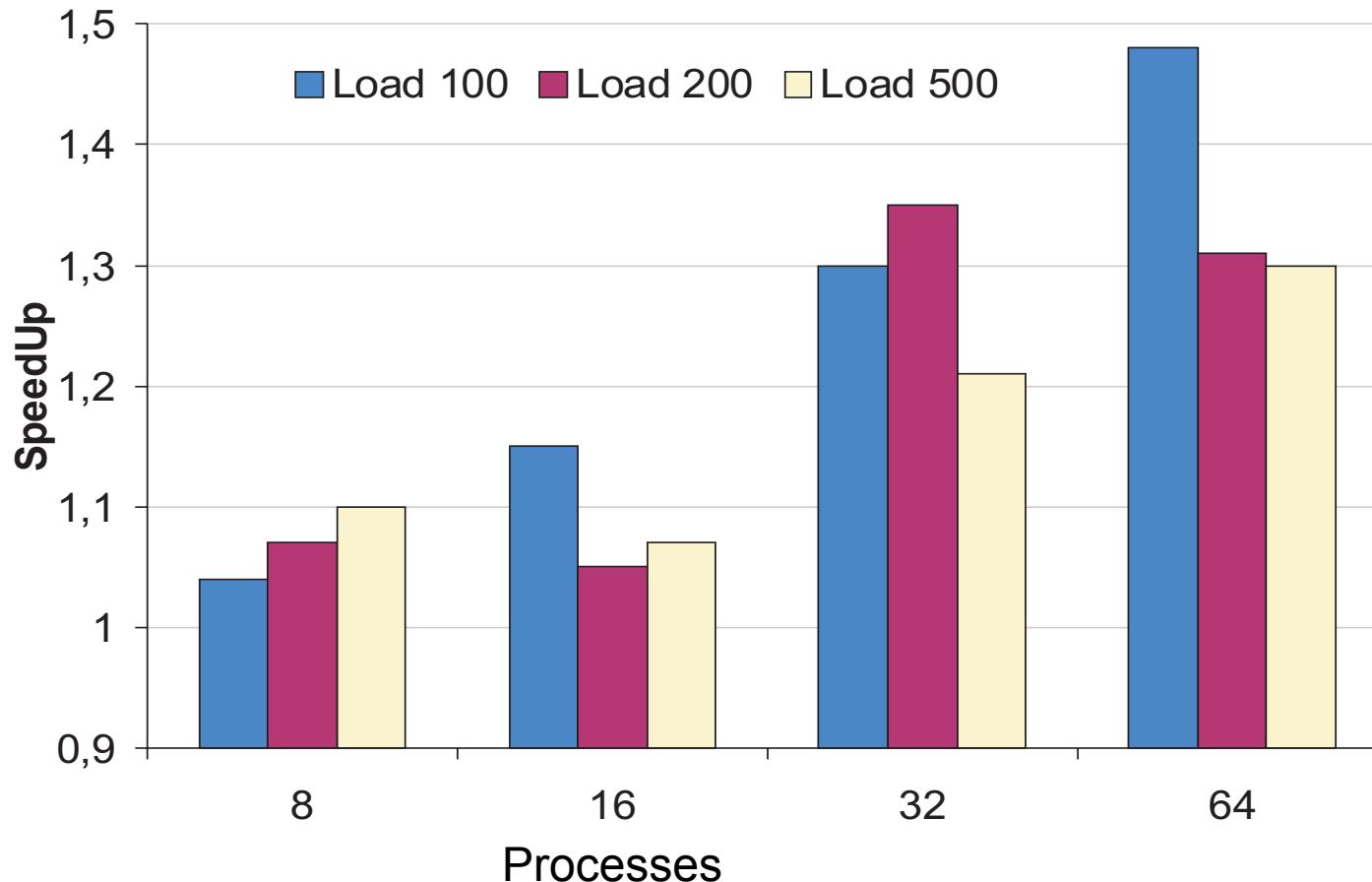
140



Evaluation LA_TwoPhase I/O BISP3D with cluster A

141

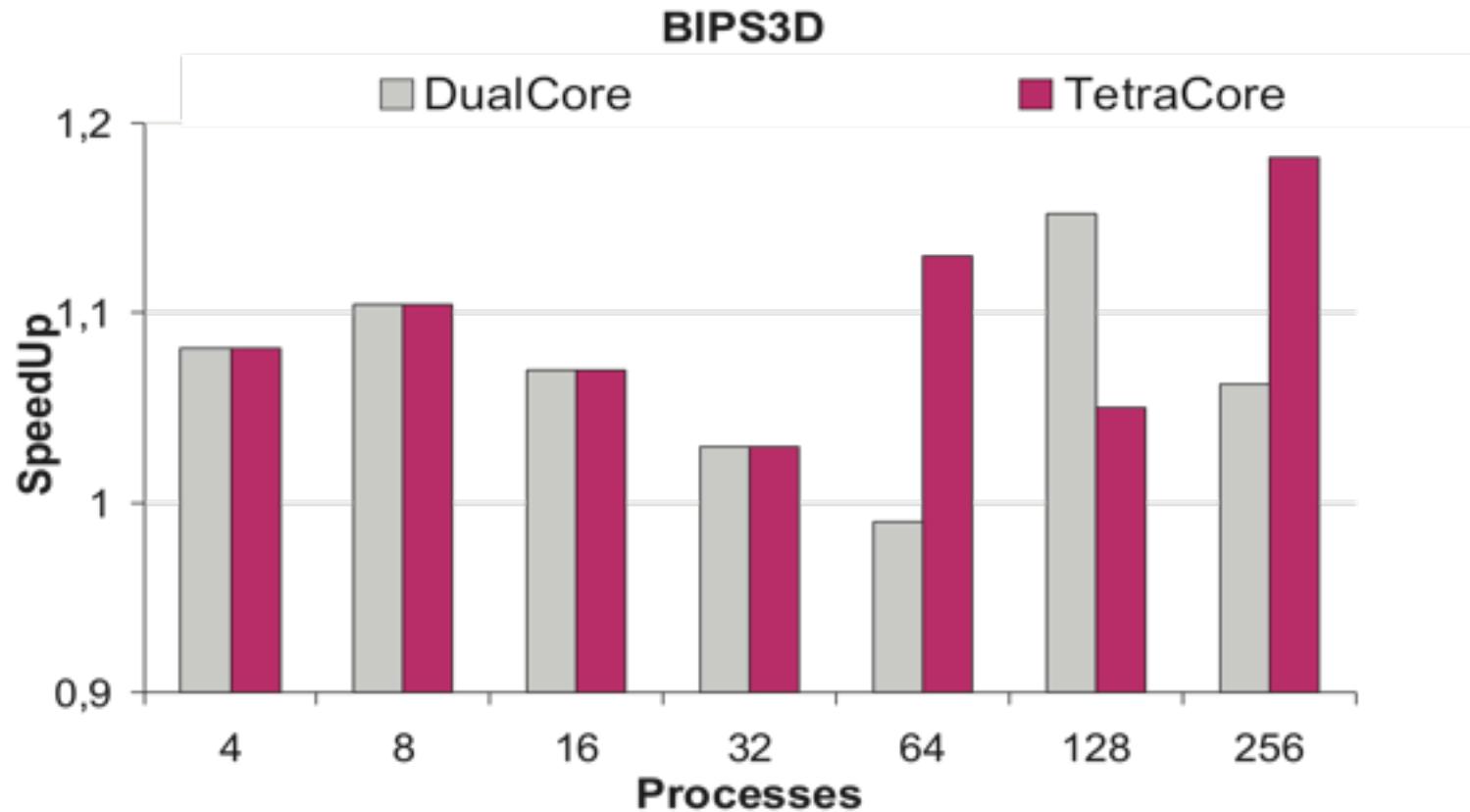
Mesh 1



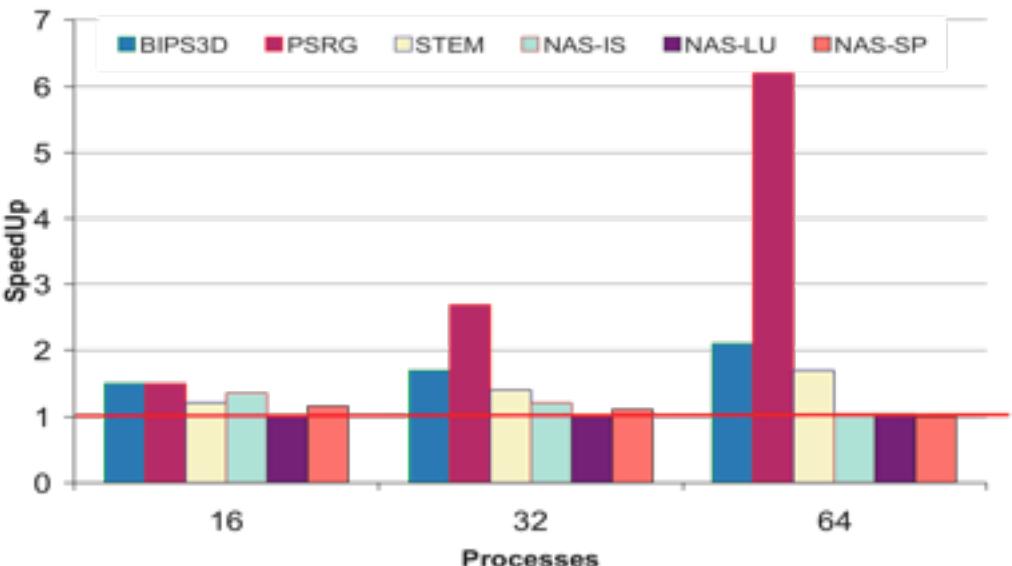
Evaluation LA_TwoPhase I/O BIPS3D with cluster B and cluster C

142

Mesh 4



Adaptive Compression Strategy (RAS)



143

RAS desactiva compresión cuando no obtiene beneficio (NAS).

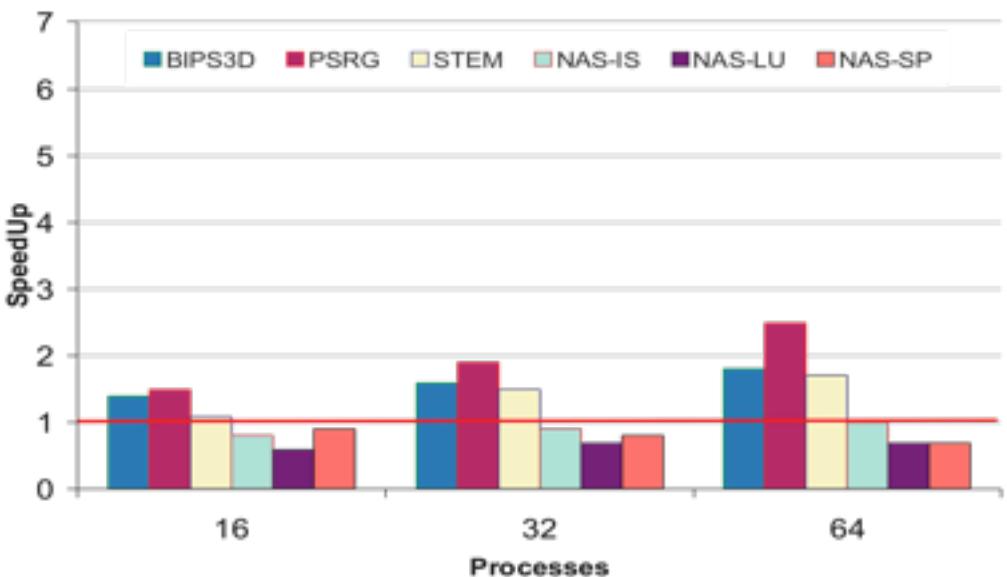
RAS elige el mejor algoritmo por mensaje.

Ejemplo PRSG utiliza LZO y RLE para comprimir en función de características de mensajes.

Mejor dejar tomar decisiones al algoritmo.

RAS no penaliza

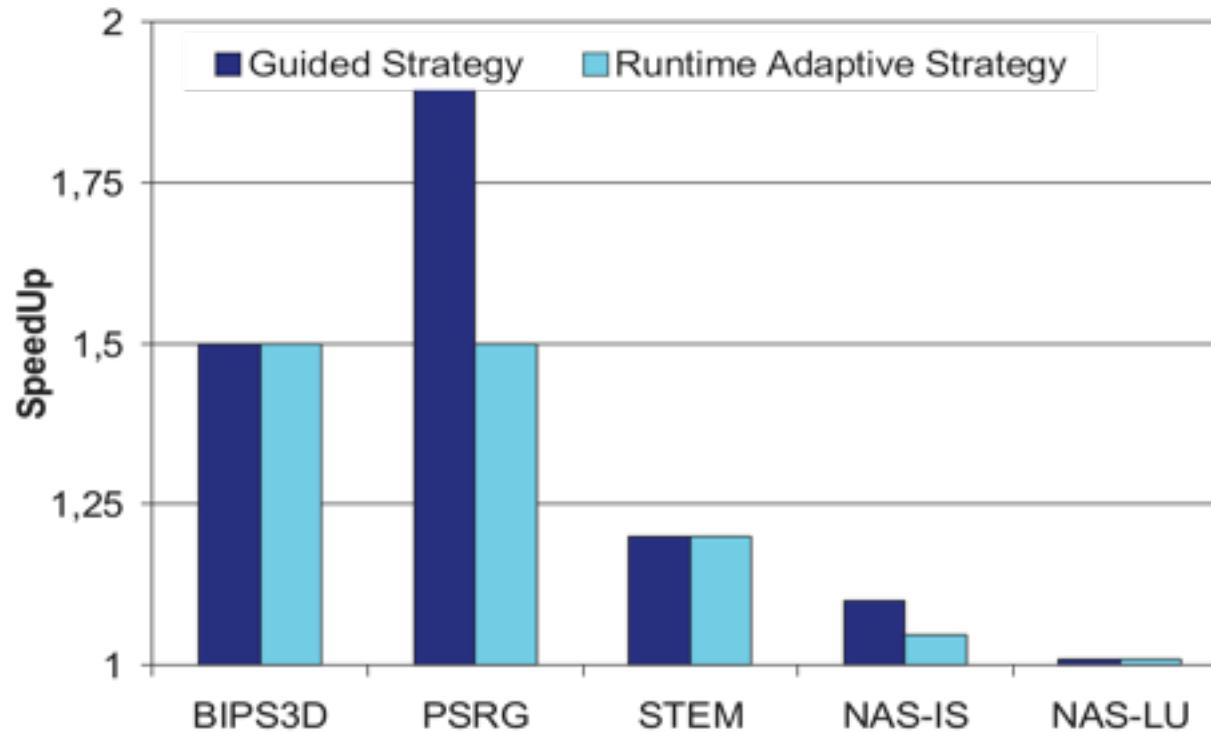
NO adaptive (RC)



Guided Strategy (GS)

144

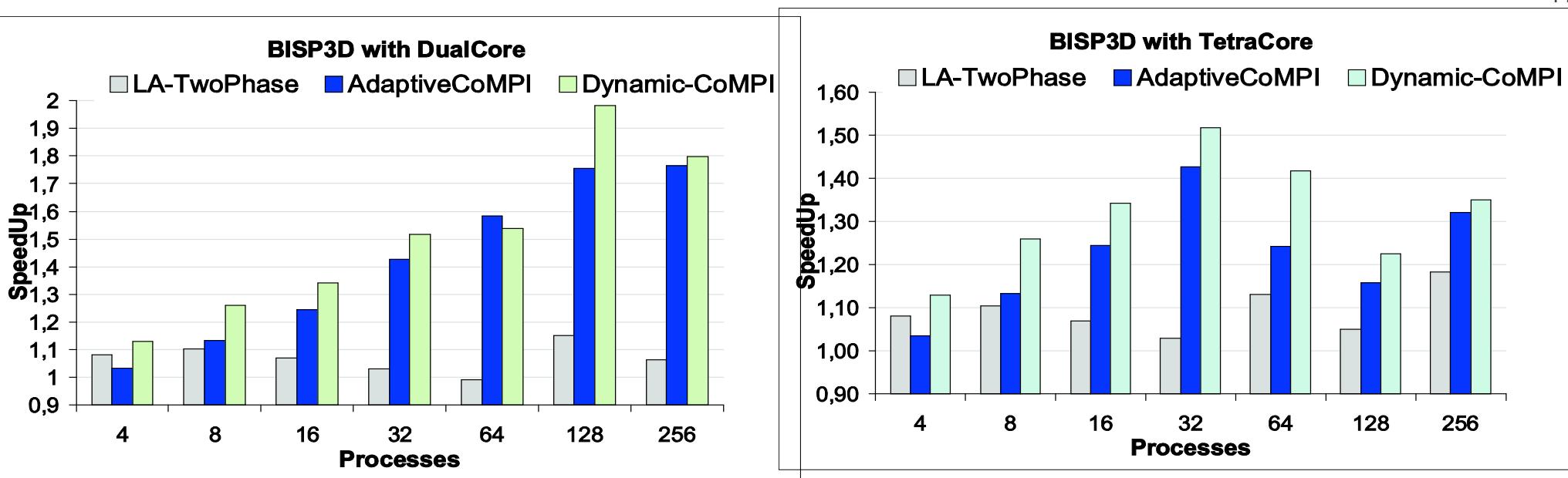
16 processes



RAS es tan buena estrategia como GS, ya que obtienen resultados similares.

Comparación de técnicas con BISP3D

145

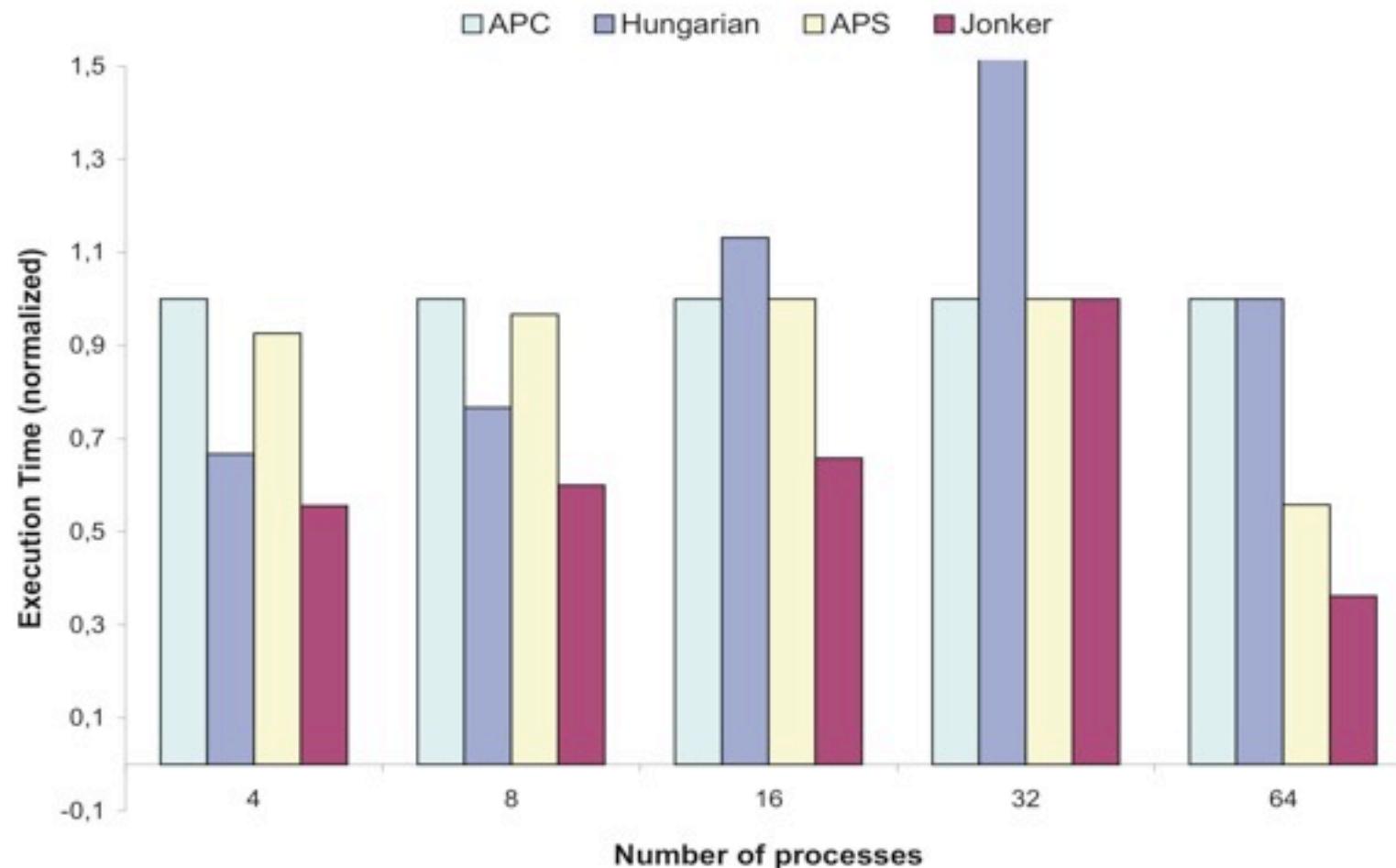


- Speedup LA-TwoPhase mayor en TetraCore → Mayor impacto el patrón de agregadores.
- Speedup Adaptive-CoMPI mayor en DualCore que TetraCore → Se mandan más mensajes comprimidos.
- Dynamic-CoMPI Speedup = LA-TwoPhase + Adaptive-CoMPI → Son complementarias ambas técnicas.

Linear Assignment Problem (I)

- LAP computes the optimal assignment of m items to n elements given an $m \times n$ cost matrix.
- Several algorithms have been developed for LAP:
 - Hungarian algorithm.
 - Jonker and Volgenant algorithm.
 - APC and APS Algorithms.
- All algorithms produce the same assignment.
- The difference is the time to compute the optimal allocation.

3.1 Linear Assignment Problem (II)



3.1 Linear Assignment Problem (II)

