



# The Lightweight Protocol CLIC: Performance of an MPI implementation on CLIC

Díaz, A. F.; Ortega, J.; Cañas, A.; Fernández, F.J.; Prieto, A.  
Departamento de Arquitectura y Tecnología de Computadores,  
Universidad de Granada

*Email:* [afdiaz@atc.ugr.es](mailto:afdiaz@atc.ugr.es), [julio@atc.ugr.es](mailto:julio@atc.ugr.es), [acanas@atc.ugr.es](mailto:acanas@atc.ugr.es),  
[jfernand@atc.ugr.es](mailto:jfernand@atc.ugr.es), [aprieto@ugr.es](mailto:aprieto@ugr.es)

# *Introduction*

- Usually Cluster Communications in Linux are based on TCP/IP
- Portability is reduced if we try to use NIC & Architecture optimizations  
(NICs can be obsoleted in a short time)
- OS can't offer full network performances:
  - Latency / Bandwidth.
  - Problems with multiple NICs. (Channel Bonding)
  - Heterogeneous Networks.
  - Reliable Broadcast/ Multicast.

## *What is CLIC ?*

- CLIC (Communication on Linux Cluster)
- Reliable Transport System optimized for Cluster Computing
- Developed on Linux (kernel module)
- Using OS resources:  
(*scheduler*, NIC drivers, kernel functions)

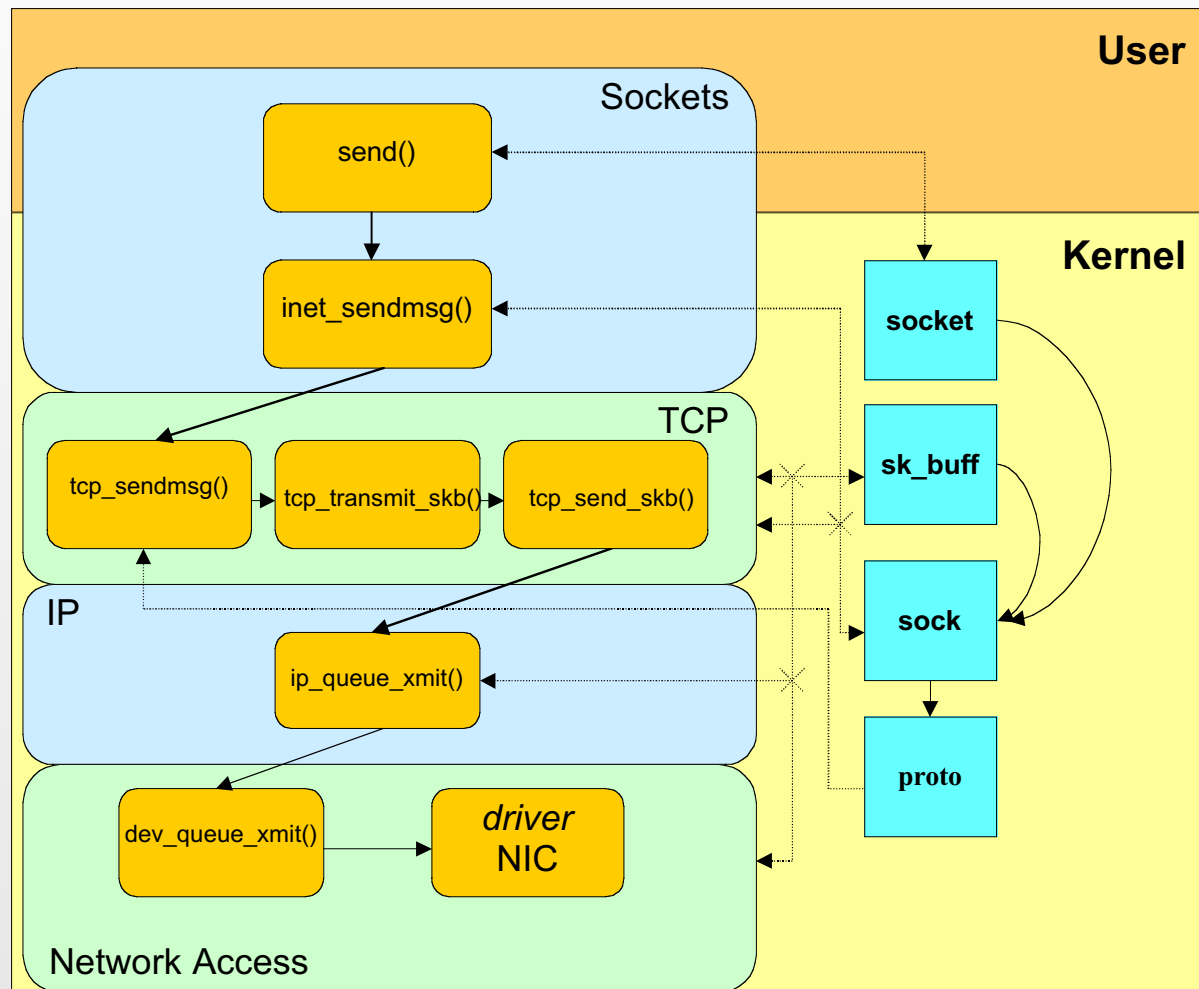
## *CLIC Features*

- Low *overhead SW*. Avoid TCP/IP stack.
- Own flow control & errors correcting mechanisms: Unicast, Multicast & Broadcast
- More Send types: *Broadcast, Asynchronous*
- Architecture Independent:
  - Autoconfiguration, Channel Bonding.*
- Heterogeneous Topologies. Routing
- Multiple Transport Characteristics.
- Multithread & SMP Support. Spinlocks & mutex

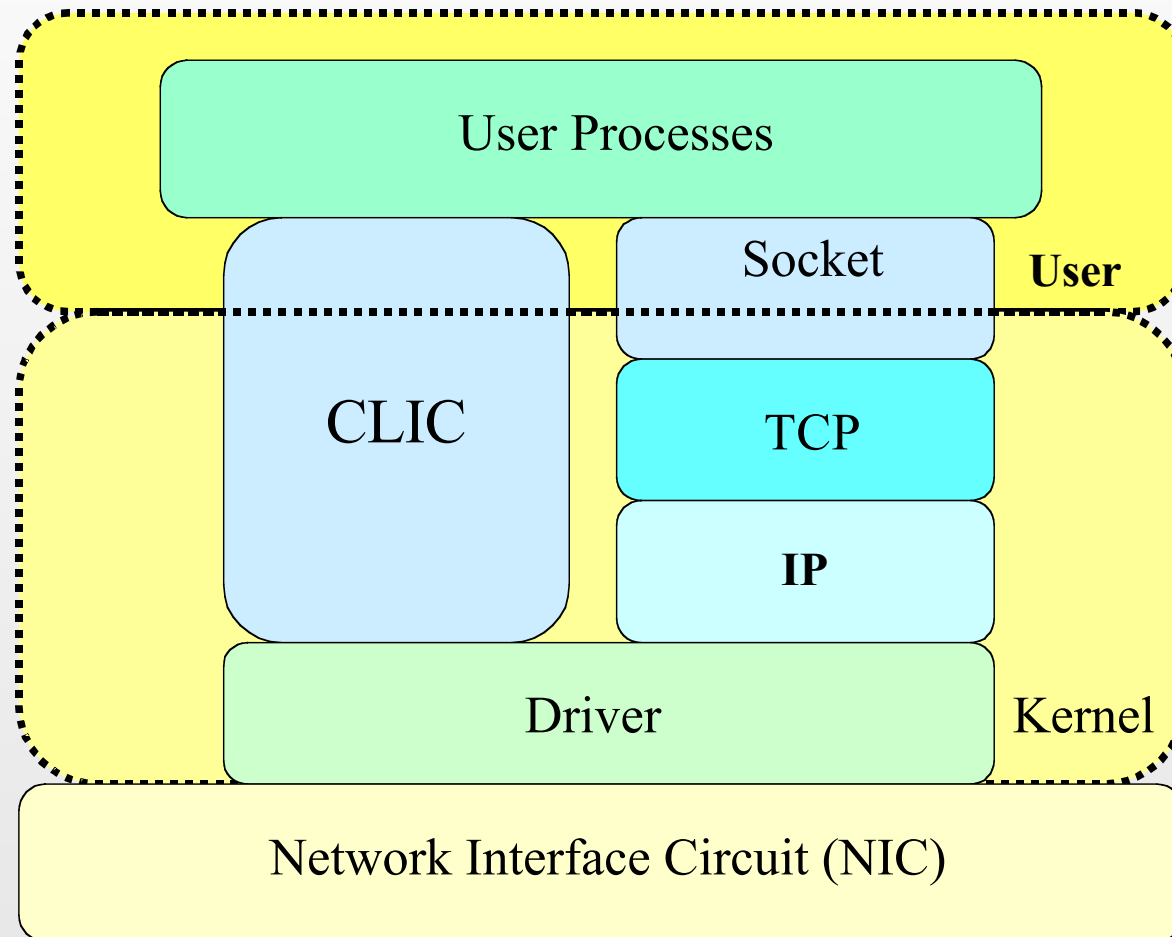
## *CLIC vs. other Systems*

Features	CLIC	GAMMA	VIA	Beowulf	Mad. II	U-NET
NIC & Arch. Indep.	YES	NO	YES	YES	YES	No
Threads, SMP	YES	NO	YES	YES	YES	YES
Multiprotocol Supp.	YES	NO	YES	Op.Syst.	YES	NO
Secure System	YES	NO	NO	NO	NO	NO
Uni. Flow-Control	YES	YES	NO	Op.Syst.	YES	NO
Broadcast Flow-Ctl.	YES	NO	NO	NO	NO	NO
Packet Recovery	YES	NO	NO	Op.Syst.	YES	NO
Channel Bonding	YES	NO	NO	YES	NO	NO
Sincr. & Asynchr.	YES	NO	NO	NO	YES	NO
Autoconfiguration	YES	NO	NO	NO	?	NO
Routing	YES	NO	NO	Op.Syst.	Het. No enr	NO
Send. same Host	YES	NO	YES	Op.Syst.	YES	NO

# IP Stack

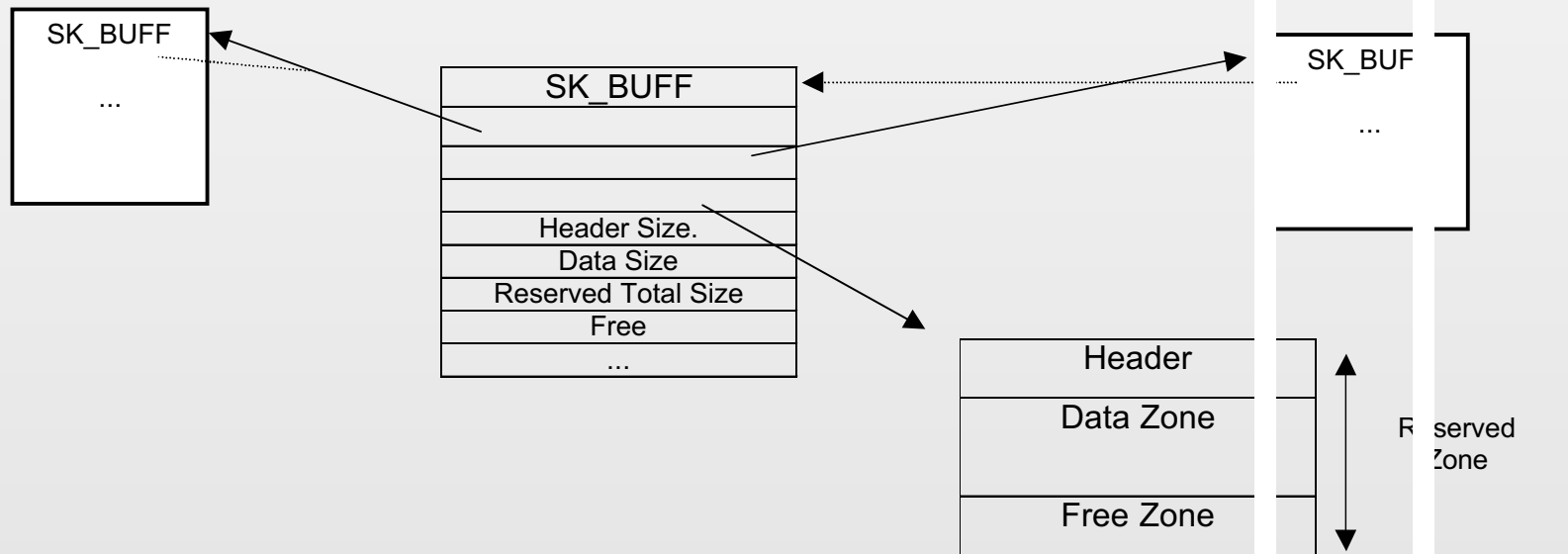


# *Avoiding IP Stack ...*



*Sk\_buff*  $\leftrightarrow$  NIC

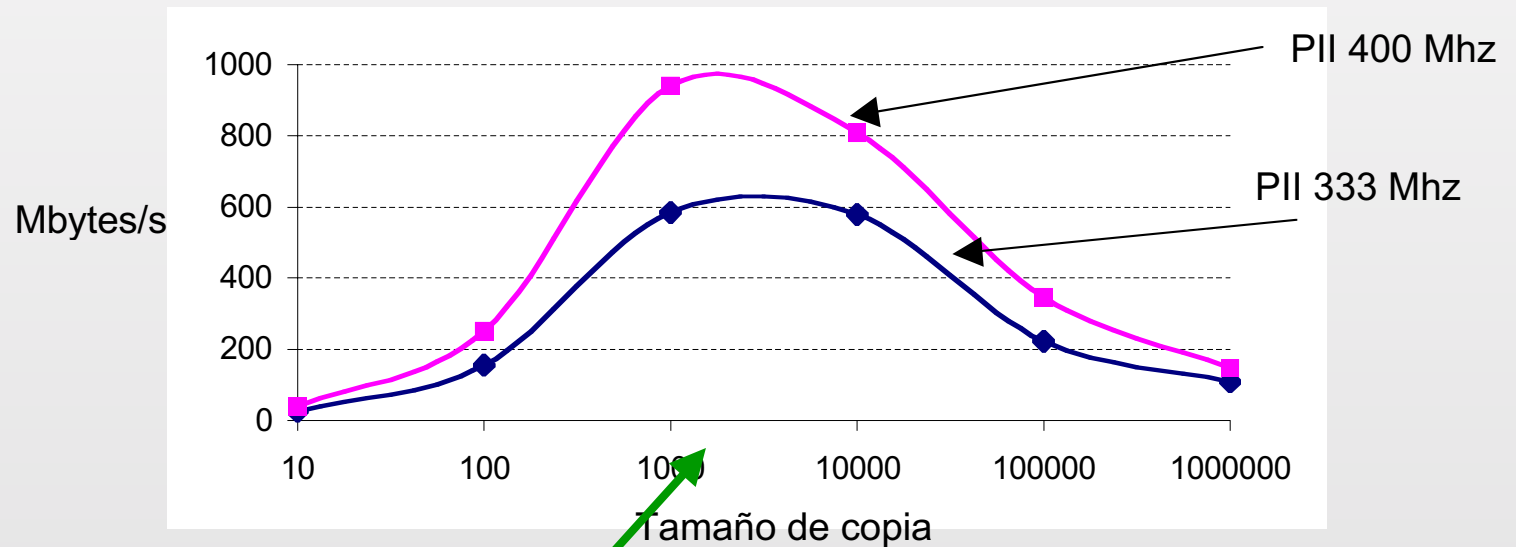
- Information exchange with NIC.
- sk\_buff* needs memory copy *sk\_buff* but:  
architecture independence, pipeline in send process





## *Bandwidth & Memory Cost*

- Minimize number of memory copys, better.
- Using 0-copy is not allways better performances.
- Bottleneck is in the network, not in the memory.



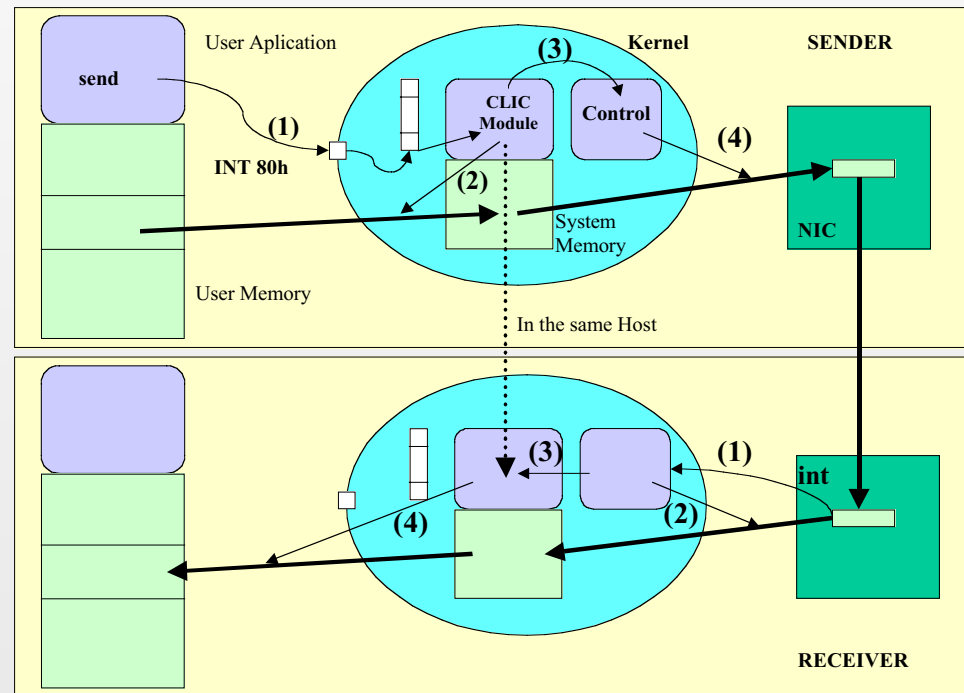
# How CLIC works ...

## SENDER

- (1) System Call
- (2) Copying to *sk\_buf*
- (3) Flow Control
- (4) Sending to NIC (*driver*)

## RECEIVER

- (1) Interrupt Routine
- (2) *Bottom Halves*
- (3) CLIC control & proc.
- (4) Send user's memory



# *Times in CLIC*

Compose & send NIC

Network (NIC out)

Network (Switch Out)

Interrupt Routine

BH + CLIC

Context Change

1500 bytes (1)

1500 bytes (2)

1500 bytes (1)

1500 bytes(2)

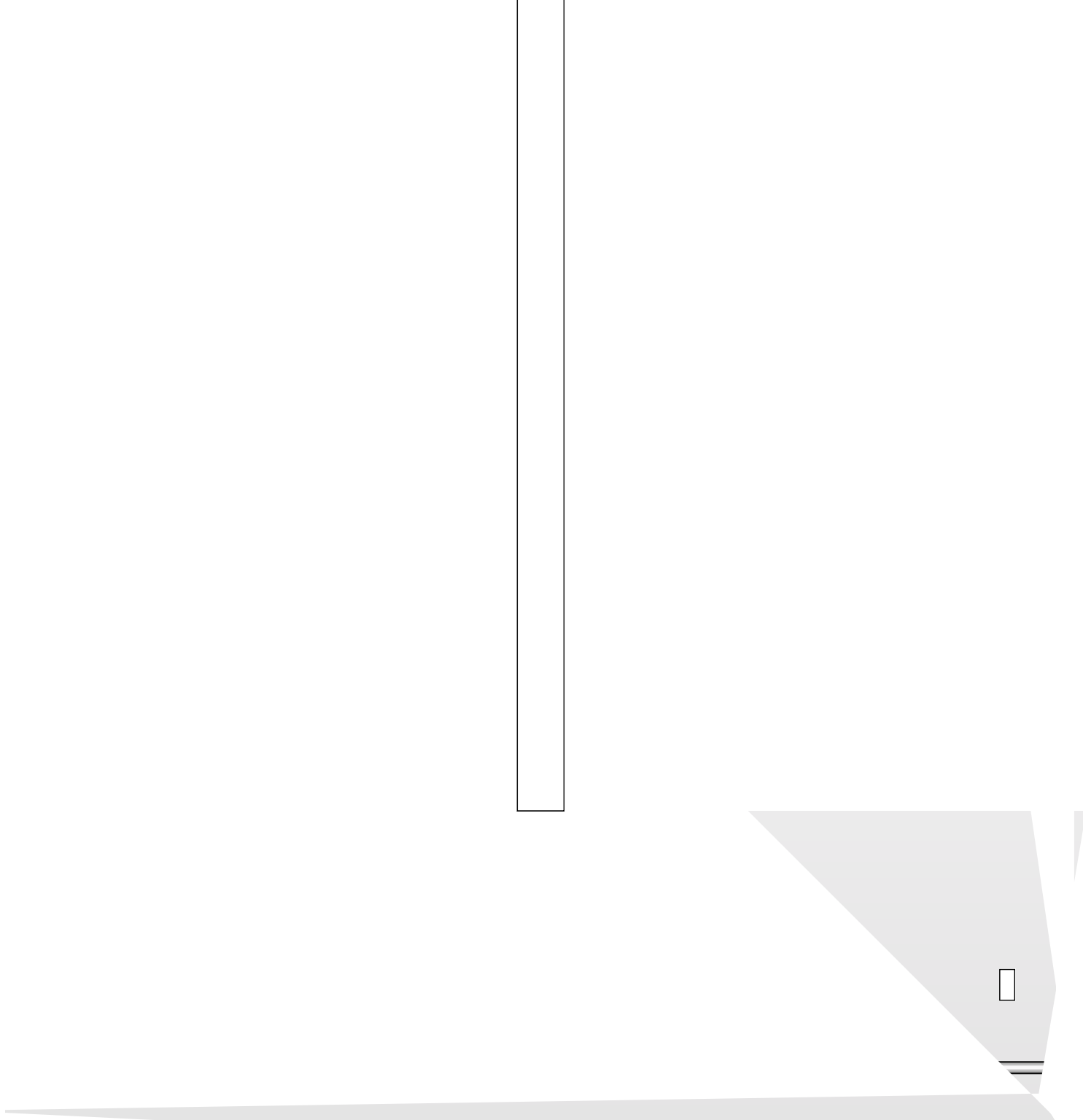
124  $\mu$ s

A

Snd.

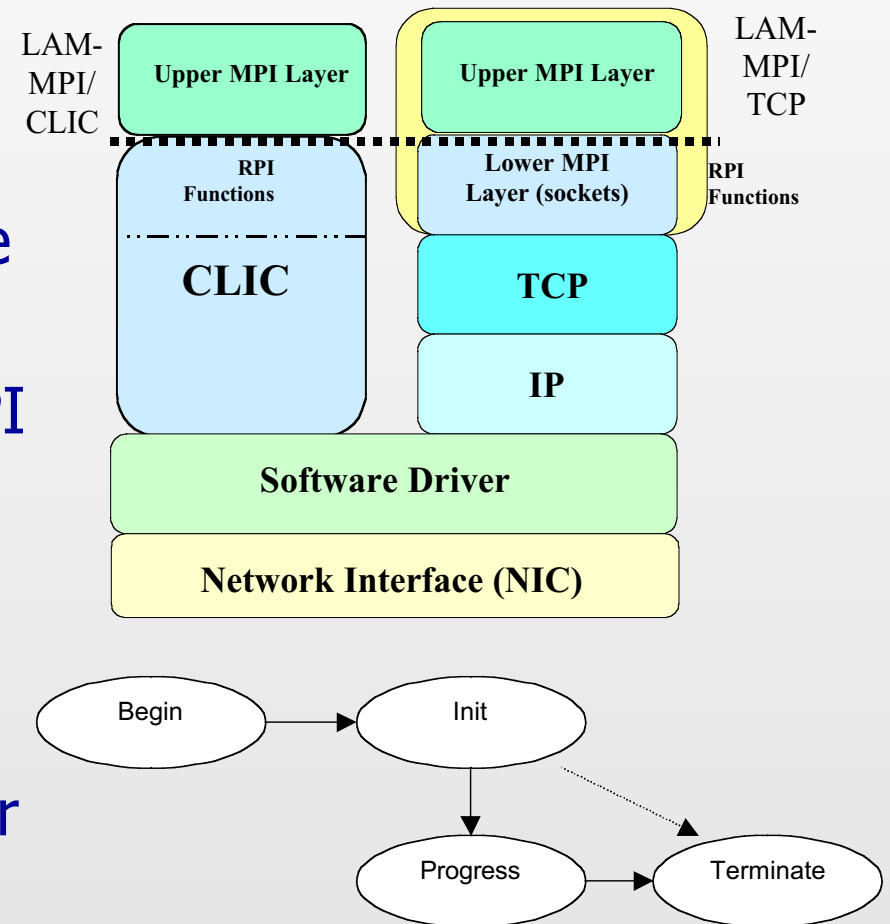
## *Other Optimizing OS CLIC...*

- SEND:
  - Memory management (*sk\_buffer* reuse).
  - Open communication ASAP.
  - Direct NIC Send.
- Receive:
  - Processing packets
- Both:
  - Flow Control optimized for full-speed.
  - Selective Lost packet re-transmissions.
  - Low overhead.



# LAM-MPI over CLIC

- RPI Functions. C2C Model.
- Multiple process, nodes and users of MPI >> CLIC\_module
- MPI Requests are converted into CLIC messages with MPI characteristics (Buffered, Sincronous, Ready) & (comm,tid, source, dest)
- Request can progres automatically informing Upper MPI Layer.



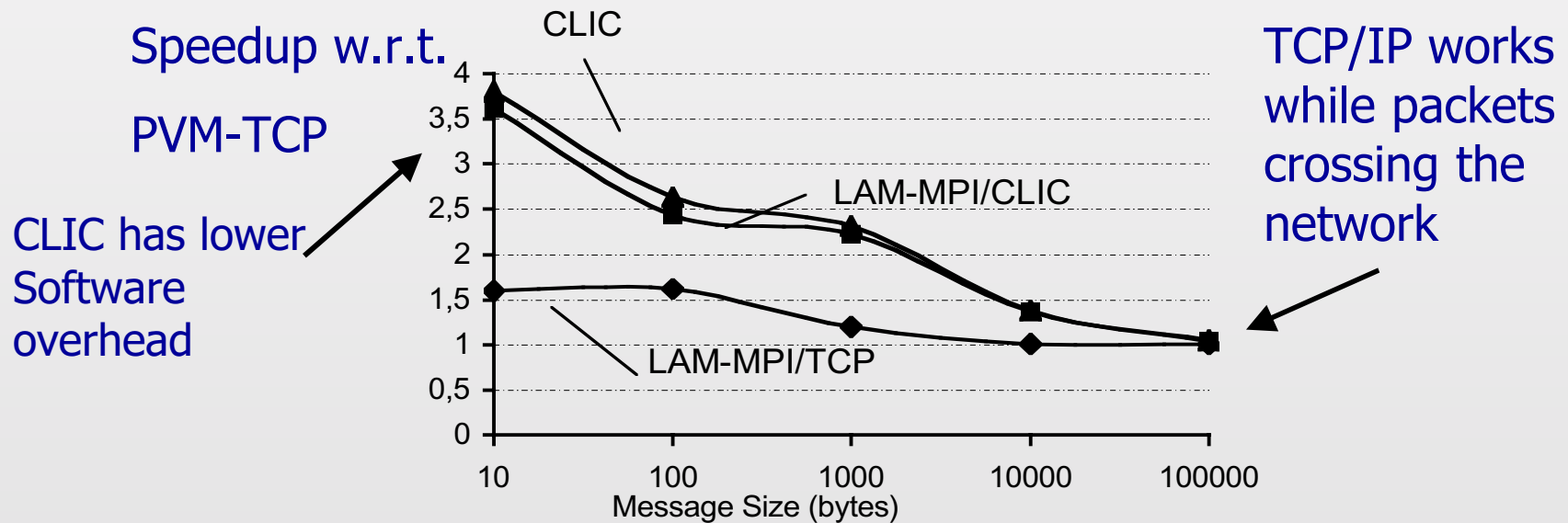
# CLIC Performances

Message Size	PVM	MPI-TCP	MPI-CLIC	CLIC
10	224	140	62	59
100	248	153	102	94
1000	454	379	204	196
10000	1285	1279	943	932
100000	9018	8990	8690	8615

Data transfers

(time in  $\mu\text{s}$ )

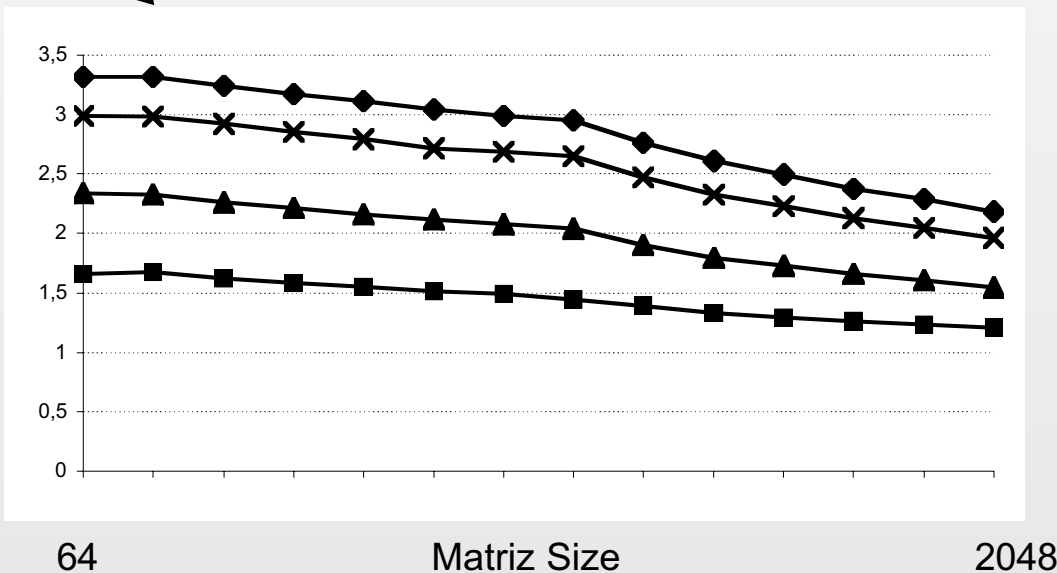
Network: Fast Ethernet  
& Switch



# CLIC Performances (II)

Matrix Product:

Latency is also important...



9 proc.  
8 proc.  
6 proc.  
4 proc.

More nodes, better performances if you use *broadcast send*

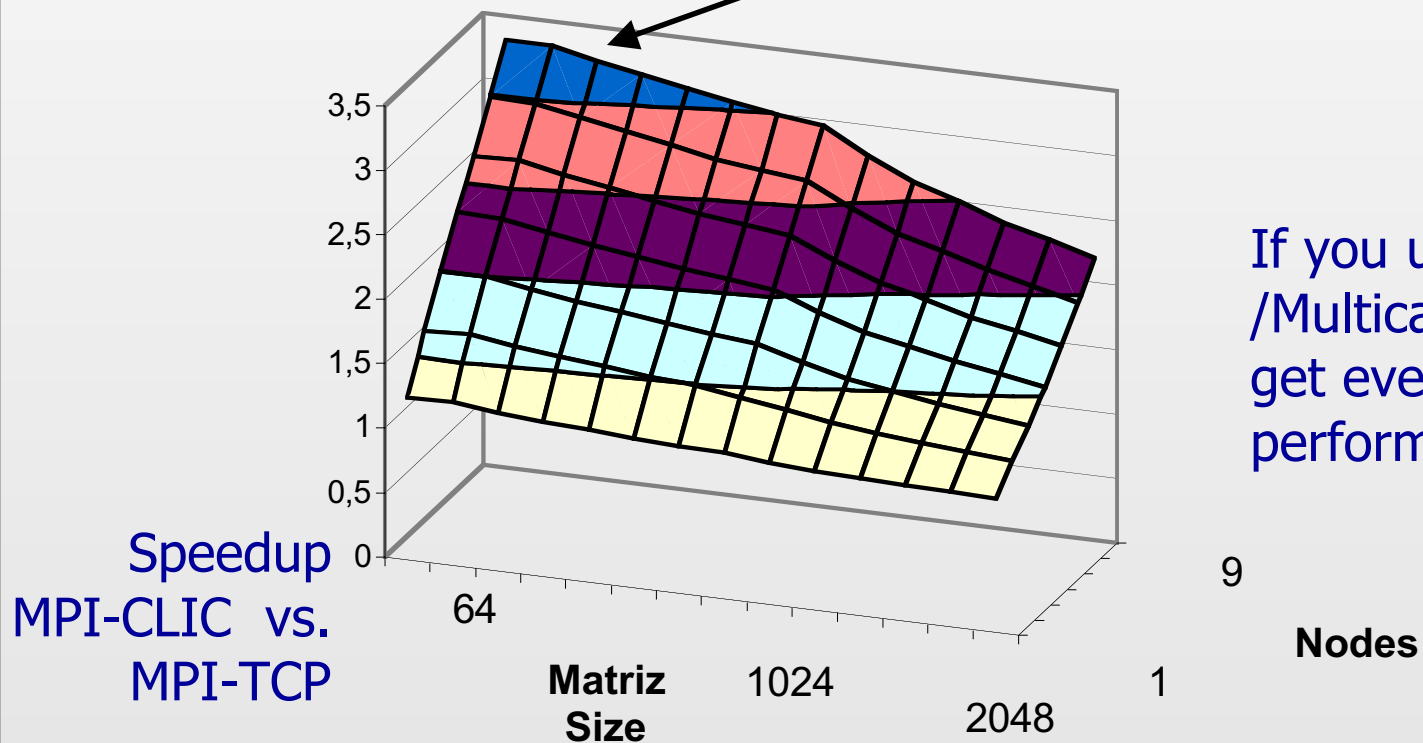
Speedup MPI-CLIC vs.  
MPI-TCP



# CLIC Performances (III)

Matrix Product:

Latency is also important...



If you use Broadcast /Multicast CLIC can get even better performances

## *Conclusions*

- A communication layer (CLIC) is proposed to use OS efficiently.
- Upper layer systems (PVM, MPI,...) can be efficiently used on top of CLIC.
- CLIC improves the performance of the communications so user-level applications can take advantage of network features. (Better latency & Bandwidth, *Broadcast, Channel Bonding*).
- This System can be upgraded to new networking technologies. (Future study)