# Message Progression in Parallel Computing - To Thread or not to Thread?

## Torsten Hoefler and Andrew Lumsdaine

Open Systems Lab
Indiana University
Bloomington, USA

IEEE Cluster 2008
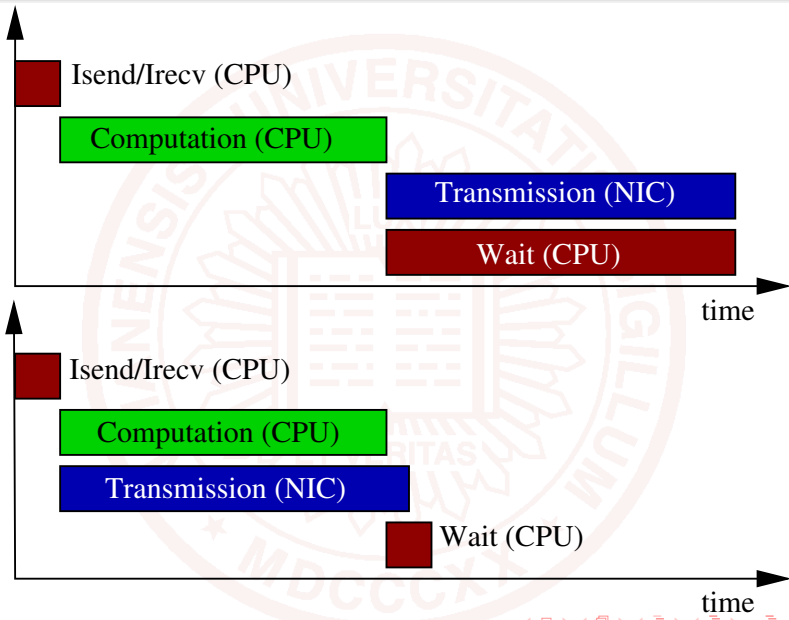
Tsukuba, Japan

September, 30th 2008

## Non-blocking Interfaces

- can help to hide latency
- mitigate effects of process skew
- reported application speedup up to 1.9
- requires much effort at the algorithm and implementation levels

## Examples

- MPI offers non-blocking point-to-point
- non-blocking collectives are discussed for MPI-3
- GASNet is fully non-blocking
- Asynchronous I/O

# Non-blocking does not mean asynchronous!

# Non-blocking Middleware Implementation

## Non-blocking Send/Receive

- eager protocol/copy for small messages
  - → uses a single message
- rendezvous protocol/synchronize for large messages
  - → uses multiple messages (two to three)
- OS bypass networking
  - → does not involve the kernel in send/receive
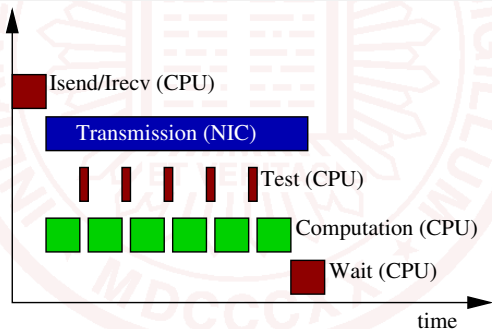  - → polling to check for messages

## Non-blocking Collectives

- similar issues as send/receive
- much more complex tasks and protocols
- multiple send/receive operations and dependencies in a single collective operation

# Progression Strategies I/II

## Manual Progression

- simplest to implement in a middleware
- user has to progress (e.g., calling MPI_Test)
- number of necessary progress calls depends on protocol
- best case: eager, worst case: pipelined protocols
- our proposed black-box scheme: $N = \lfloor \frac{size}{interval} \rfloor + 1$



time

# Progression Strategies II/II

## Hardware-based Progression

- need to run full protocol in NIC
- complicated to implement
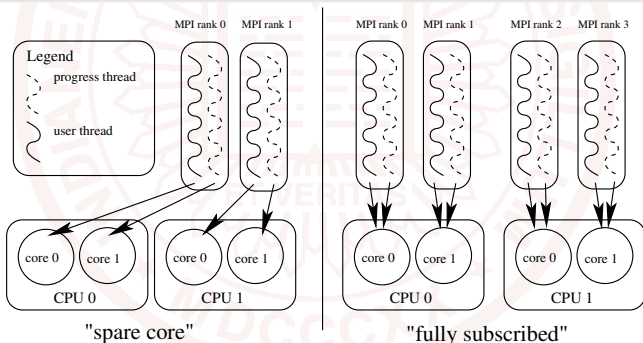- full benefits to the user
- mostly not supported

## Threaded Progression

- asynchronous threads
- often stated as "silver bullet" but not widely used (?)
- problem with manual progresssion: "fire at the right time"
- threads could solve this problem (woken up correctly)
- could enable fully asynchronous progression
- OS influence (scheduler) is significant
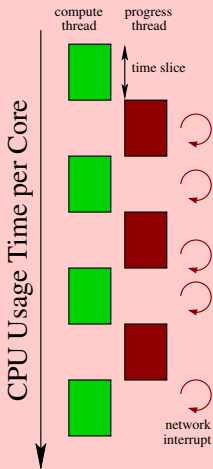
# Threading Configurations

## Spare Core vs. Fully Subscribed

- two extreme scenarios
- spare core: min P/2 cores are idle (one per process)
  $\rightarrow$ used in I/O or memory-bound applications
- fully subscribed: no cores idle
  $\rightarrow$ used in compute-bound applications



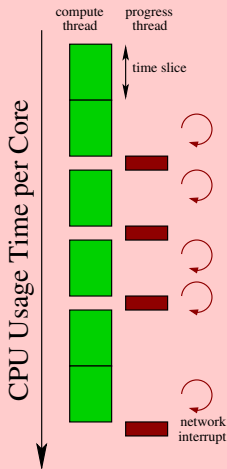"spare core"    "fully subscribed"
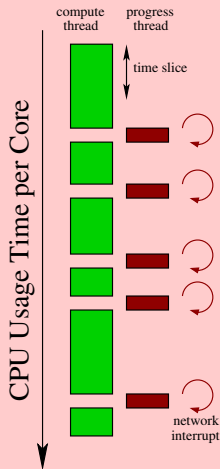
# Implementation Possibilities

## Polling vs. Interrupt vs. Real Time



a) polling

a) interrupt/normal

a) interrupt/real time

# Non-blocking Collectives and InfiniBand

## Issues with Non-blocking Collectives

- NBCs introduce dependencies
  - $\rightarrow$ e.g., sending a message in a tree
- dependencies might lead to synchronization

## Case-study InfiniBand

- supports polling and interrupt
- polling bad without non-spare core, else fastest
- interrupts are slow and cause overhead
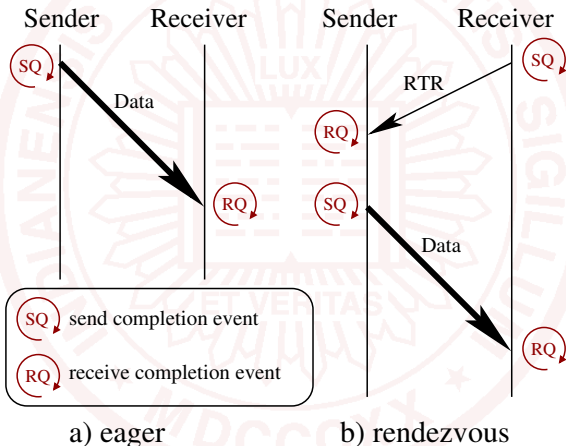- scheduler issues! (timeslice $4ms$, latency $3\mu s$)

## Real-Time Threads in Linux

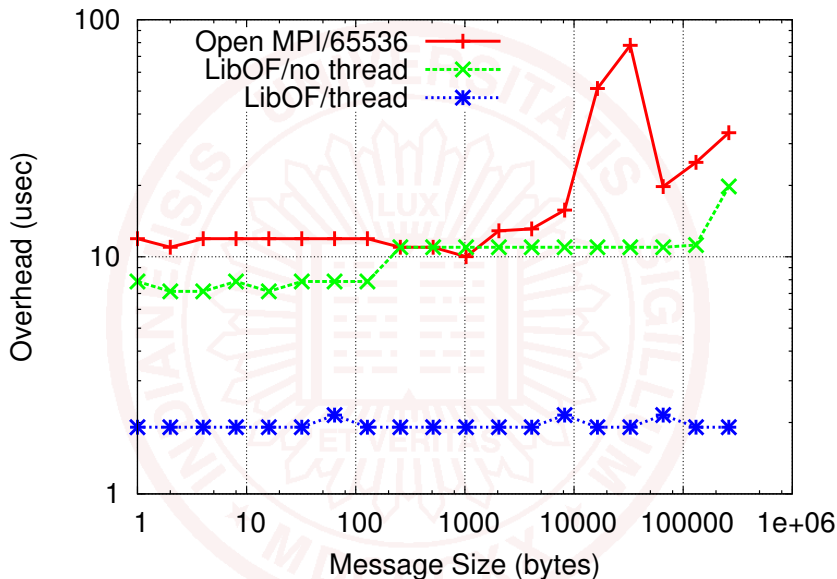- highest priority
- scheduled immediately
- no preemption
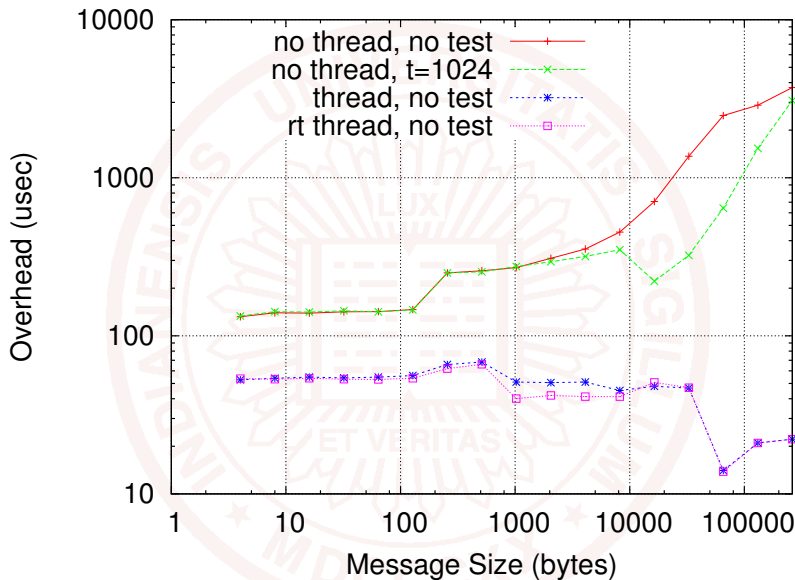
# Overhead of Threading

## Real time vs. Normal

- normal threads: interrupts coalesc, low (no) progression
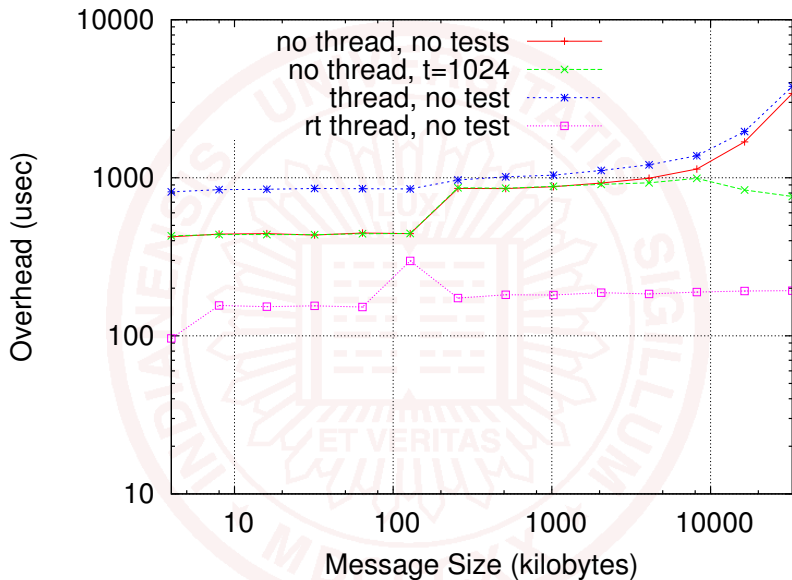- RT-threads: each interrupt pays full overhead



a) eager                    b) rendezvous

# Point-to-point overhead

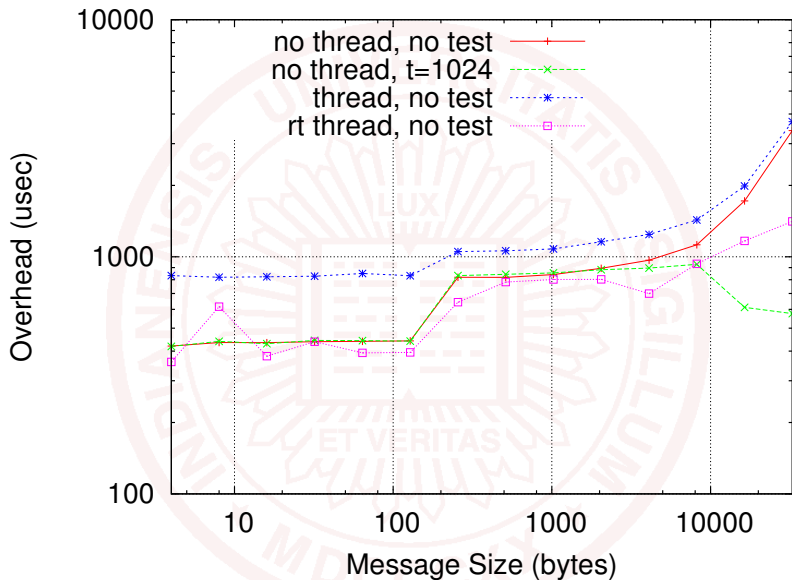## Wrong Metric?

- often used time-based benchmark
- hides interrupt overhead costs!
- overhead $\approx 3.4\mu s$ per interrupt
- many many interrupts in an NBC; 1016 in pipelined case

## Work-based benchmark

- compute fixed work quantum
- results account for interrupt overhead
- should be used for any threaded progression analysis!

## Summary

- we developed fully threaded LibNBC for IB
  $\rightarrow$ high overhead
- tested implementation with RT threads
  $\rightarrow$ lower overhead (better than theory?)
- implemented new work-based benchmarking metric
  $\rightarrow$ realistic (high) overhead

## Conclusions and Future Work

- threaded implementation makes sense with spare cores
- very tricky without spare cores $\rightarrow$ manual again?
- investigate other options
  $\rightarrow$ signalled progression (not safe/realistic!)
  $\rightarrow$ OS involvement (opposite to OS bypass)
  $\rightarrow$ hardware progression

# Summary and Future Work

## Summary

- we developed fully threaded LibNBC for IB
  $\rightarrow$ high overhead
- tested implementation with RT threads
  $\rightarrow$ lower overhead (better than theory?)
- implemented new work-based benchmarking metric
  $\rightarrow$ realistic (high) overhead

## Conclusions and Future Work

- threaded implementation makes sense with spare cores
- very tricky without spare cores $\rightarrow$ manual again?
- investigate other options
  $\rightarrow$ signalled progression (not safe/realistic!)
  $\rightarrow$ OS involvement (opposite to OS bypass)
  $\rightarrow$ hardware progression