

# Energy Characterization and Optimization in Heartbeat Failure Detection Systems

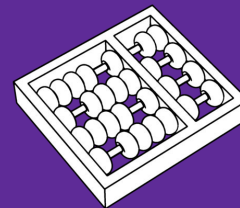
**Pedro Henrique Rosso**

p233687@dac.unicamp.br

*University of Campinas (Unicamp)*

*MO632/MC972 – Energy Efficient Computing*

**Final Project Presentation** • November 29<sup>th</sup> 2021



UNICAMP

# Overview

- Introduction
  - Problem
  - Objectives
- Heartbeat Systems
- Methodology
  - Proposals
  - Experiments
- Results and Discussions

# Introduction



- High-Performance Computing
- Fault Tolerance
  - Failure detection mechanisms
    - Heartbeat (ULFM, **OCFTL**)
    - Daemon processes (Reinit++)
- Fault Tolerance in IoT

- During my Master
  - OCFTL
  - MPI-Based FT library for OmpCluster
- No concern about energy efficiency
- Characterize and think about it focusing energy efficiency

- Energy Characterization
  - High resource usage - Parallel Applications
  - Low resource usage
- Impact of heartbeat parameters
- Different communication backends
- Optimizations in current algorithm
- Relationship between Energy Efficiency and Heartbeat Efficiency

# Heartbeat Systems



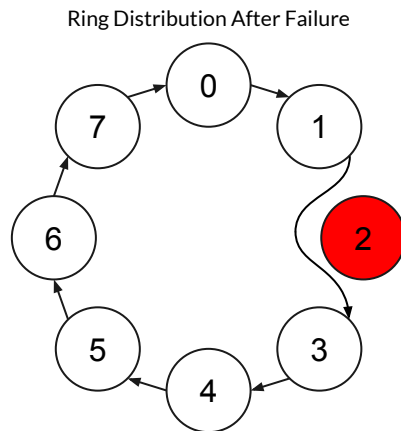
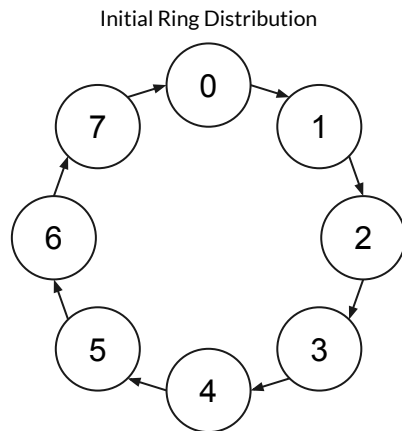
- Inspired by ULFM's heartbeat
- Broadcast Propagation

## Properties

- Period
- Timeout

## Definitions

- Emitter
- Observer





# Heartbeat - Implementation (Main Loop)

---

**Algorithm 1:** Standard OCFTL heartbeat algorithm

---

```
1 while (!hb_done) do
2   if timeout expires then
3     broadcast(emitter);
4     find new emitter ;
5     rearrange the ring;
6   end
7   if period achieved then
8     alive_beat(observer);
9   end
10  if (alive msg received) then
11    resets timeout;
12  end
13  if (new_observer msg received) then
14    observer  $\leftarrow$  new_observer;
15  end
16  if (broadcast received) then
17    do related procedure;
18    if (first time this broadcast) then
19      broadcast(bc_message);
20    end
21  end
22  sleep for (timestep)
23 end
```

---

# Methodology



- Evaluate the standard algorithm
- Proposes modifications to improve the main loop
- Test with other backend (Not MPI)
- Characterize all options for different HB parameters
- Evaluate relationship between HB efficiency and energy

# Proposals



```

1 Main-Thread:
2 while (!hb_done) do
3   if (timeout expires) then
4     find new emitter ;
5     broadcasts the failure;
6     rearrange the ring;
7   end
8   if (any message received) then
9     if (type == alive) then
10      resets timeout;
11    end
12    if (type == new_observer) then
13      observer ← new_observer;
14    end
15    if (type == broadcast) then
16      do related procedure;
17      if (first time this broadcast) then
18        replicates the broadcast;
19      end
20    end
21  end
22  sleep_for(timestep)
23 end
24
25 Alive-Thread:
26 while (!hb_done) do
27   alive_beat(observer);
28   sleep_for(period)
29 end


```

## NEW - OCFTL

- Checks only one time for a message
- Standard calls 4 times MPI calls
- This approaches calls 1 time
- Extra thread to send beats

## NNG - OCFTL

```

1 Main-Thread:
2 while (!hb_done) do
3   if (timeout expires) then
4     find new emitter ;
5     broadcasts the failure;
6     rearrange the ring;
7   end
8   for ( i = 0; i < size; i = i + 1 ) { 
9     if (any message received in socket[i]) then
10      if (type == alive) then
11        resets timeout;
12      end
13      if (type == new_observer) then
14        observer ← new_observer;
15      end
16      if (type == broadcast) then
17        do related procedure;
18        if (first time this broadcast) then
19          replicates the broadcast;
20        end
21      end
22    end
23  }
24  sleep_for(timestep)
25 end
26
27 Alive-Thread:
28 while (!hb_done) do
29   alive_beat(observer);
30   sleep_for(period)
31 end

```

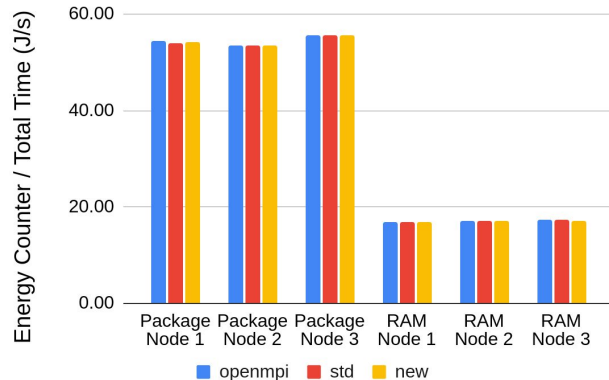
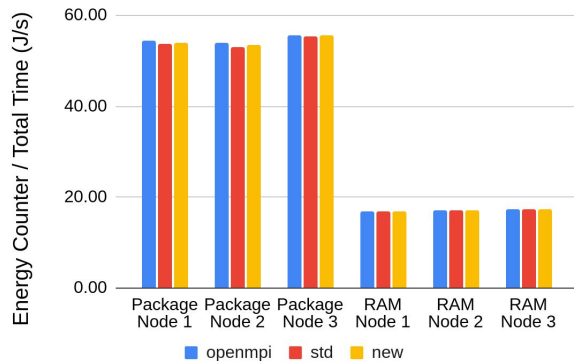
- Uses TPC based communication
- Using NNG backend
- TCP Socket between each pair of processes
- Extra thread to send beats

- Execution on Sorgan (OmpCluster's small cluster)
- 4 Applications (5 samples)
  - LULESH (<https://github.com/LLNL/LULESH>) - Hydrodynamics
  - HPCCG (<https://github.com/Mantevo/HPCCG>) - Conjugate Gradients
  - MINIMD (<https://github.com/Mantevo/miniMD>) - Molecular Dynamics
  - LIBONLY (Custom made) - Low resource usage program
- Based on OpenMPI V4.0
- Range from 27 to 36 processes divided in 3 nodes
- Evaluation
  - Package (Total Energy / Time Spent)
  - RAM (Total Energy / Time Spent)

# Results and Discussions





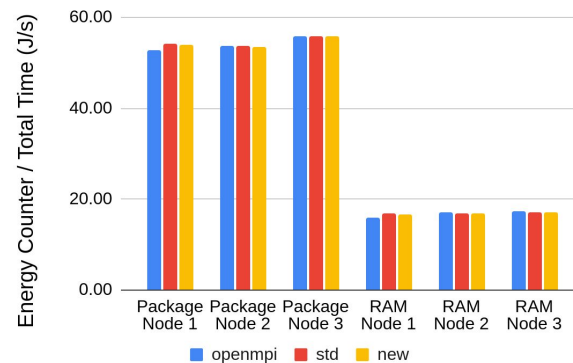


Timeout = 10s  
Period = 1s  
Timestep = 20ms

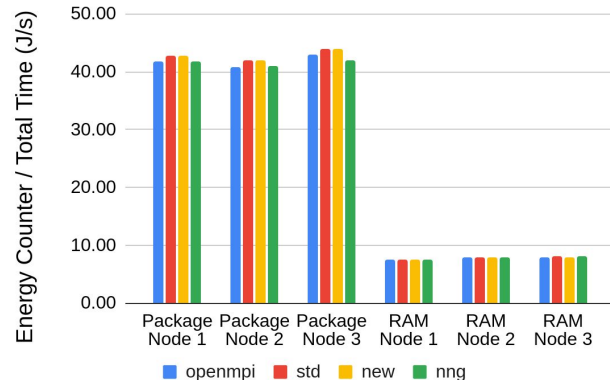
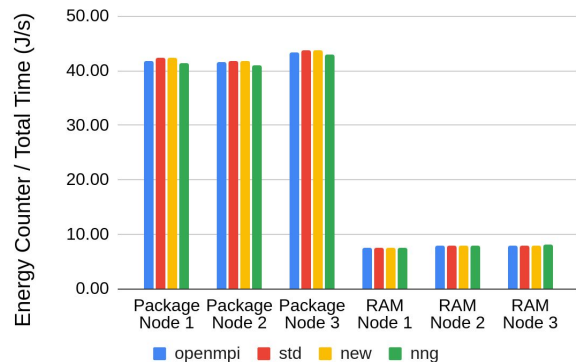
Timeout = 2s  
Period = 100ms  
Timestep = 10ms

## Keynotes:

- NNG showed false-positives in every test
- Results are very close



Timeout = 1s  
Period = 10ms  
Timestep = 1ms

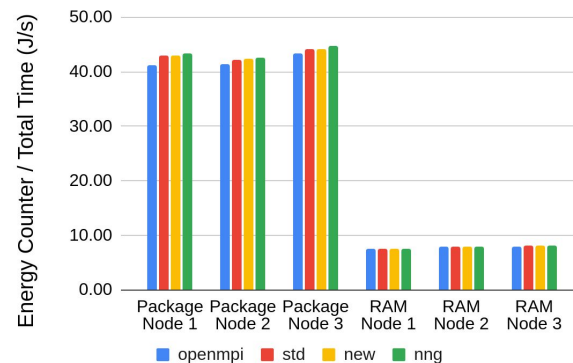


Timeout = 10s  
Period = 1s  
Timestep = 20ms

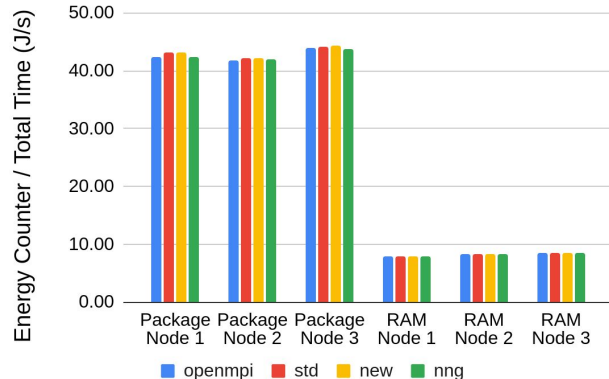
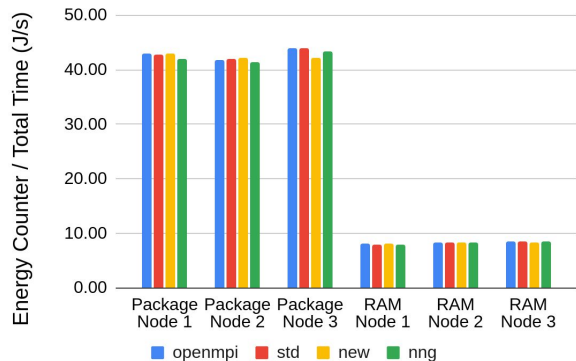
Timeout = 2s  
Period = 100ms  
Timestep = 10ms

## Keynotes:

- NNG performs a bit better for relaxed HB parameters. Lose in low value parameters
- STD and NEW keep very close results



Timeout = 1s  
Period = 10ms  
Timestep = 1ms

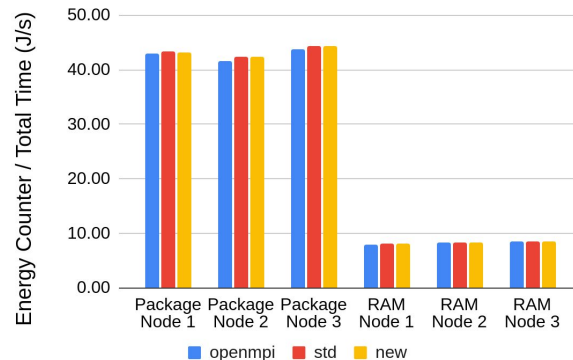


Timeout = 10s  
Period = 1s  
Timestep = 20ms

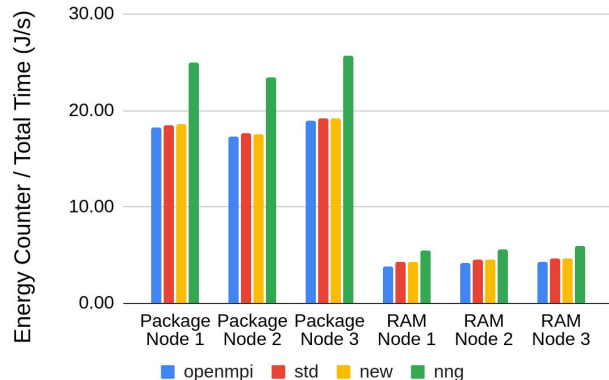
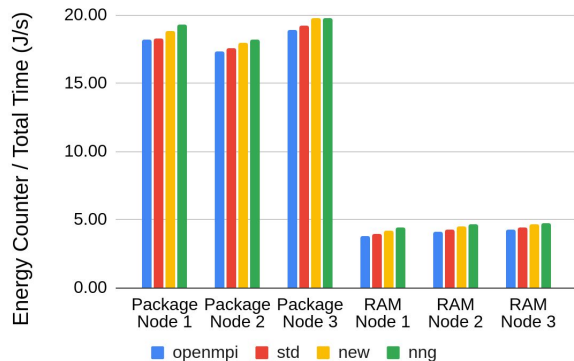
Timeout = 2s  
Period = 100ms  
Timestep = 10ms

## Keynotes:

- NNG can not perform low HB parameters
- NEW performs marginally better than STD

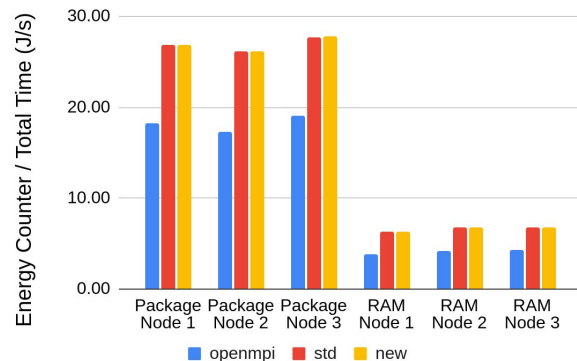


Timeout = 1s  
Period = 10ms  
Timestep = 1ms



## Keynotes:

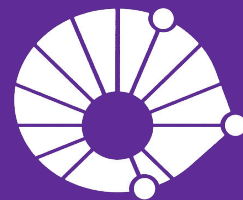
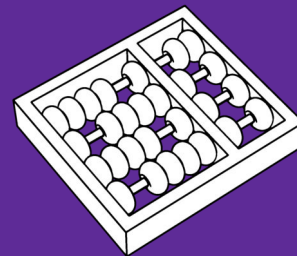
- Pure FT lib overhead
- NNG worst in all cases. Can not perform low HB parameters
- STD and NEW very close results



# Conclusions

- Characterized the HB system
  - Resource Intensive
  - Pure lib overhead
  - Energy vs HB efficiency
  - Different Backends
    - NNG is not a good option
  - What about other FD systems?
  - Algorithm
-

# Thank you!



UNICAMP