

AutoNUMA

Red Hat, Inc.

Andrea Arcangeli
aarcange at redhat.com

30 Aug 2012

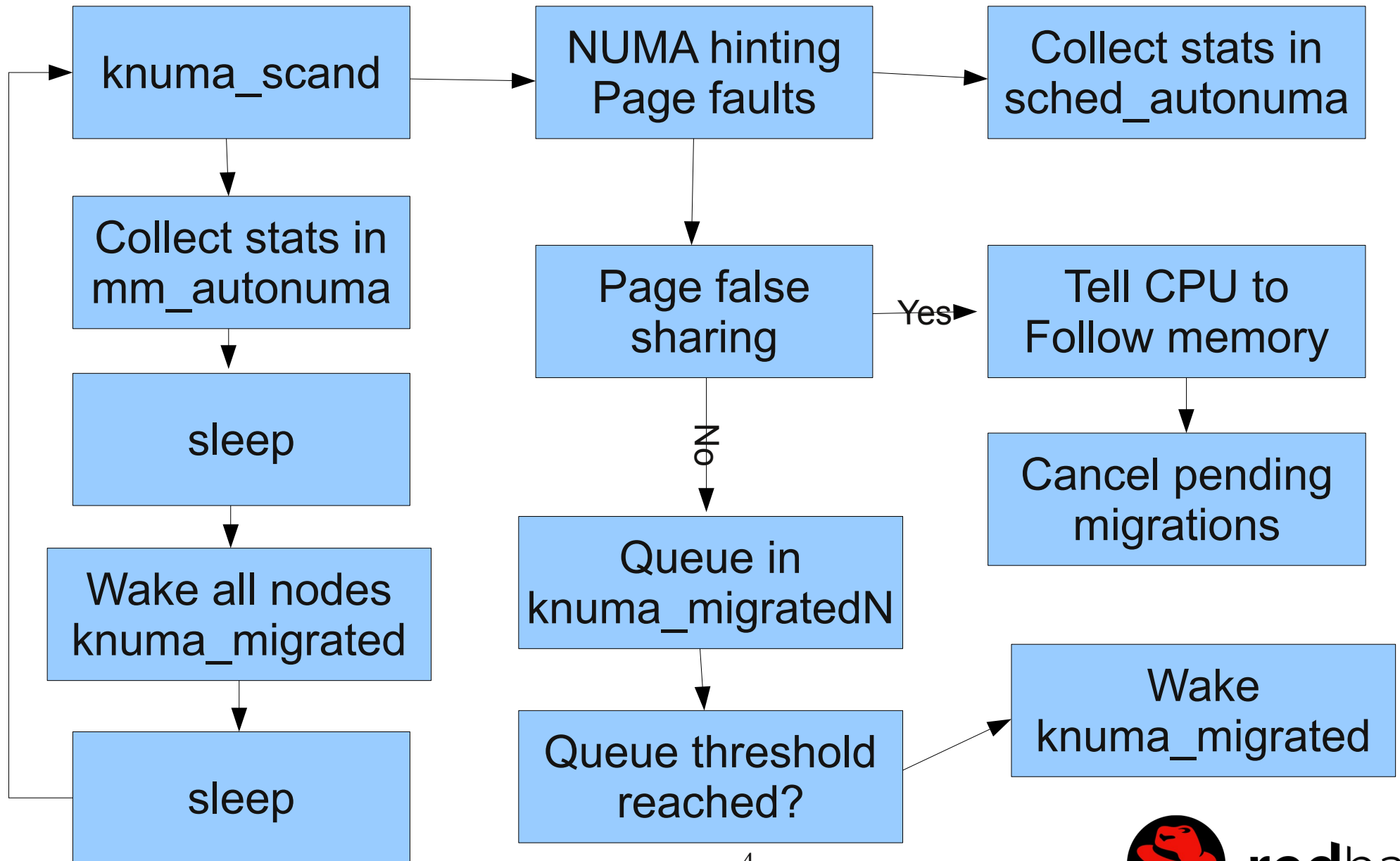
AutoNUMA components

- knuma_scand
 - If stopped, everything stops
 - Triggers the chain reaction when started
- NUMA hinting page faults
- knuma_migratedN (per node)
- scheduler (CPU follow memory & active idle balancing)
- Memory follow CPU (NUMA hinting page faults)
- False sharing detection

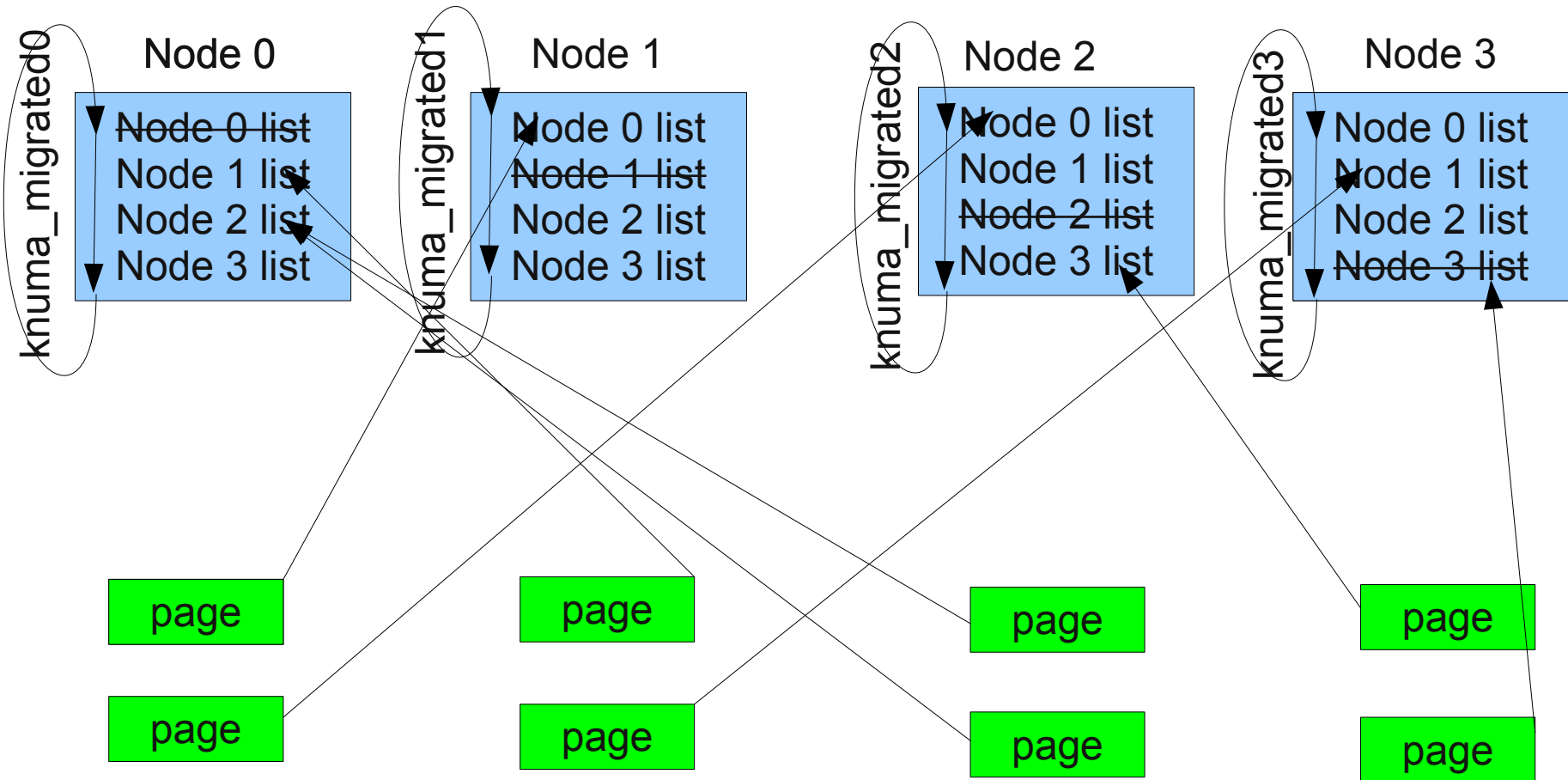
AutoNUMA data

- sched_autonuma
 - task_struct
- mm_autonuma
 - mm_struct
 - Working set or ~RSS

AutoNUMA logic



AutoNUMA knuma_migratedN

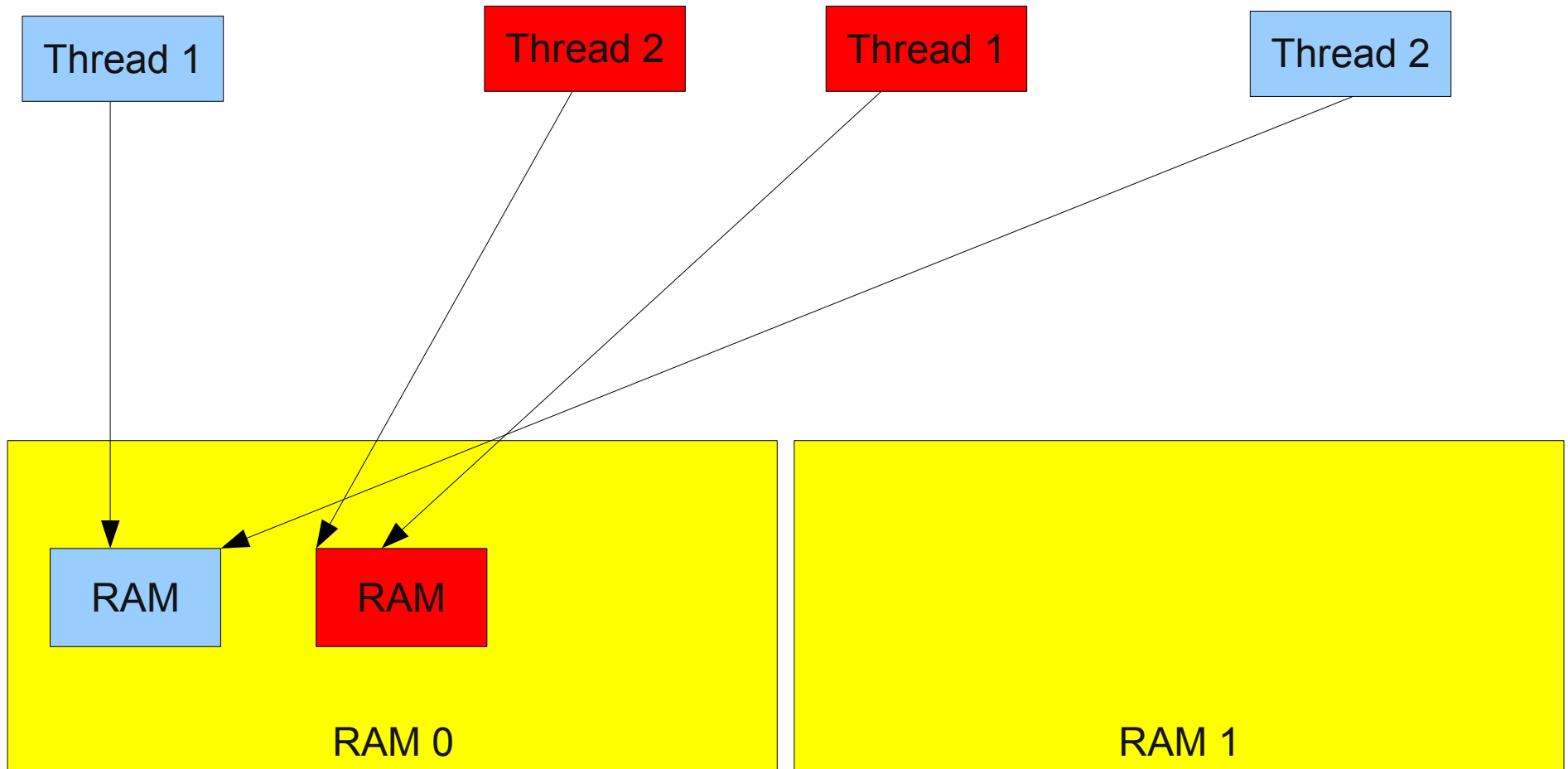


AutoNUMA-benchmark tests

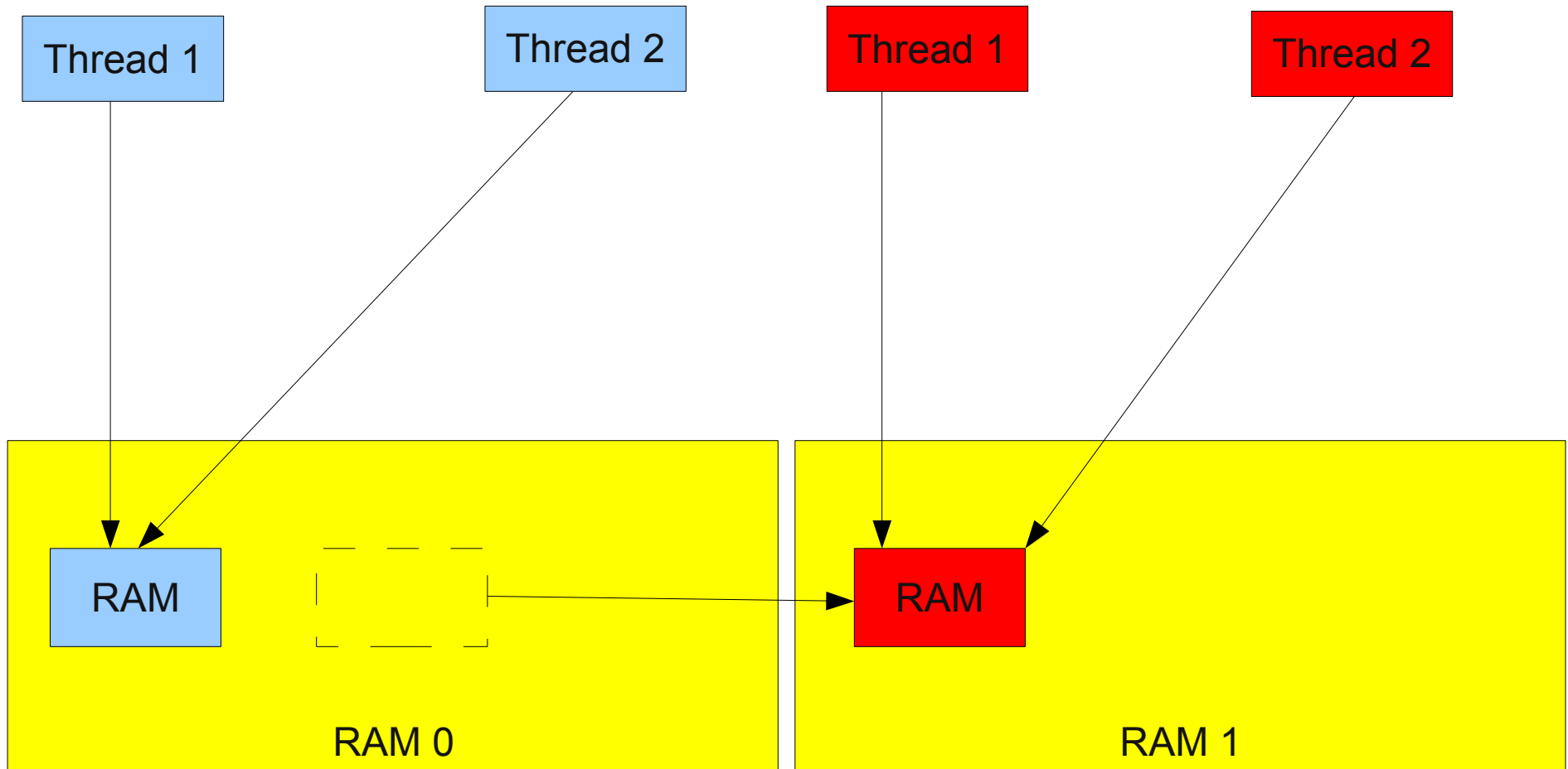
- P =# of processes, T =# of threads per process, M =memory per process
- numa01
 - $M=3\text{GiB}$ $P=2$, $T=\text{nr_cpus}/2$, all threads share all process memory
- numa01_THREAD_ALLOC
 - $M=3\text{GiB}$, $P=2$, $T=\text{nr_cpus}/2$, all threads use per-thread local memory
- numa02
 - $M=1\text{GiB}$, $P=1$, $T=\text{nr_cpus}$, all threads use per-thread local memory
- numa02_SMT
 - $M=1\text{GiB}$, $P=1$, $T=\text{nr_cpus}/2$, all threads use per-thread local memory
 - The kernel to get this right must not use more than one HT thread per core, and in turn it must decide to split the load over two NUMA nodes even if the load would fit in a single NUMA node
 - Testing with $T=\text{nr_cpus}/4$ and smaller T values, would also be interesting, but if the kernel behaves well with $T=\text{nr_cpus}/2$ there's a good chance it'll behave sanely with more than half of the CPUs idle too
- More will be added...



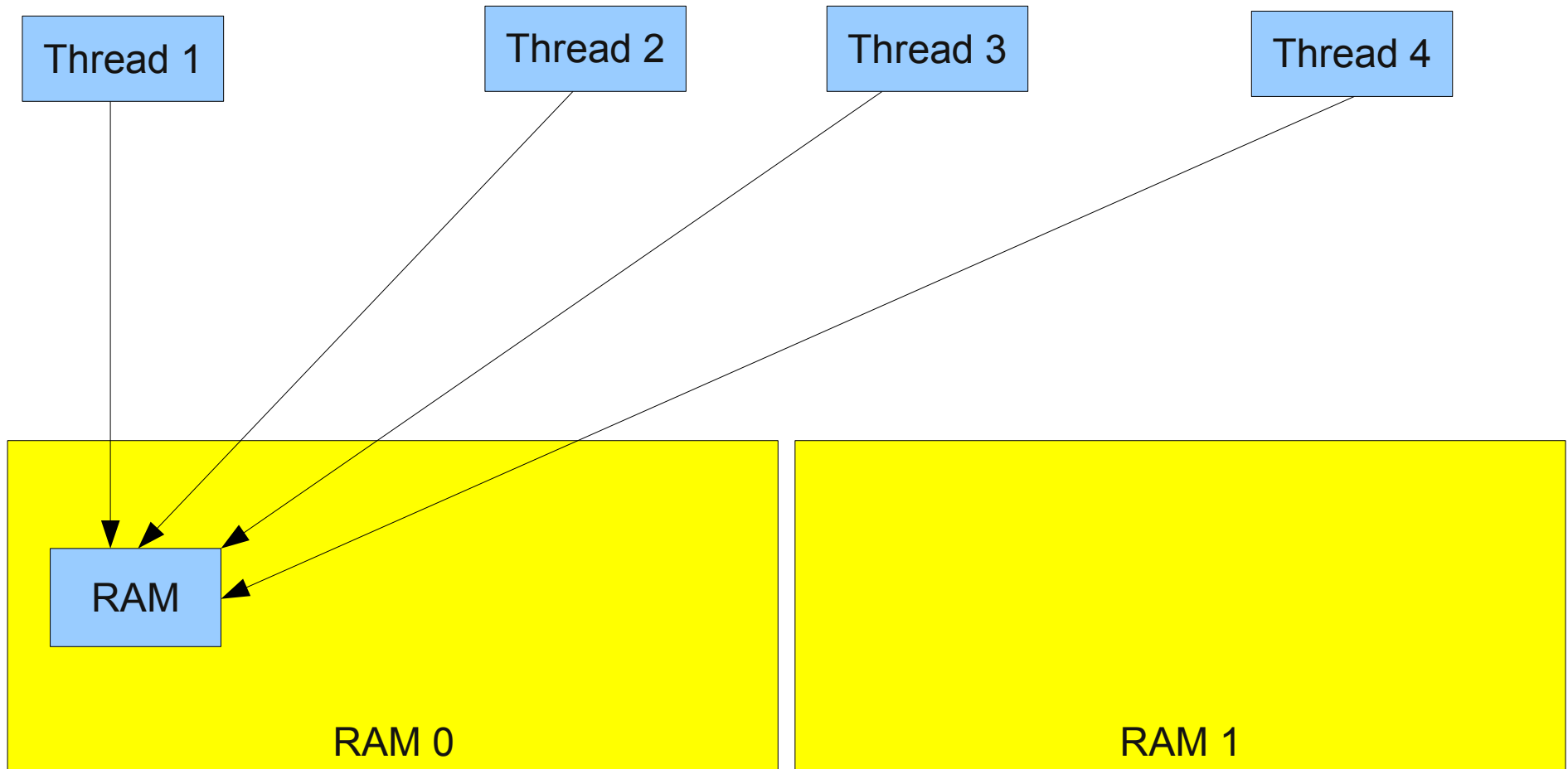
numa01 (startup)



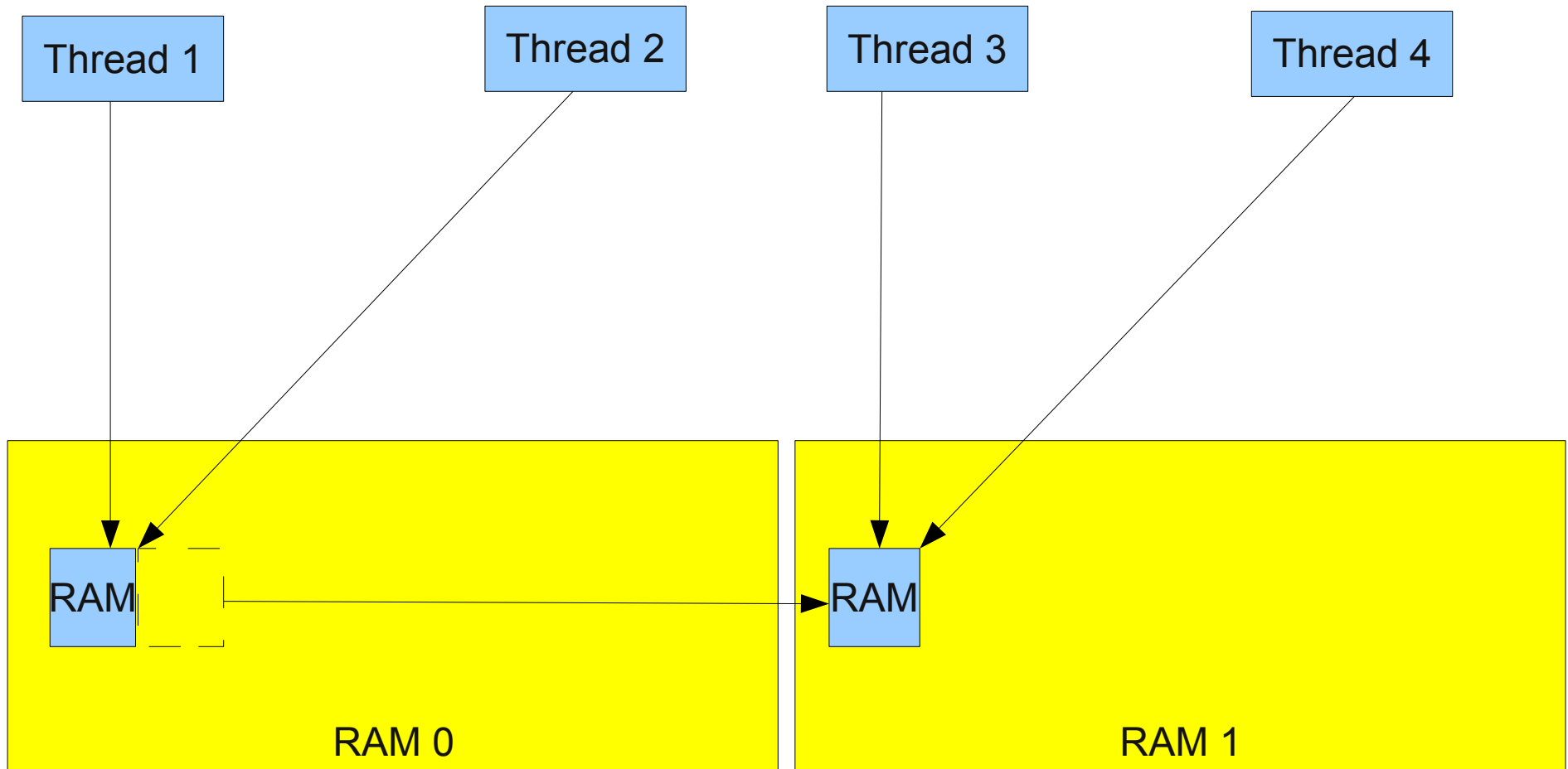
numa01 (converged)



numa02 (startup)



numa02 (converged)



autonuma-benchmark

```
$ git clone git://gitorious.org/autonuma-benchmark/autonuma-benchmark.git
```

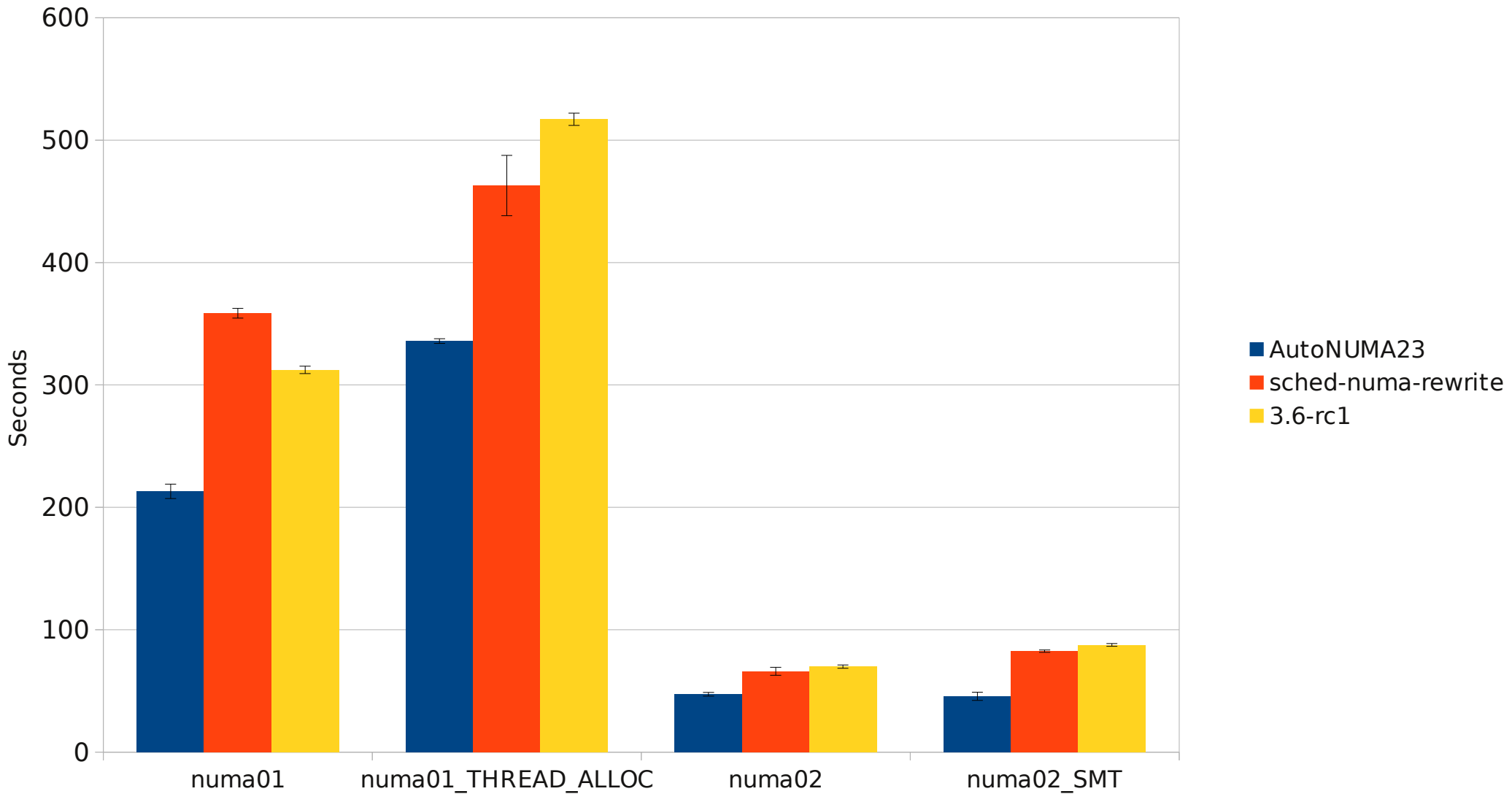
```
$ cd autonuma-benchmark
```

```
$ sudo ./start_bench.sh -s -t
```



4 runs, average and stdev

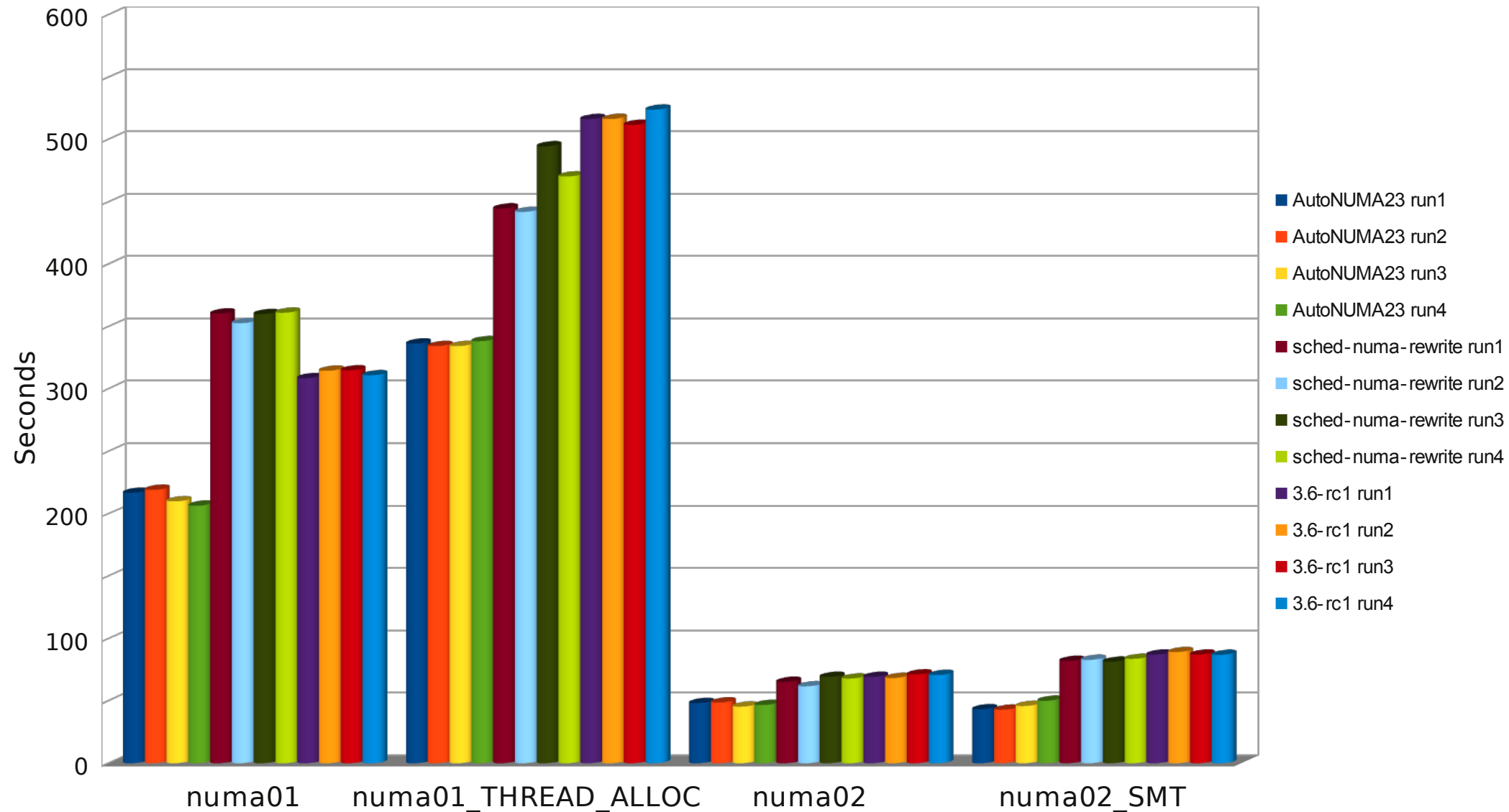
Lower is better



2 nodes

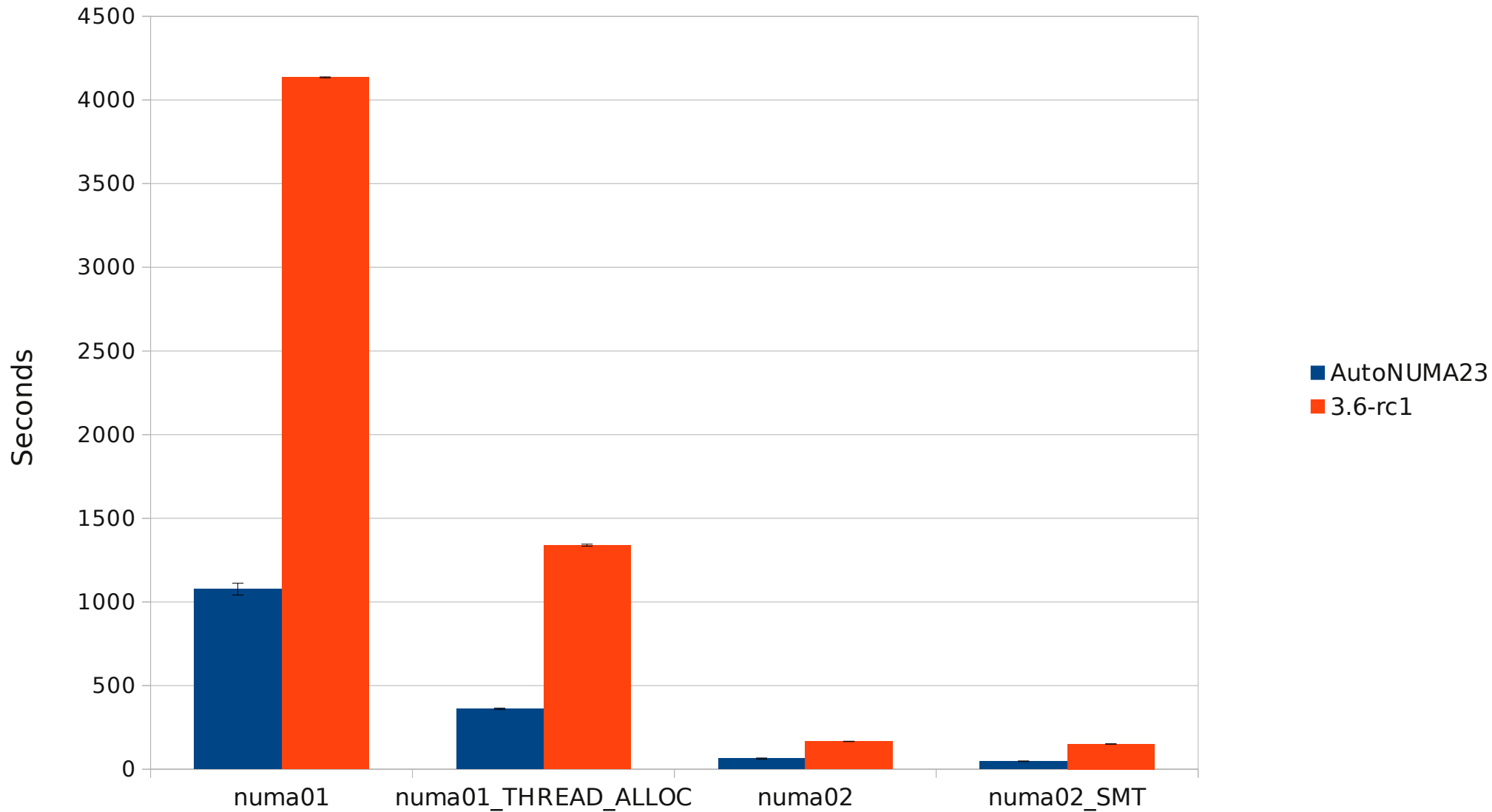
4 runs of each test, all results

Lower is better



4 runs, average and stdev

Lower is better

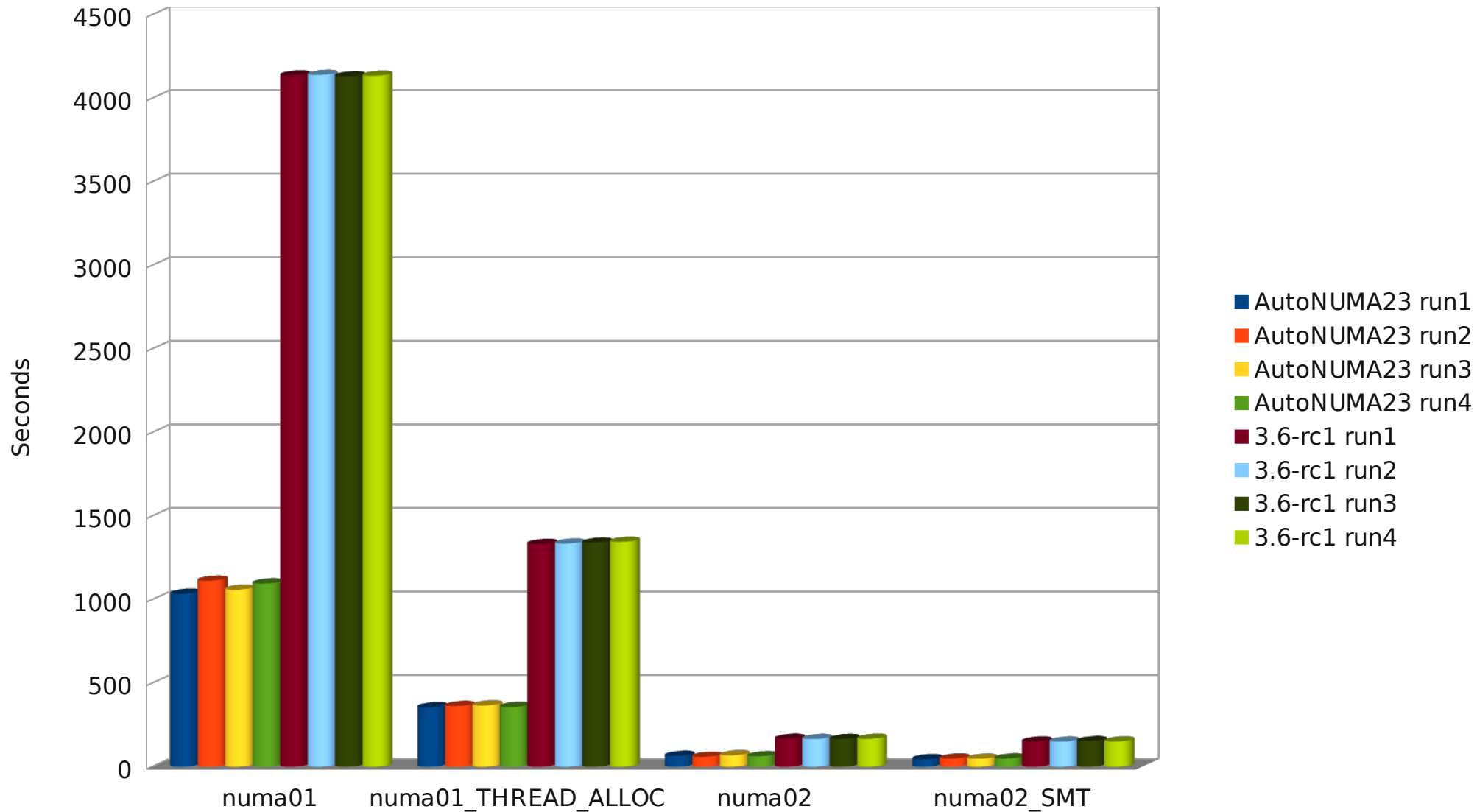


8 nodes



4 runs of each test, all results

Lower is better



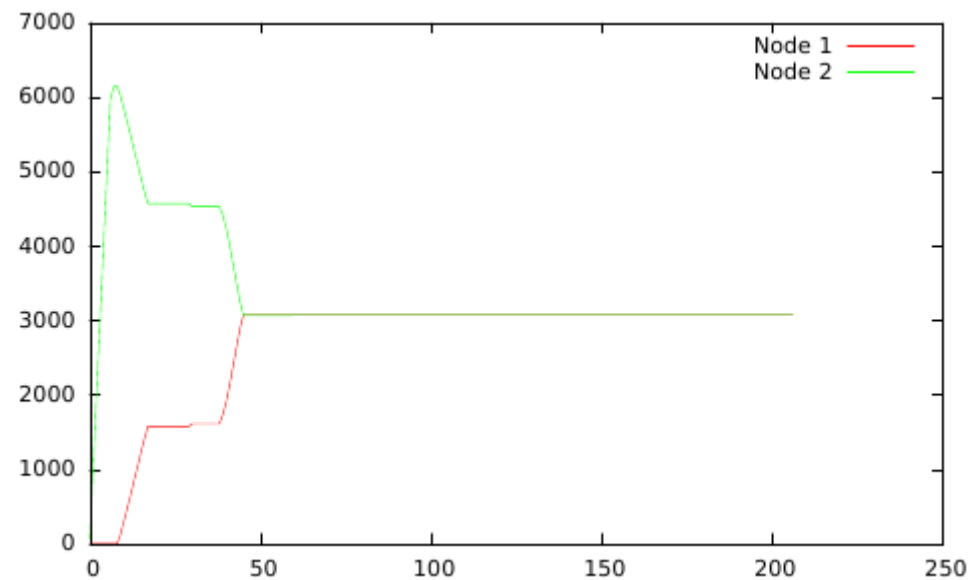
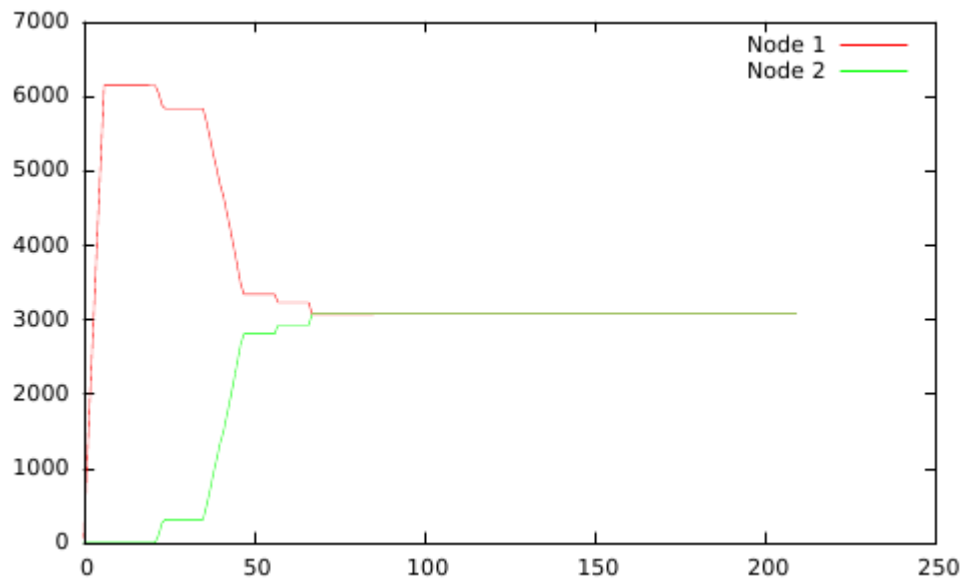
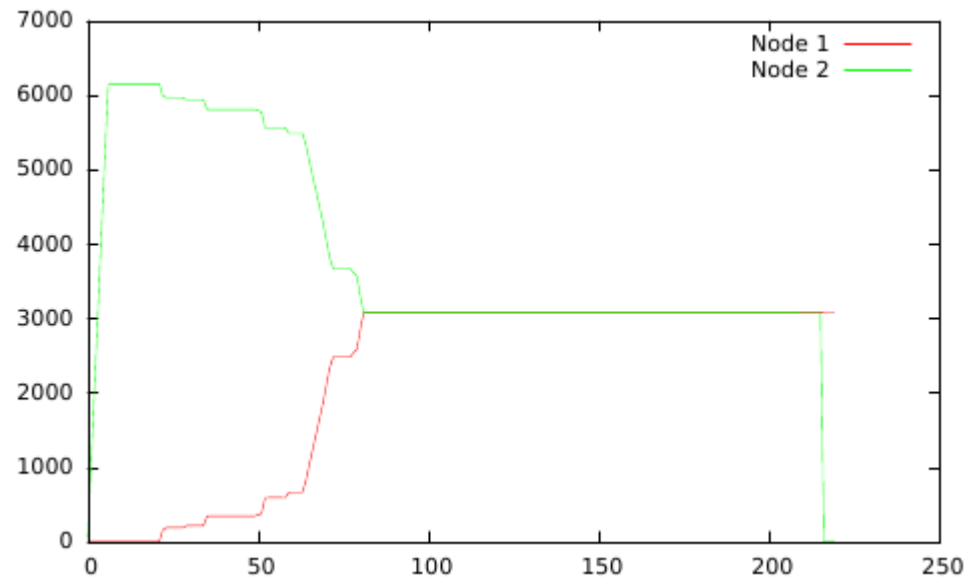
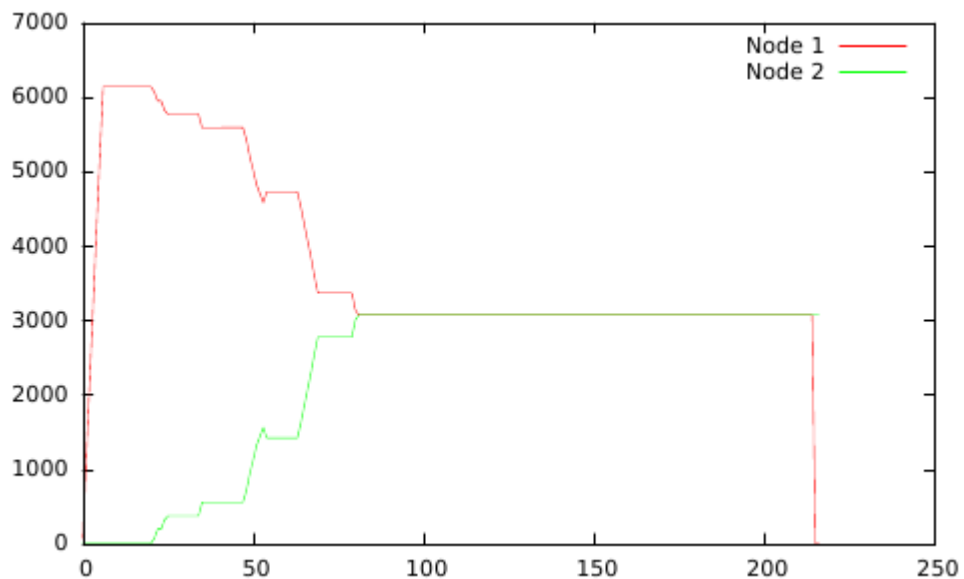
8 nodes

Convergence charts

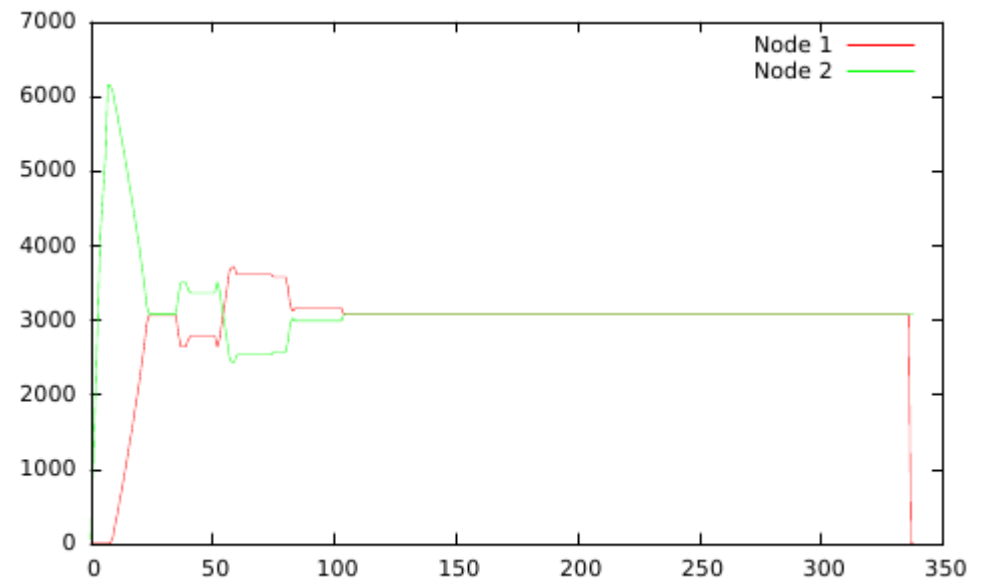
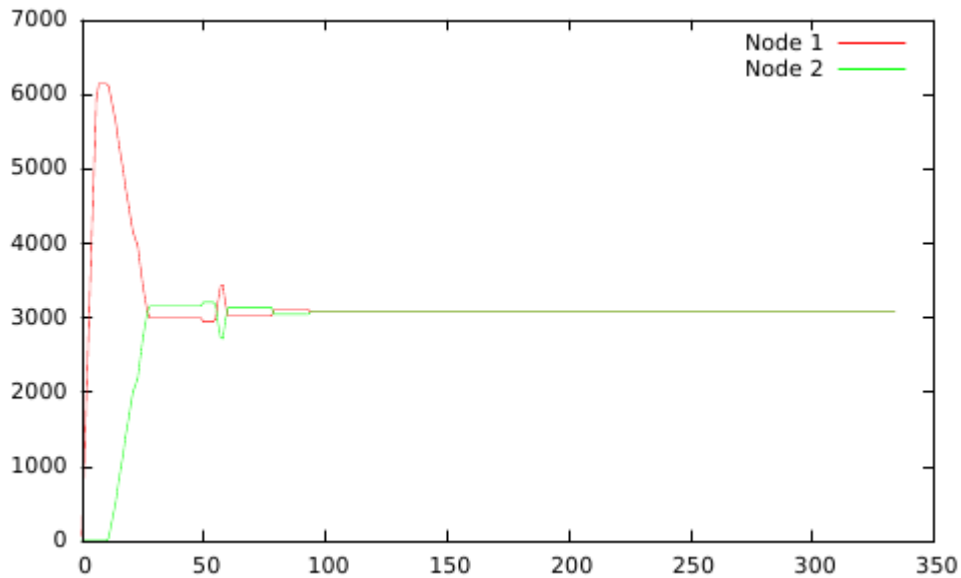
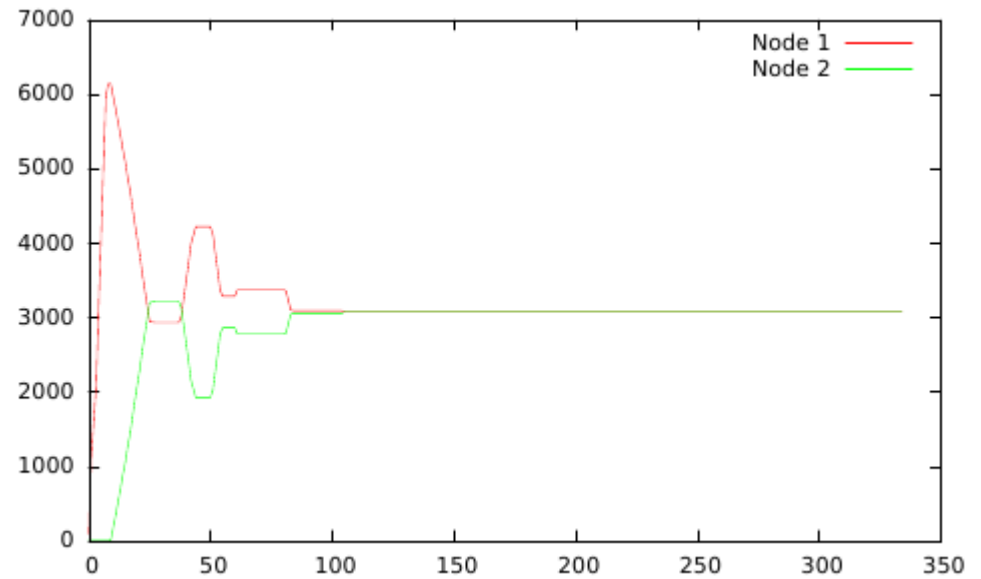
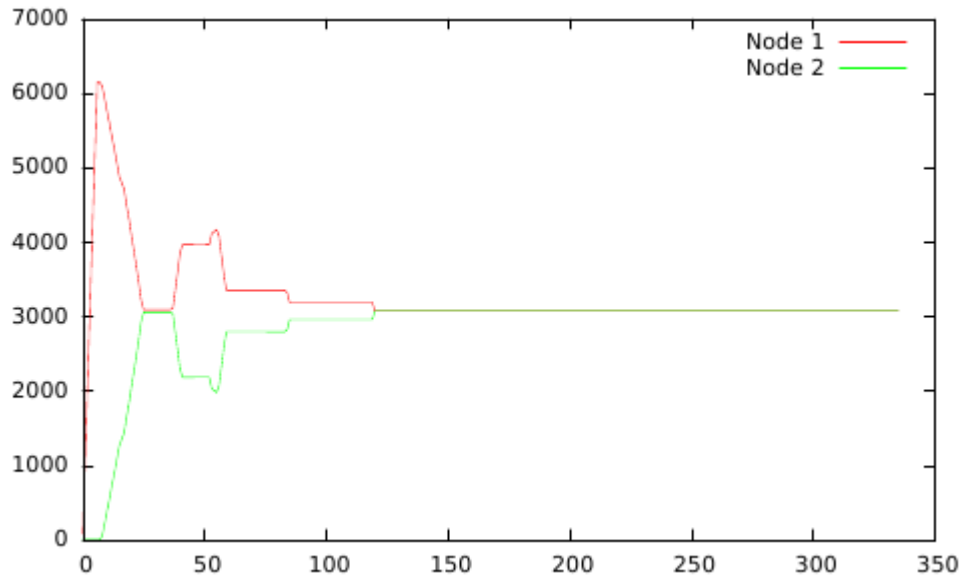
- The AutoNUMA-benchmark produces a chart for each test run (in pdf format).
 - In each chart there is one two dimensional line per NUMA node (node1, node2, etc.).
 - X=time (seconds)
 - Y=memory (MiB)
 - Each line represents how much of the test process's memory is in that NUMA node over time, for the duration of the test. Because all memory starts in one node, it illustrates how the memory migrates over time.
 - A workload converges when the memory levels are equal in all NUMA nodes
 - Note: with `nr_nodes > 2`, numa01 may not fully converge because half of the cpus will thrash on the memory of half of the nodes, but it should get close enough
 - In the future we plan to add a new numa01 "PER_NODE" test with `P=nr_nodes` and `T=nr_cpus/nr_nodes`



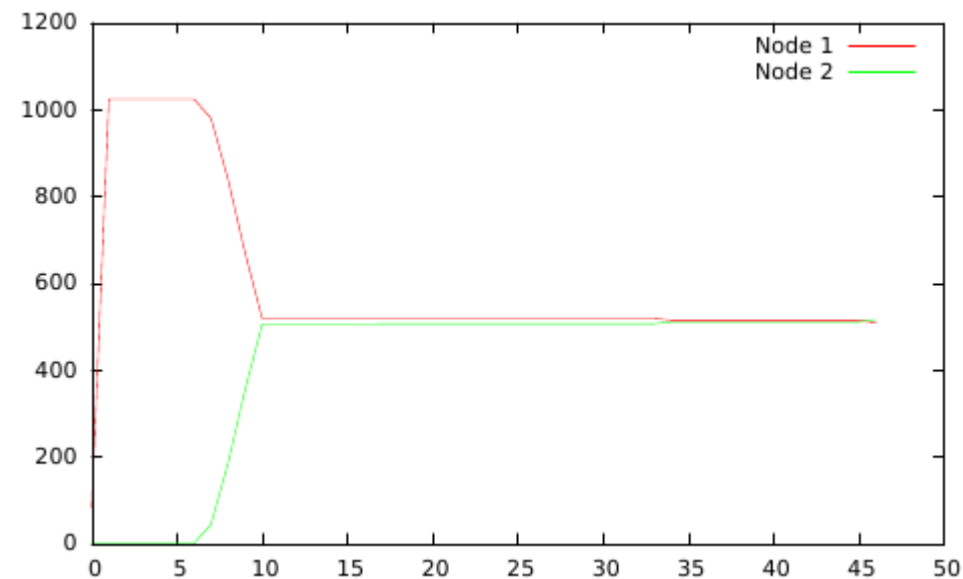
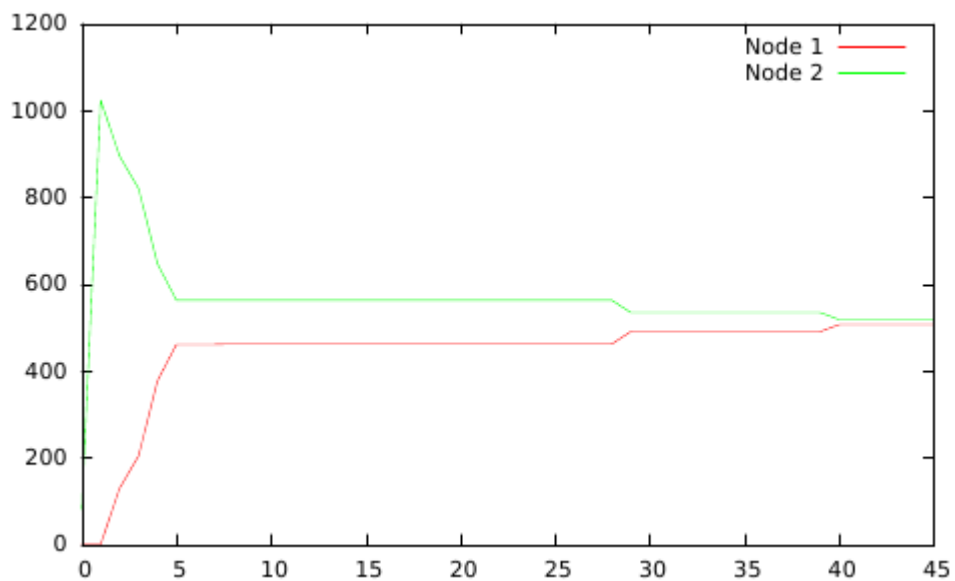
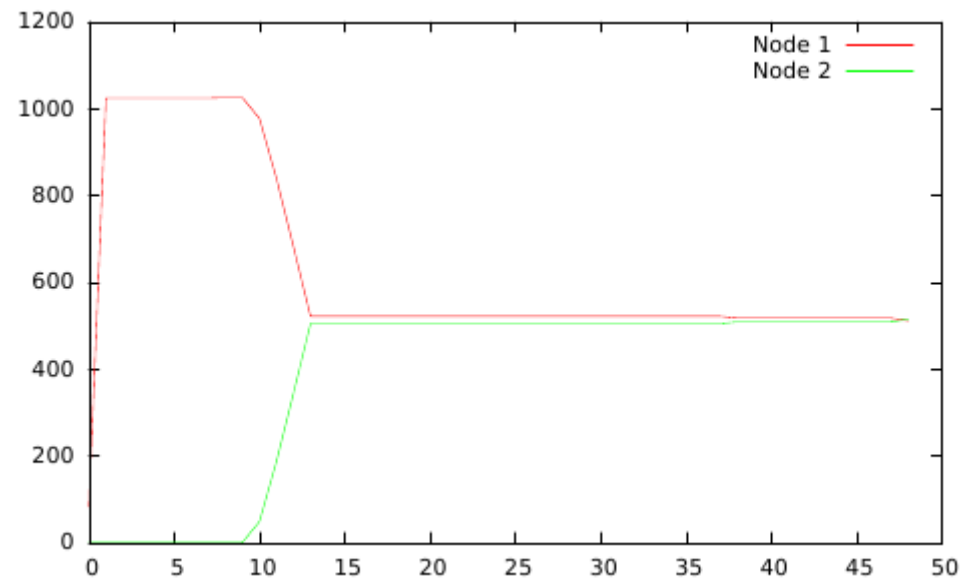
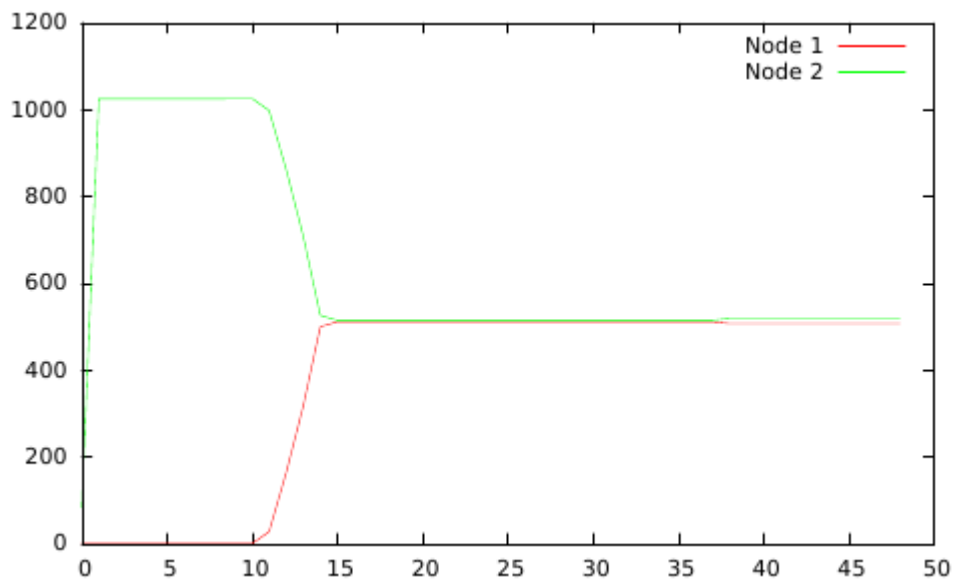
AutoNUMA23 numa01



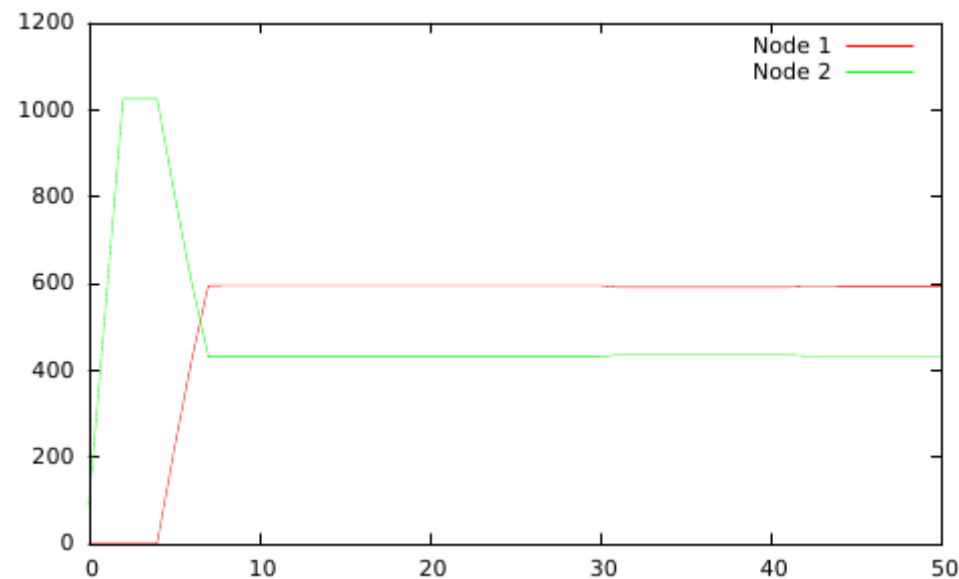
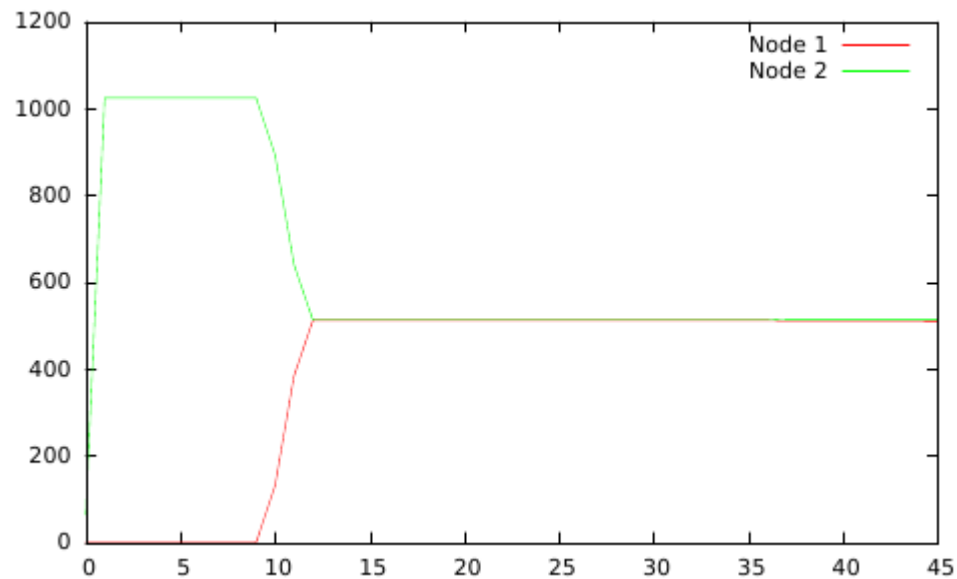
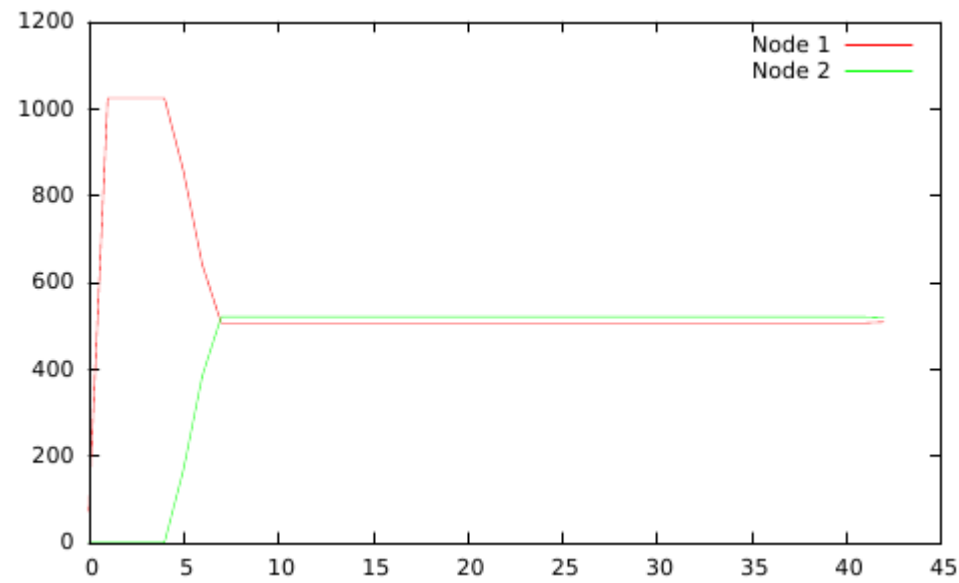
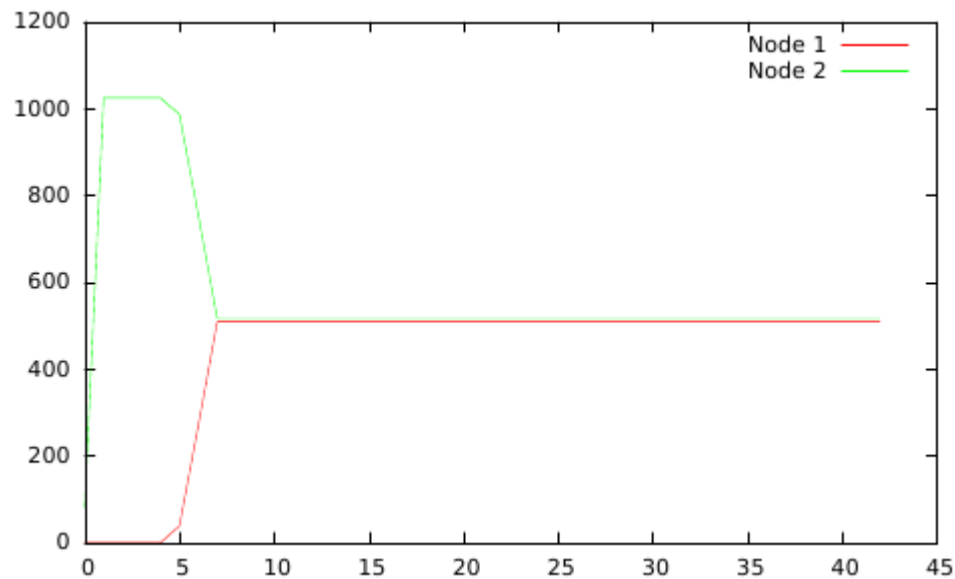
AutoNUMA23 numa01_THREAD..



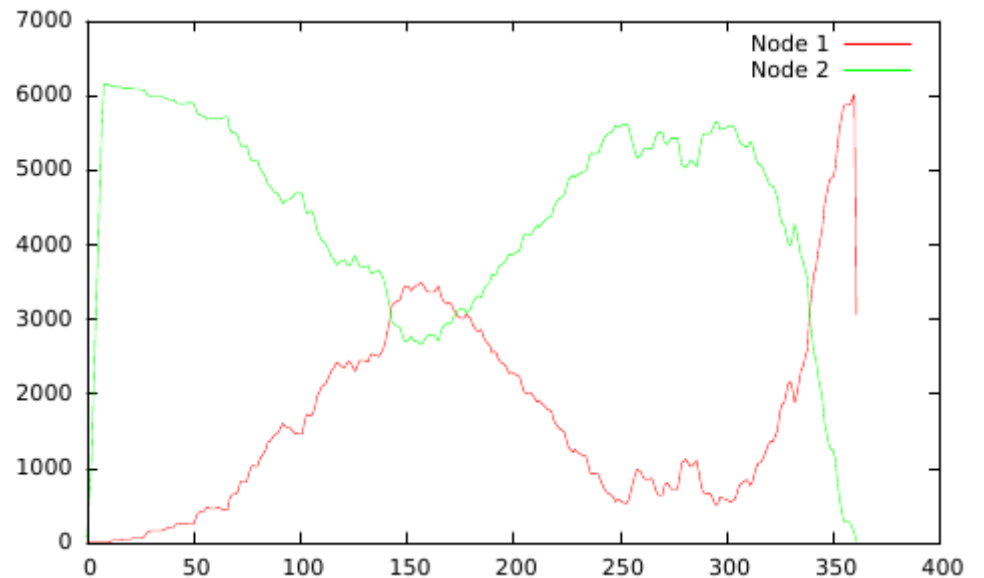
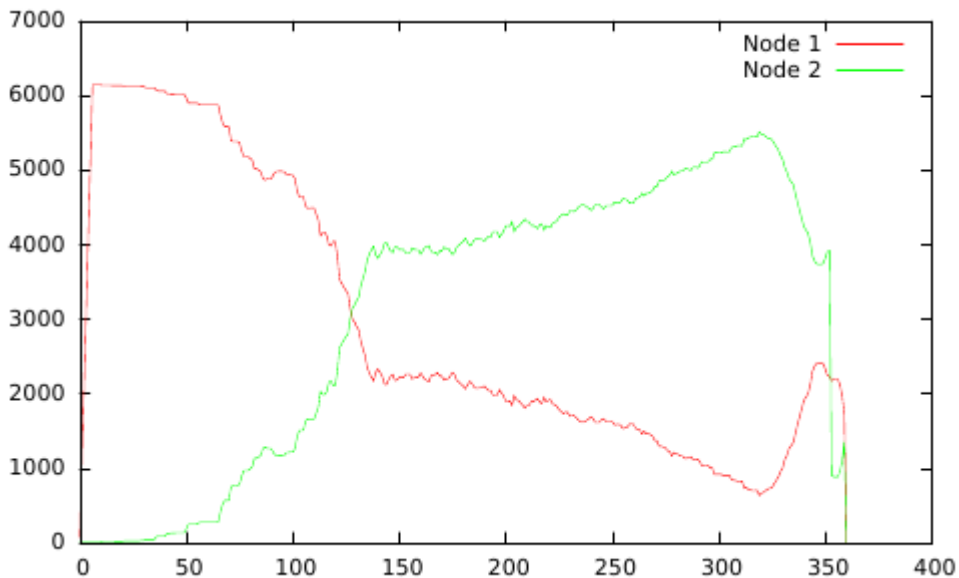
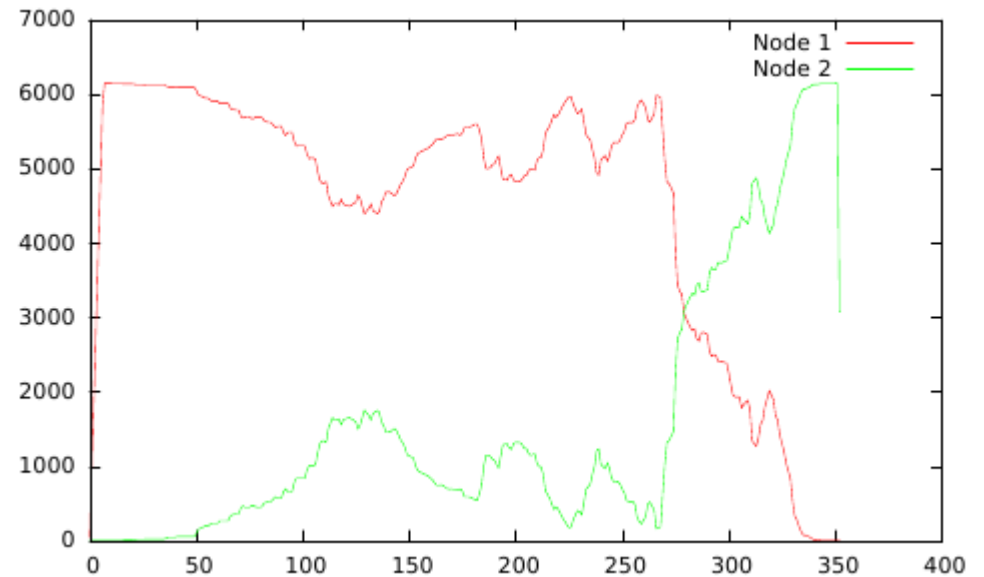
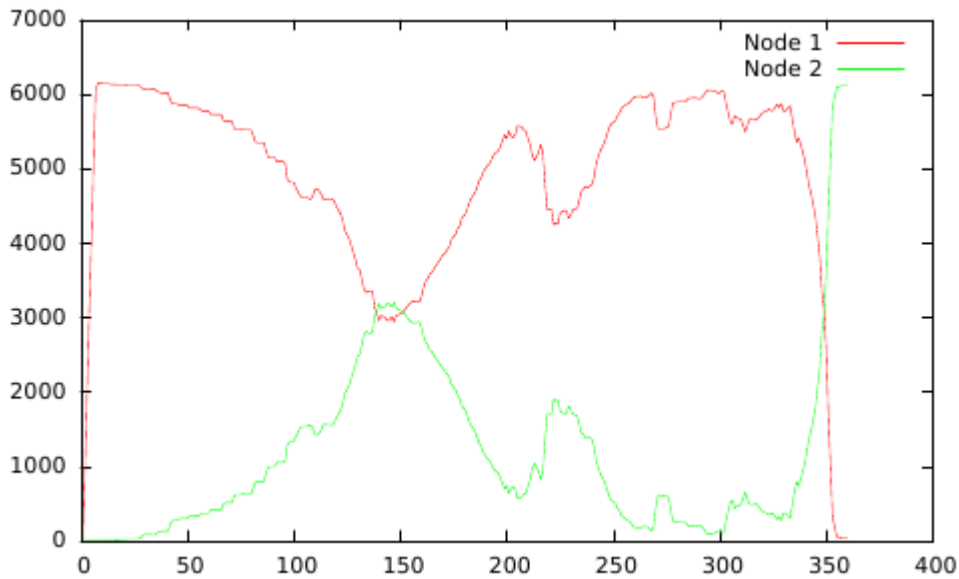
AutoNUMA23 numa02



AutoNUMA23 numa02_SMT

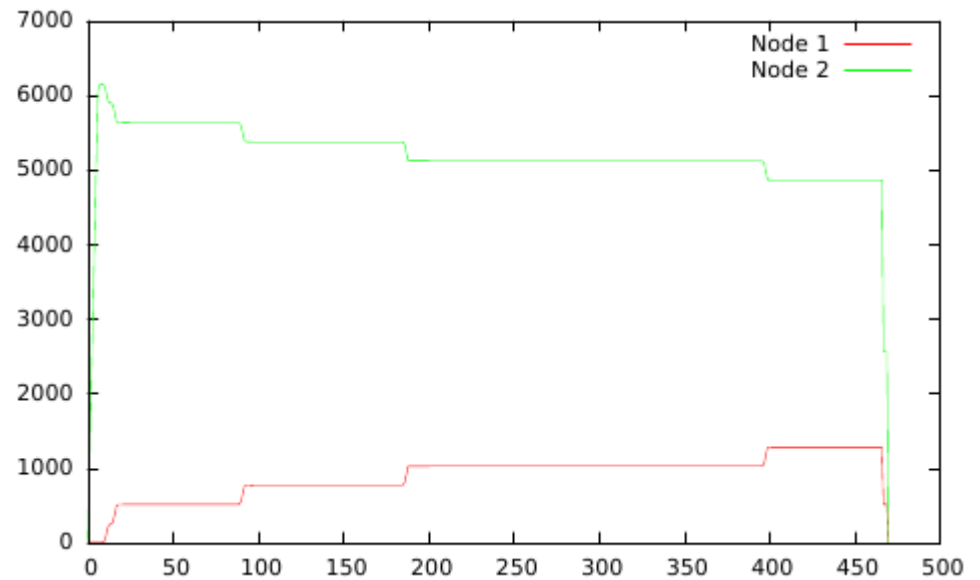
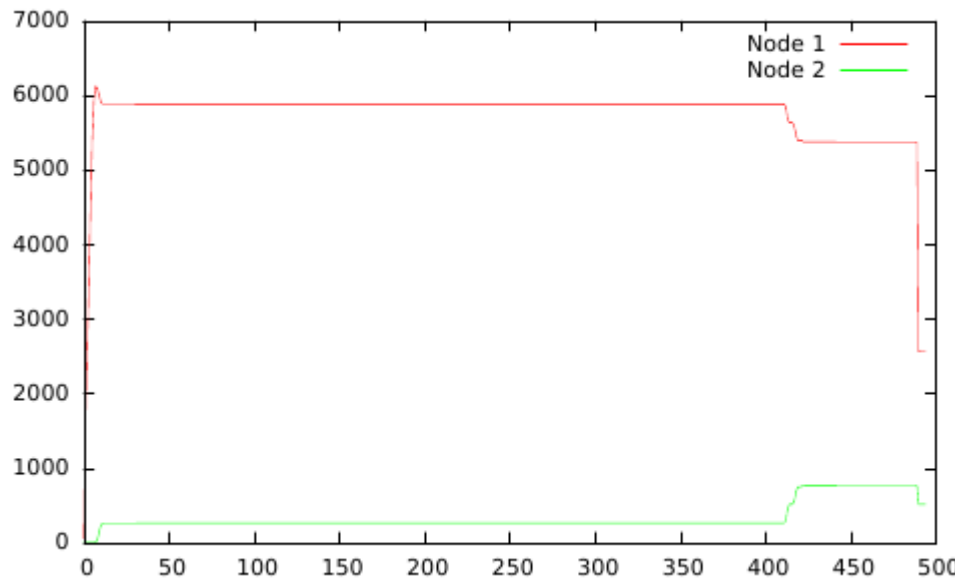
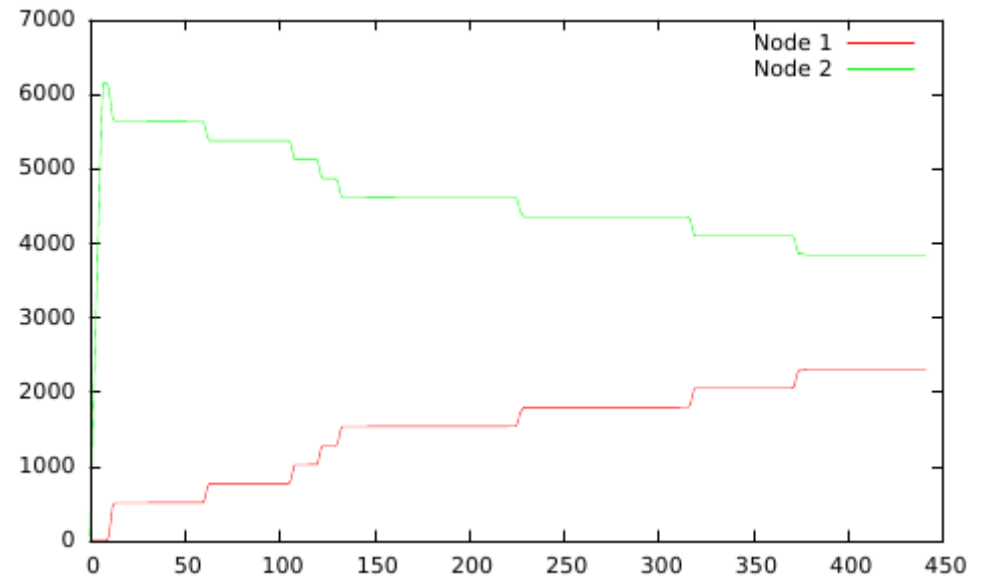
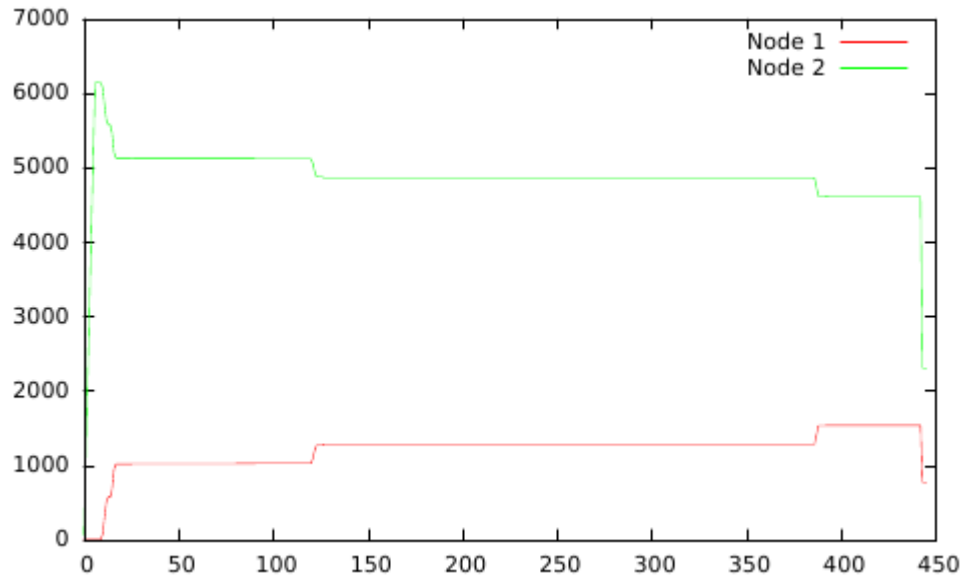


sched-numa-rewrite numa01

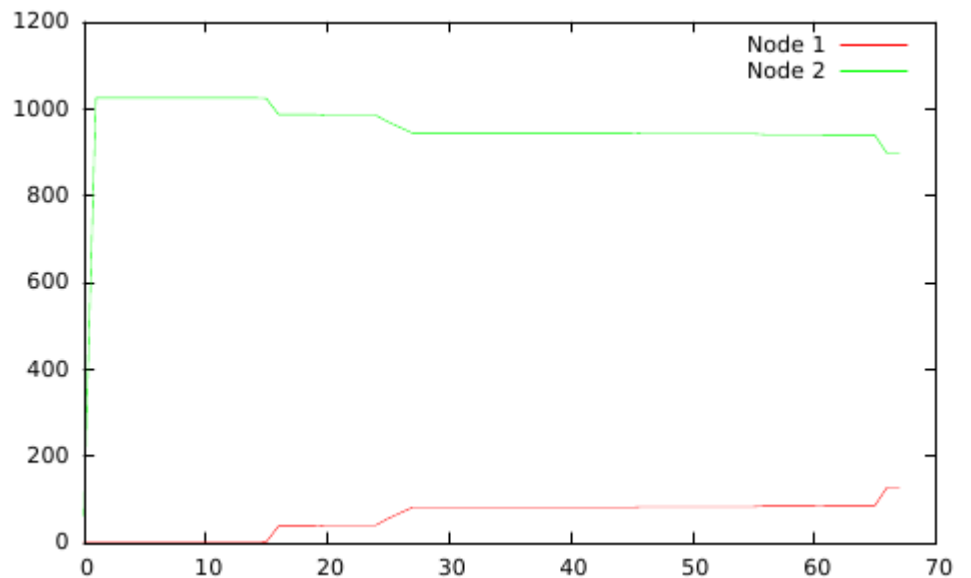
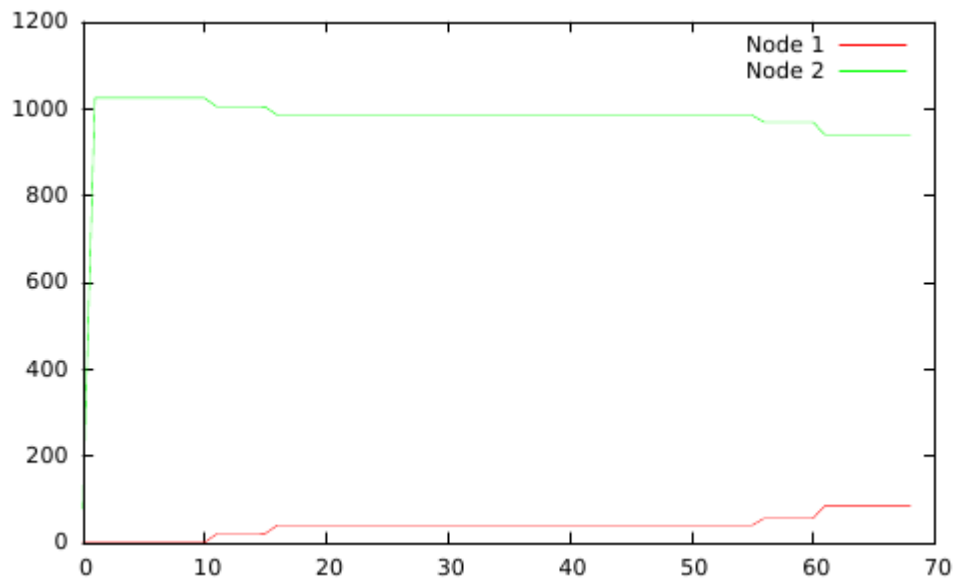
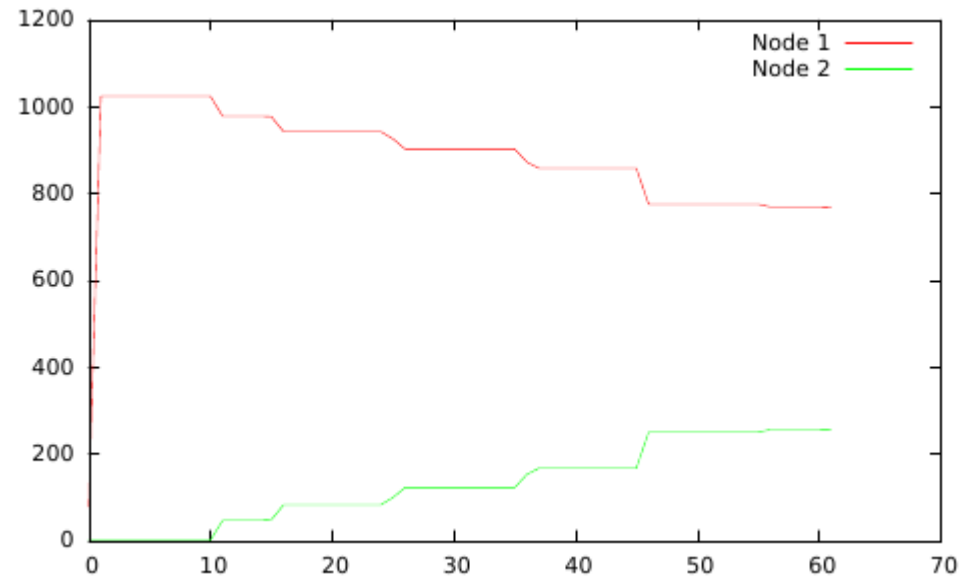
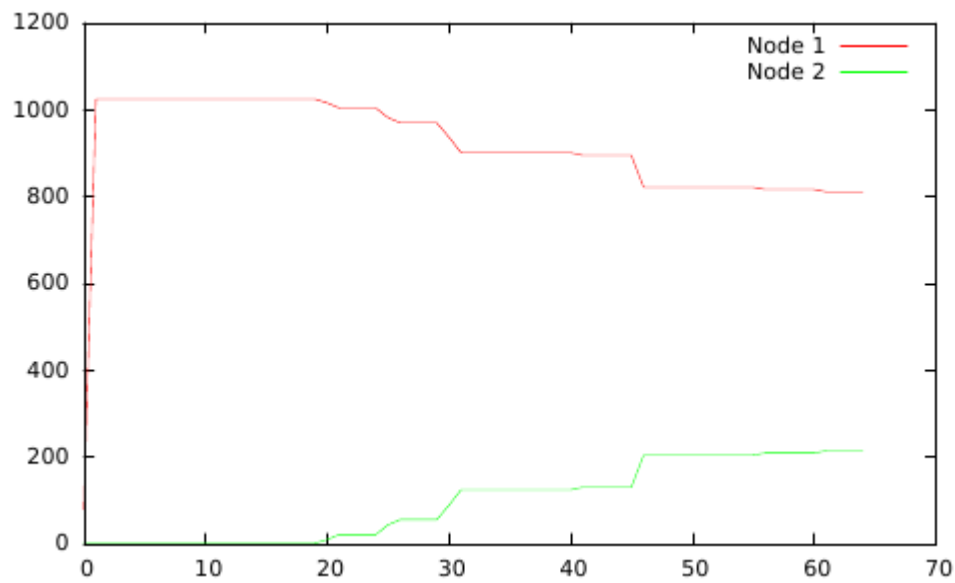


2 nodes

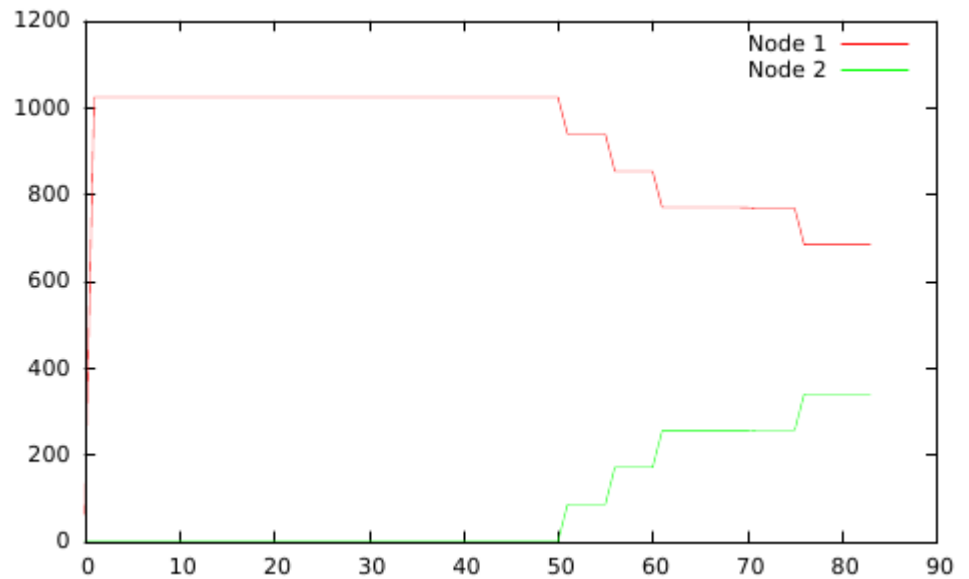
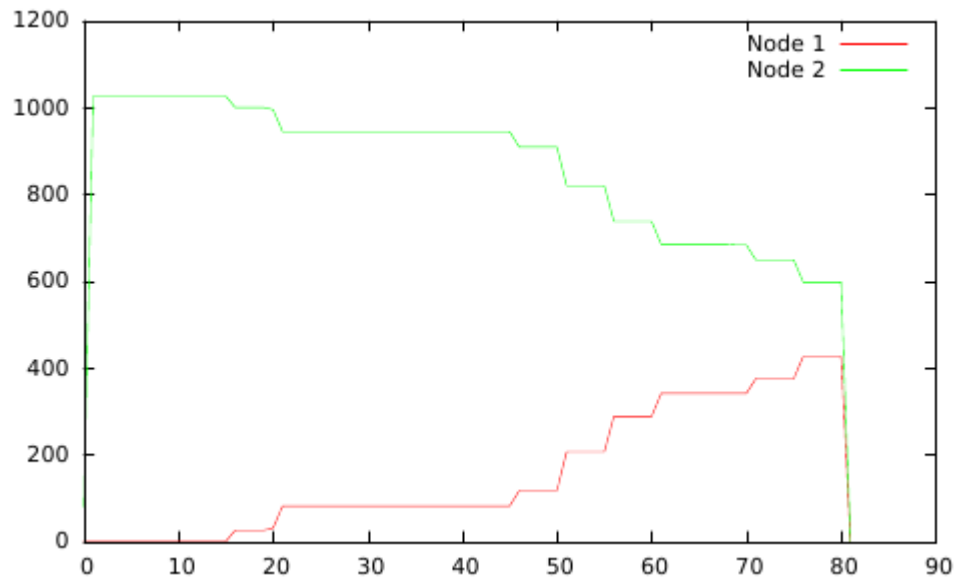
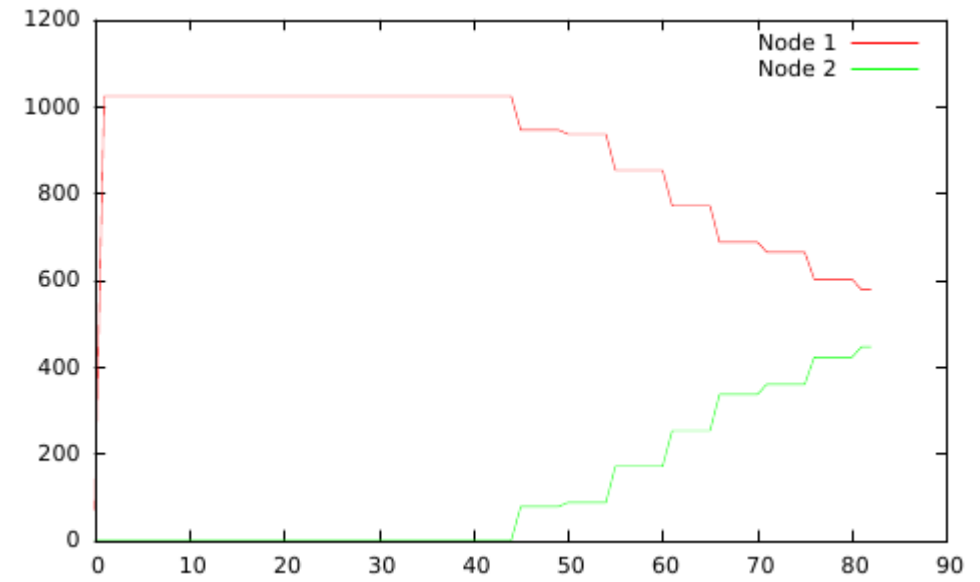
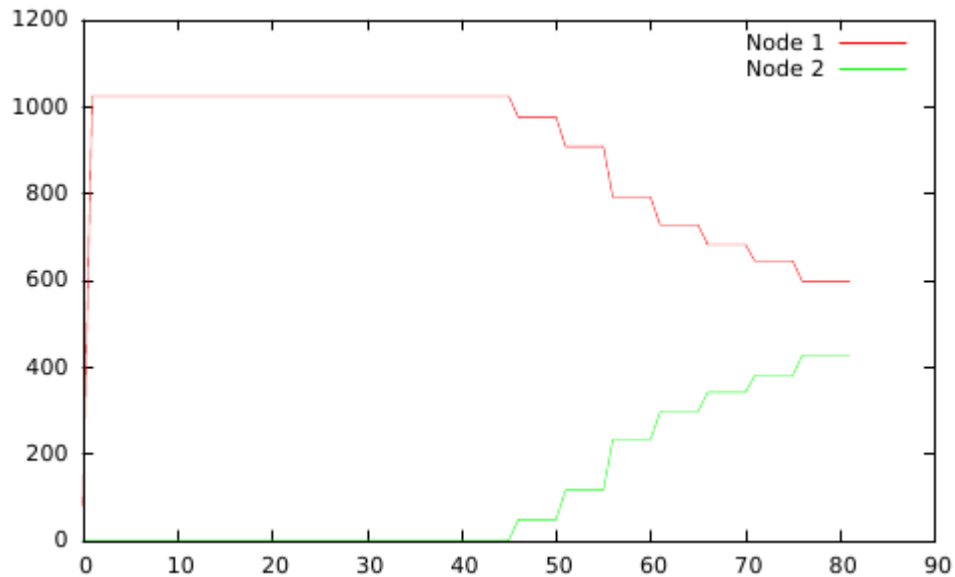
sched-numa-rewrite numa01_TH..



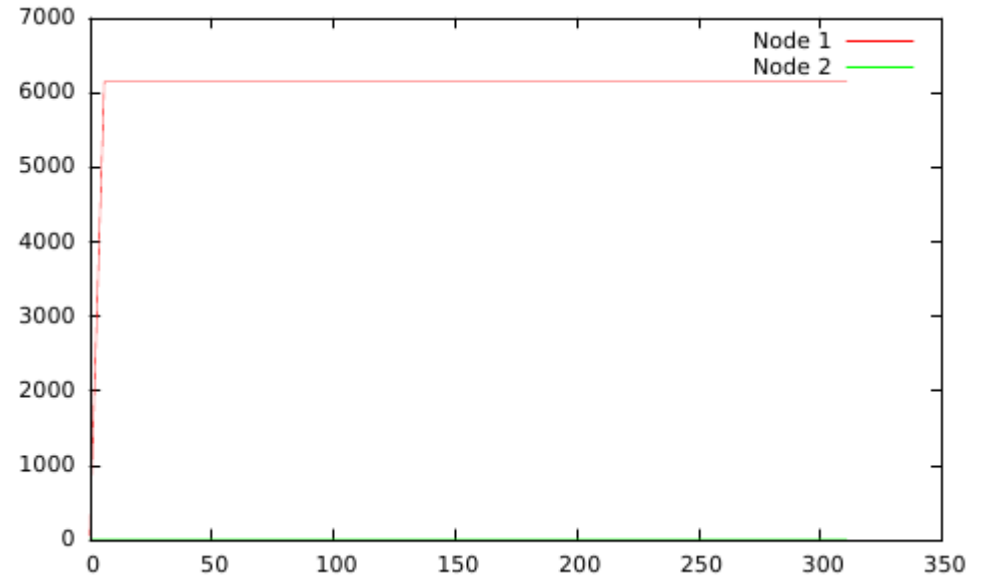
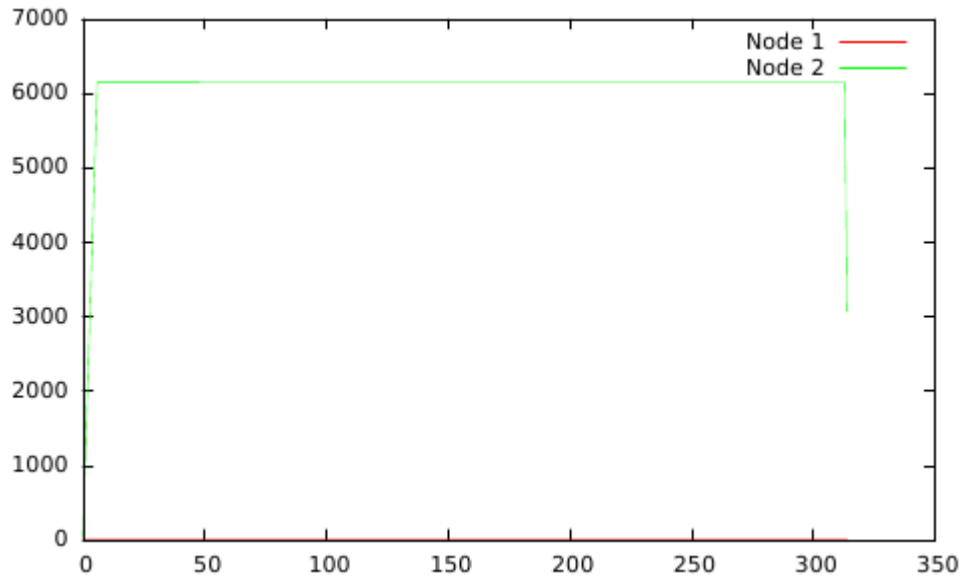
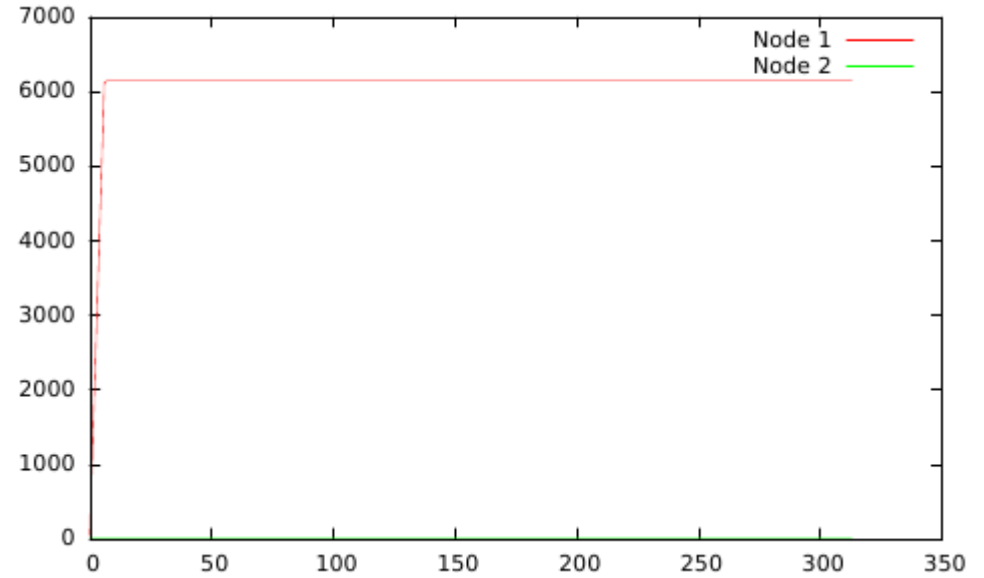
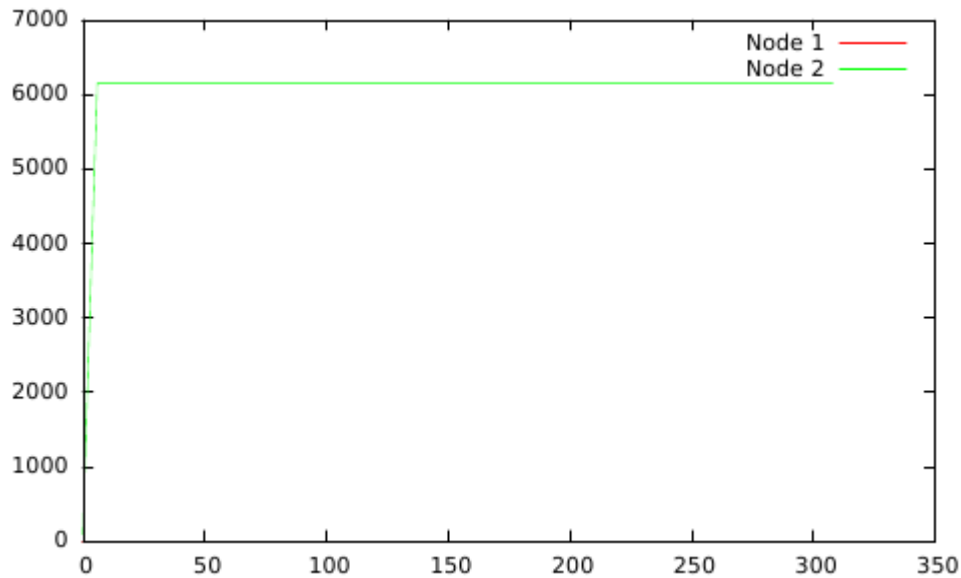
sched-numa-rewrite numa02



sched-numa-rewrite numa02_SMT

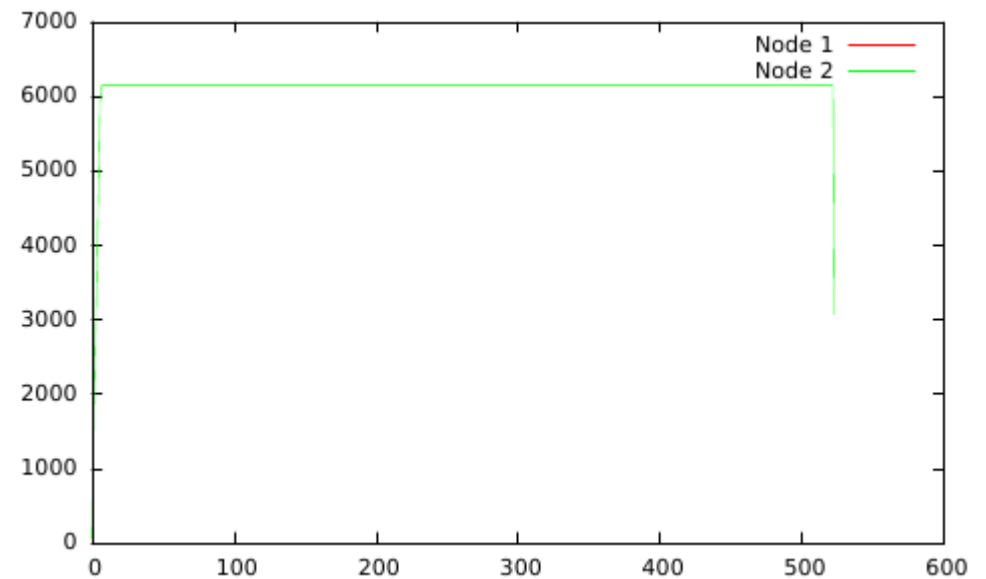
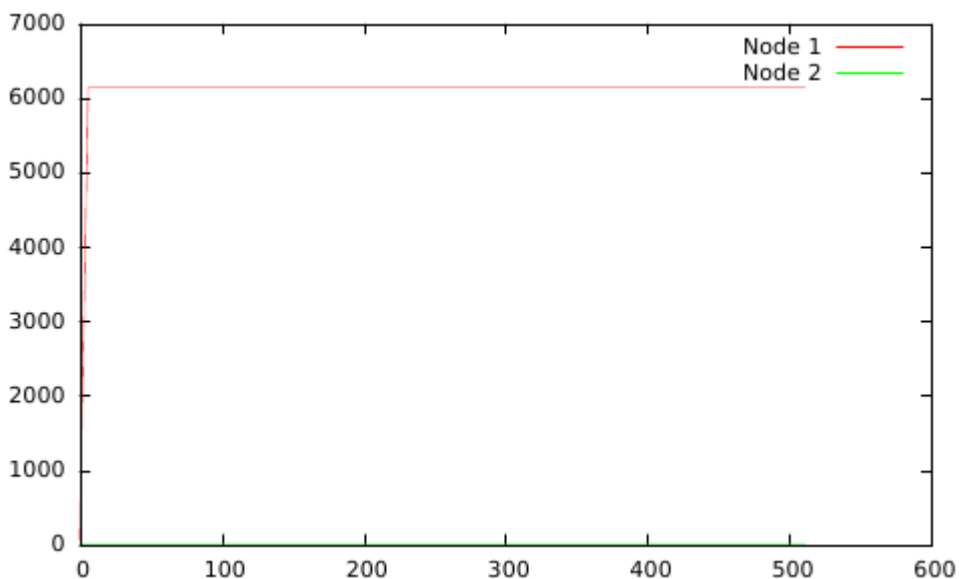
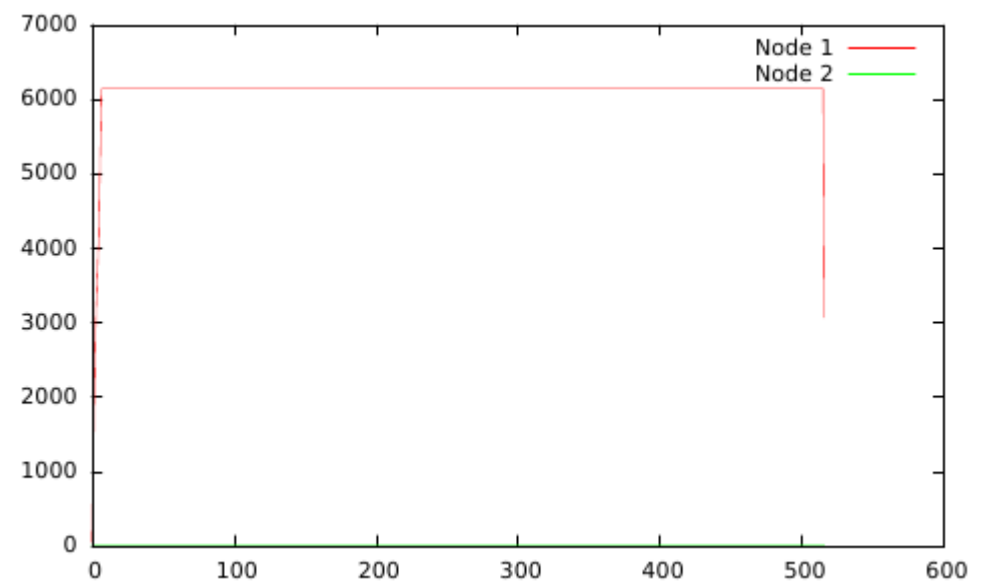
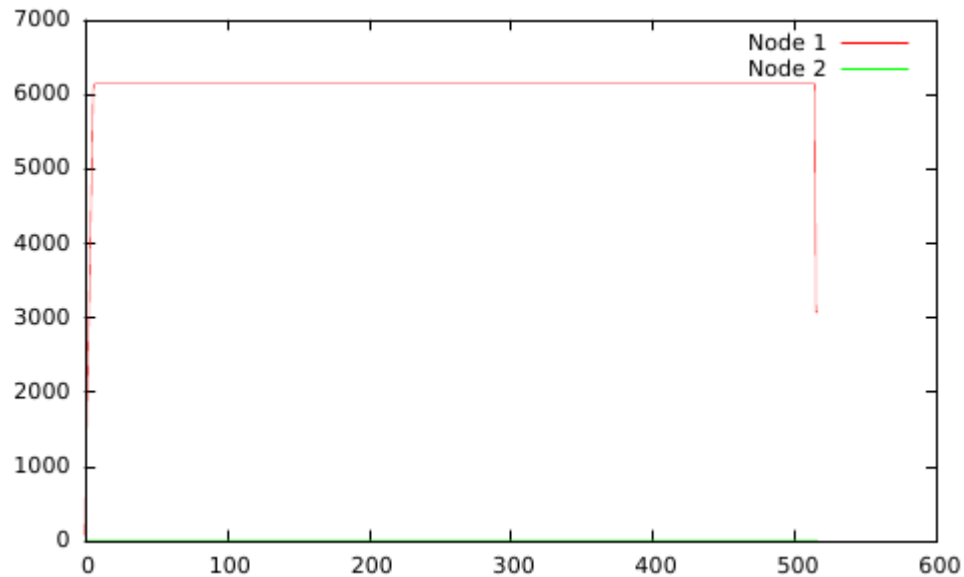


3.6-rc1 numa01



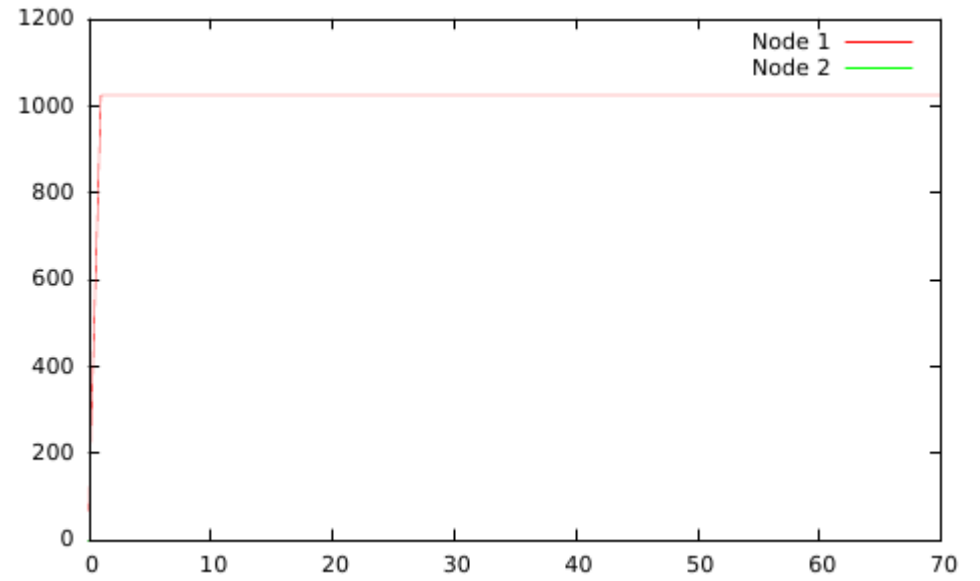
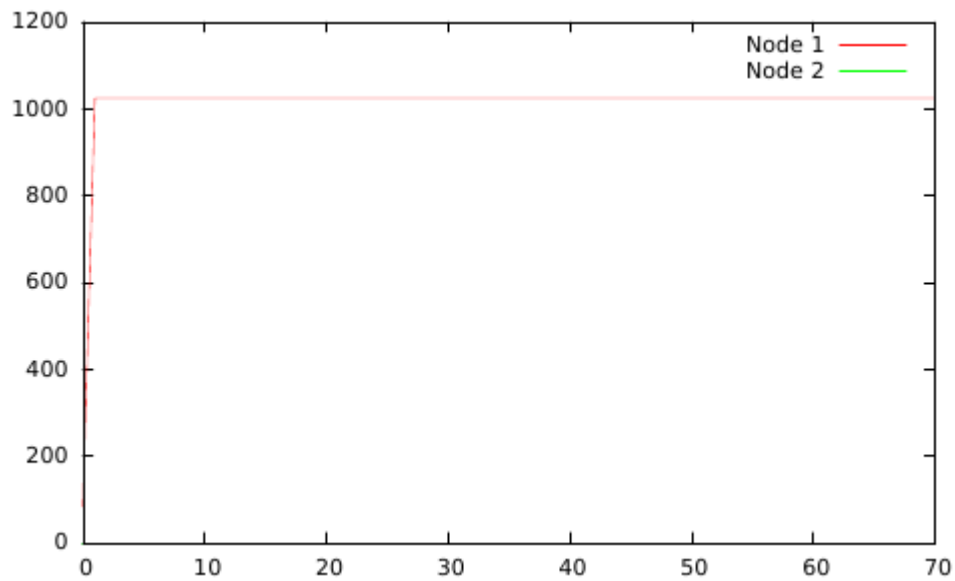
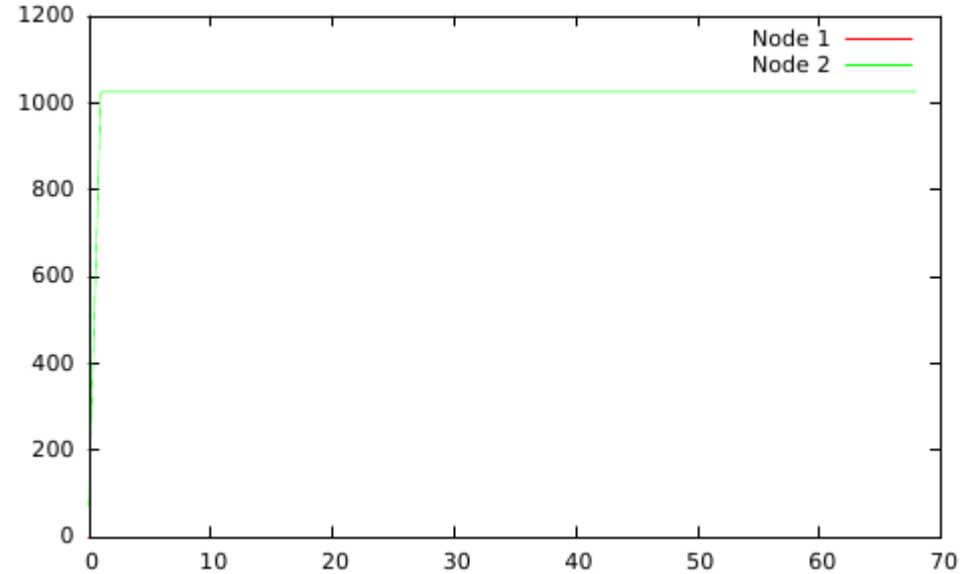
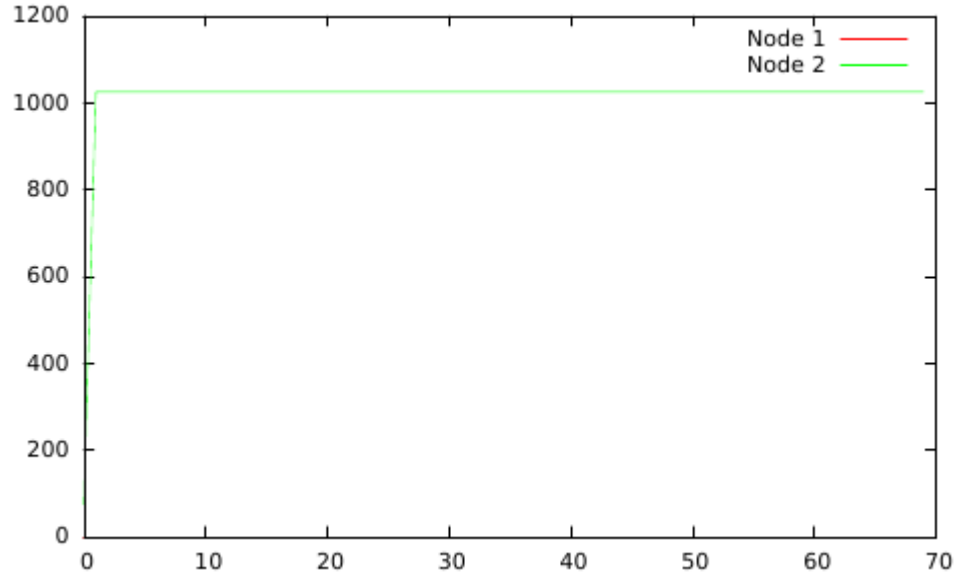
2 nodes

3.6-rc1 numa01_THREAD_ALLOC



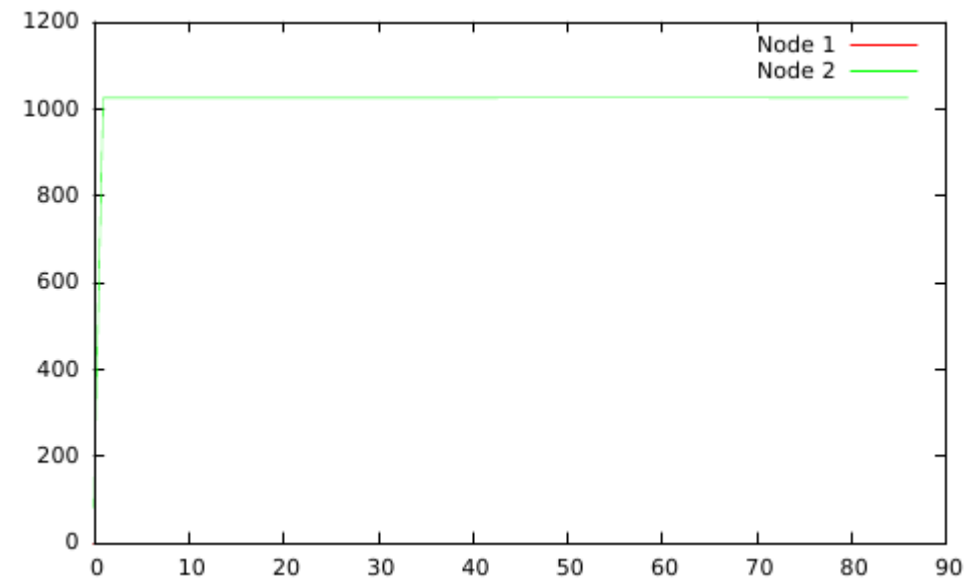
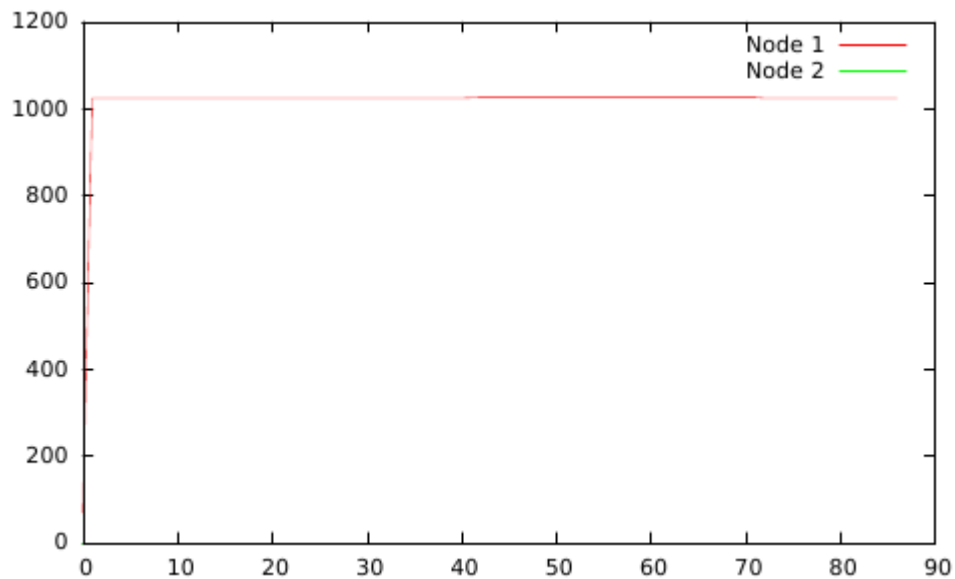
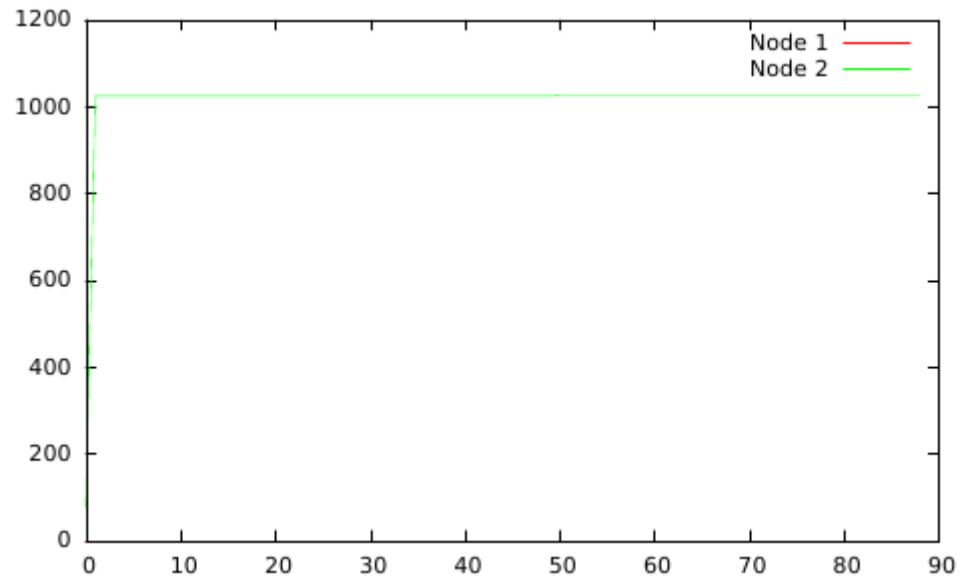
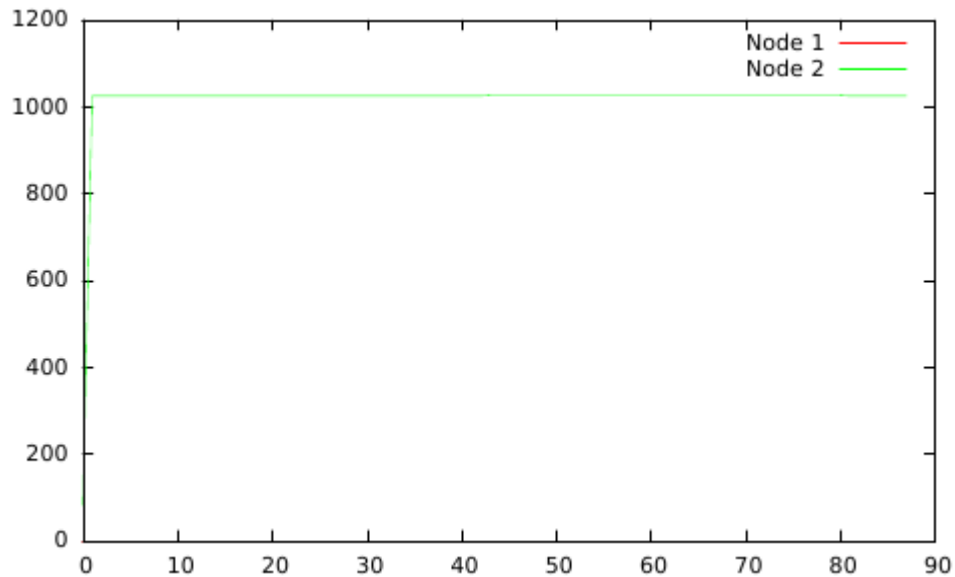
2 nodes

3.6-rc1 numa02



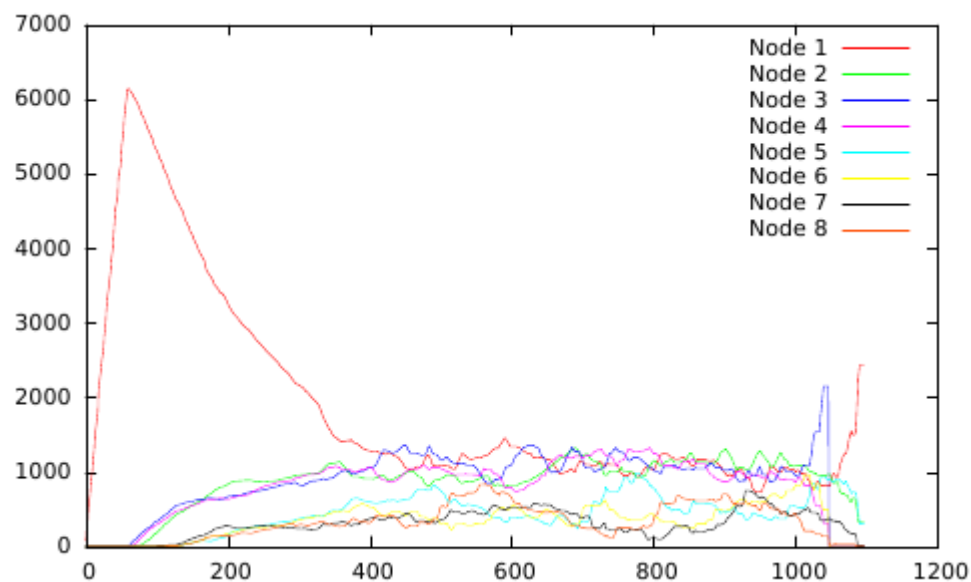
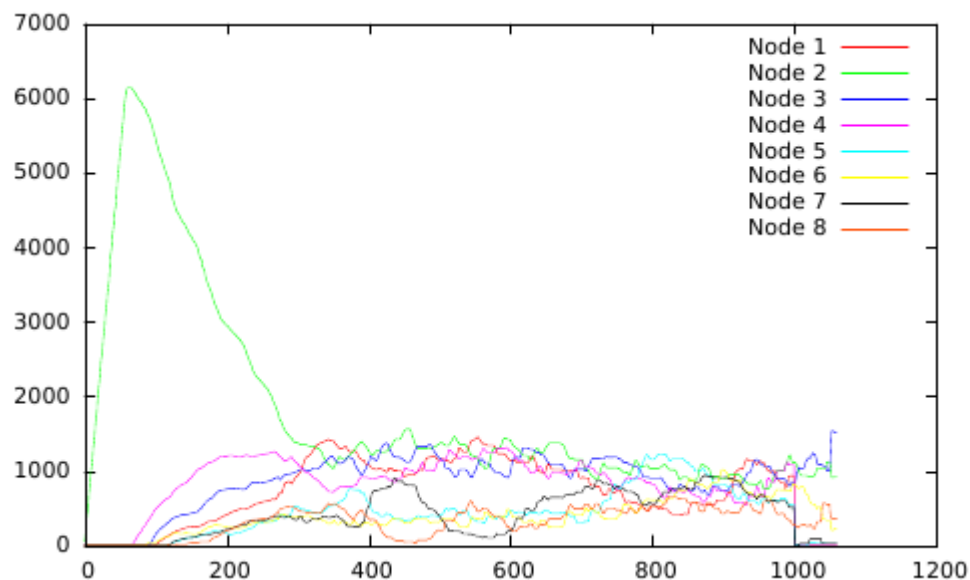
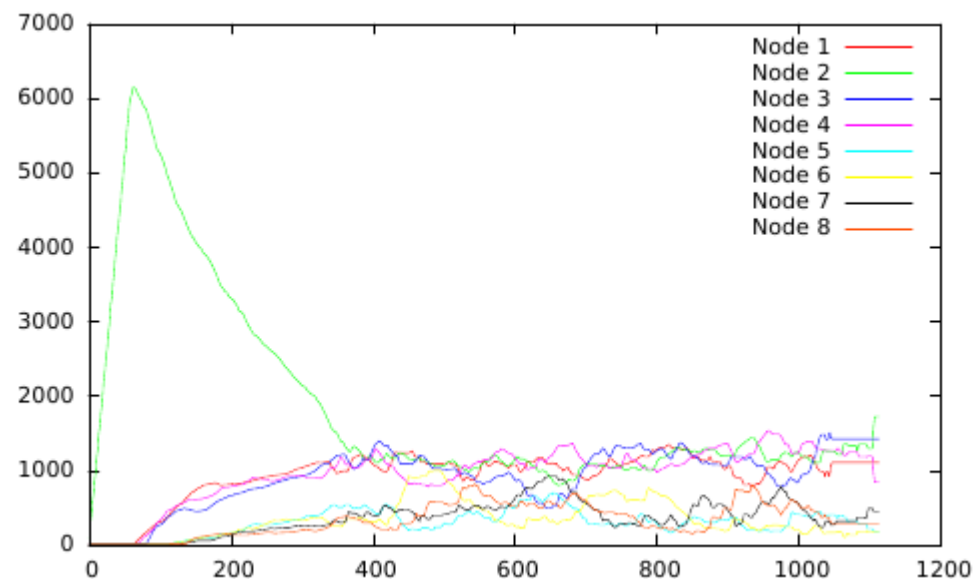
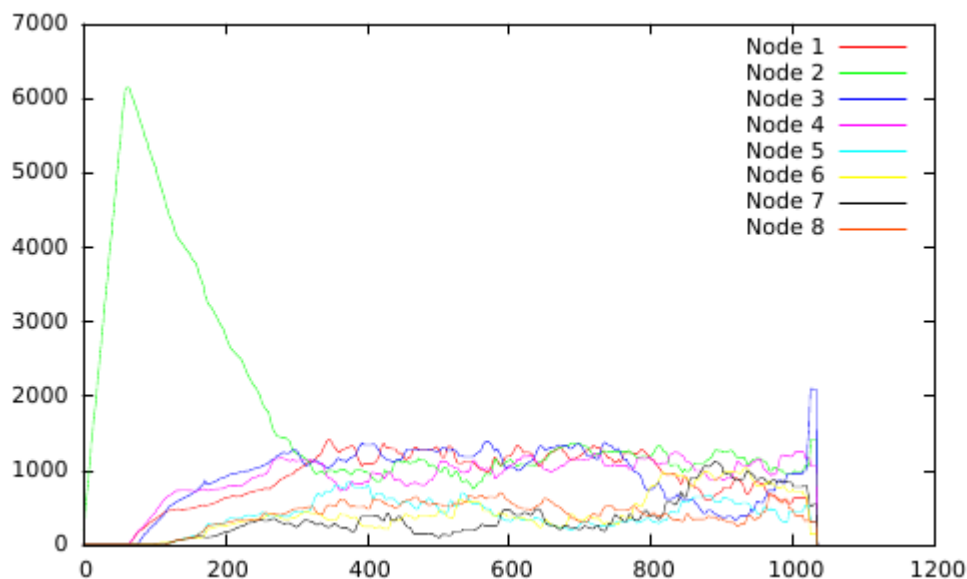
2 nodes

3.6-rc1 numa02_SMT

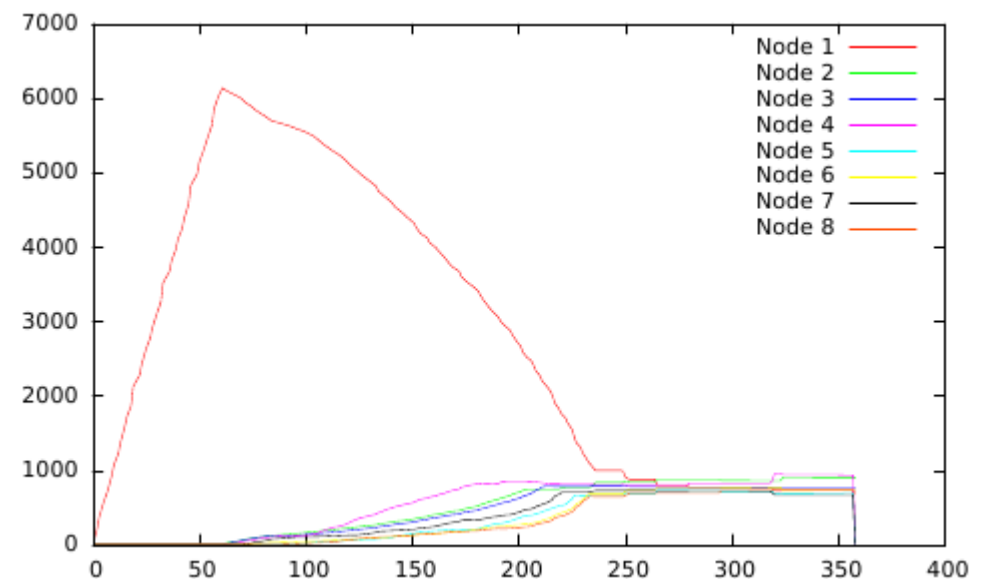
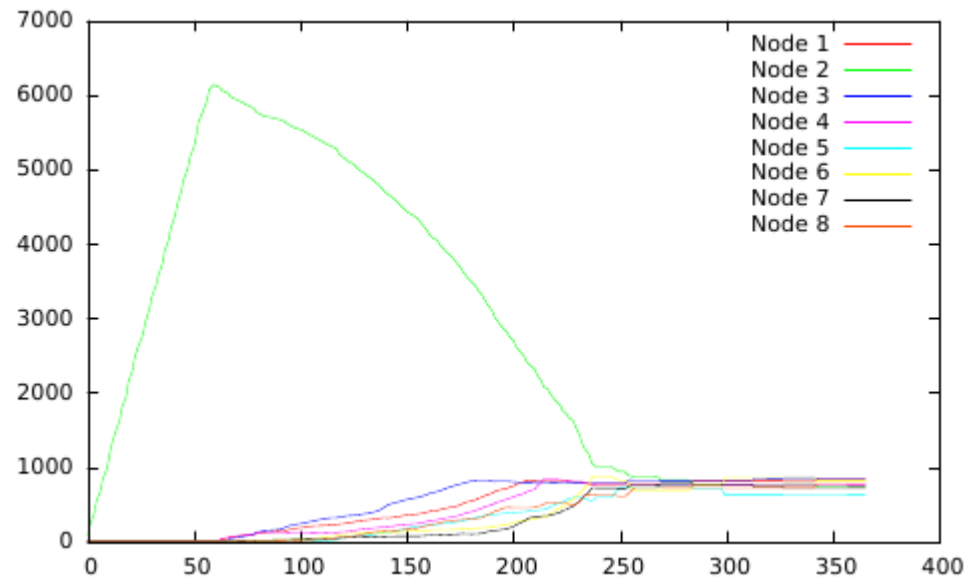
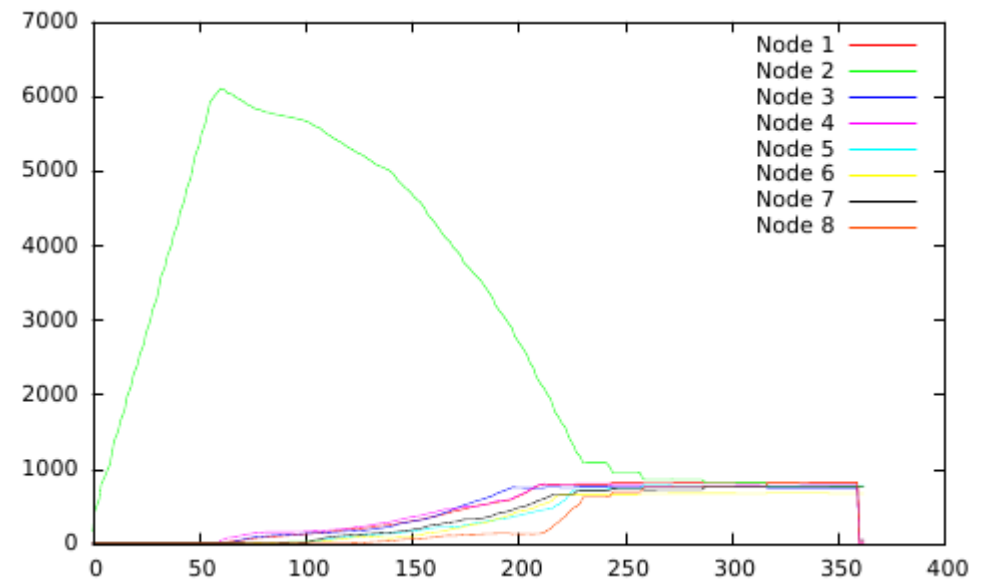
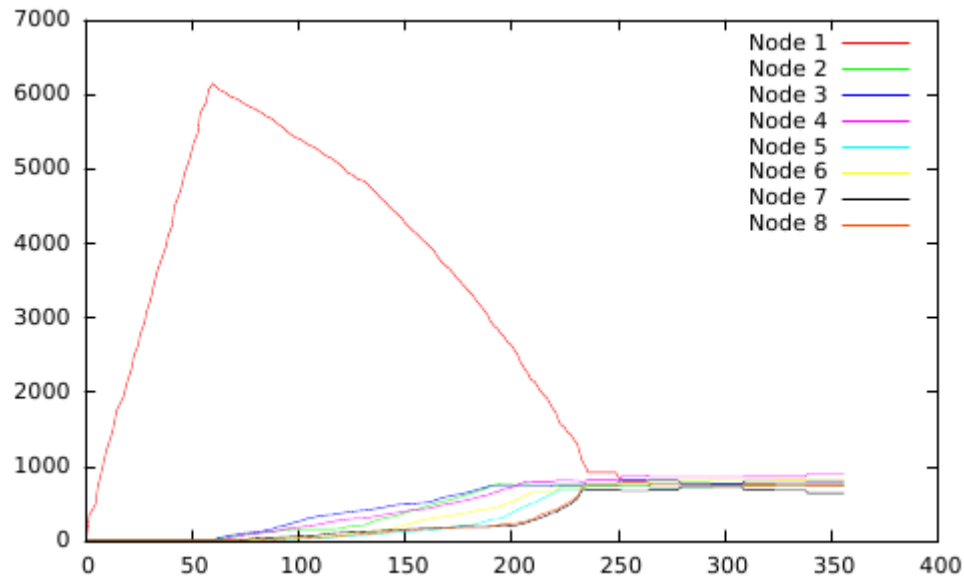


2 nodes

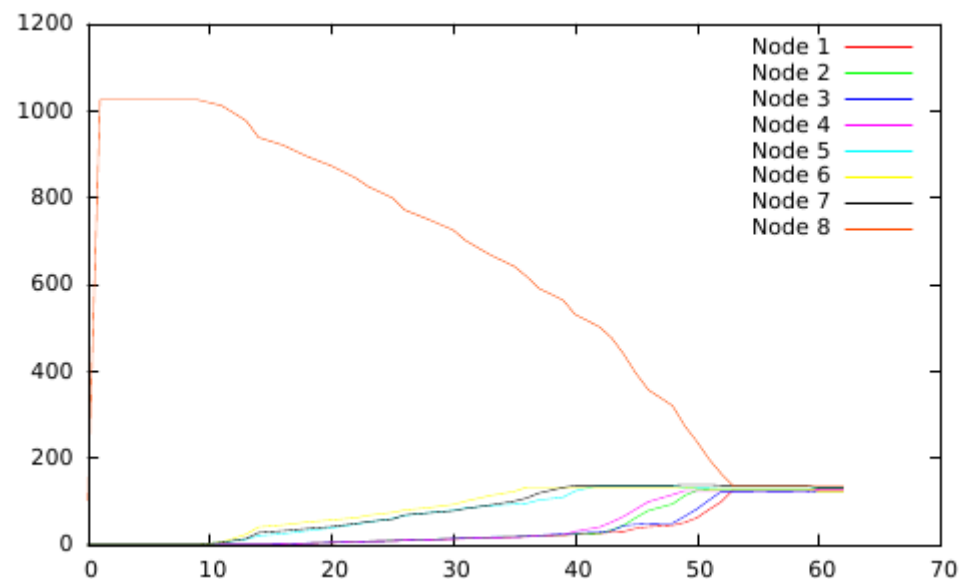
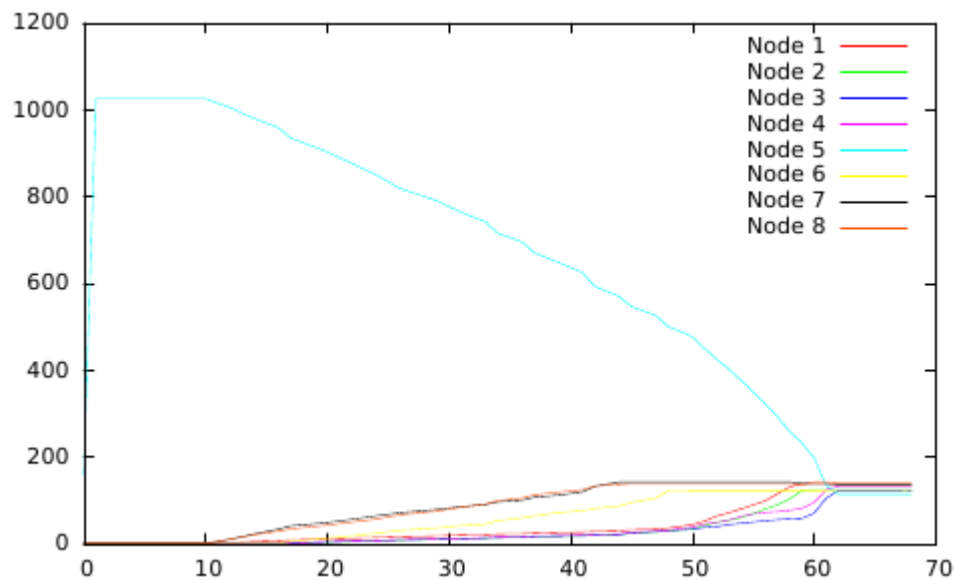
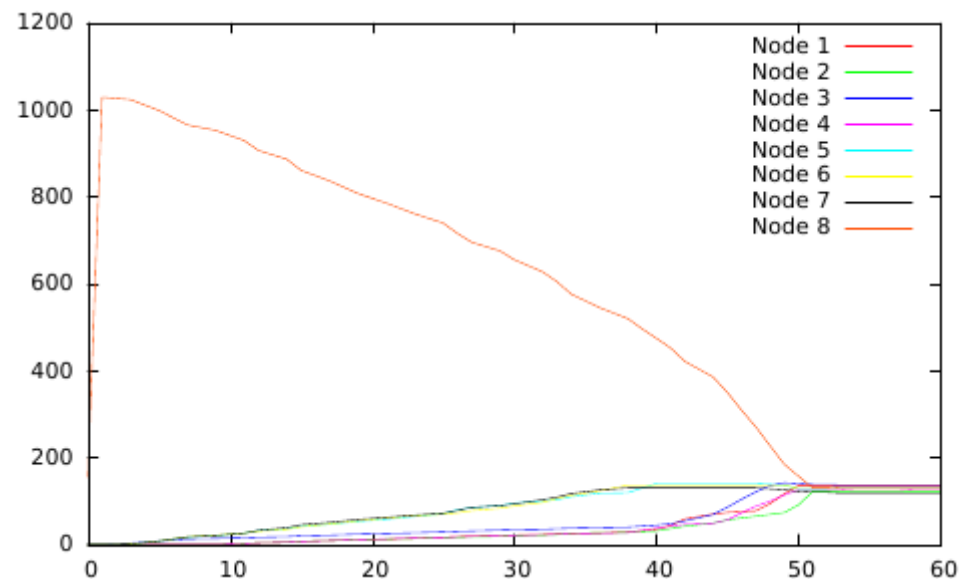
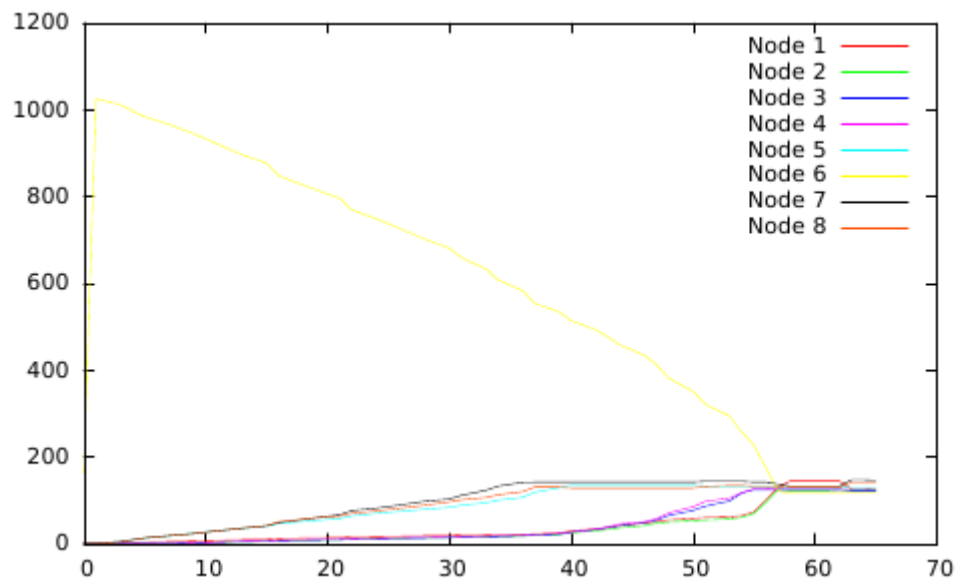
AutoNUMA23 numa01



AutoNUMA23 numa01_THREAD..

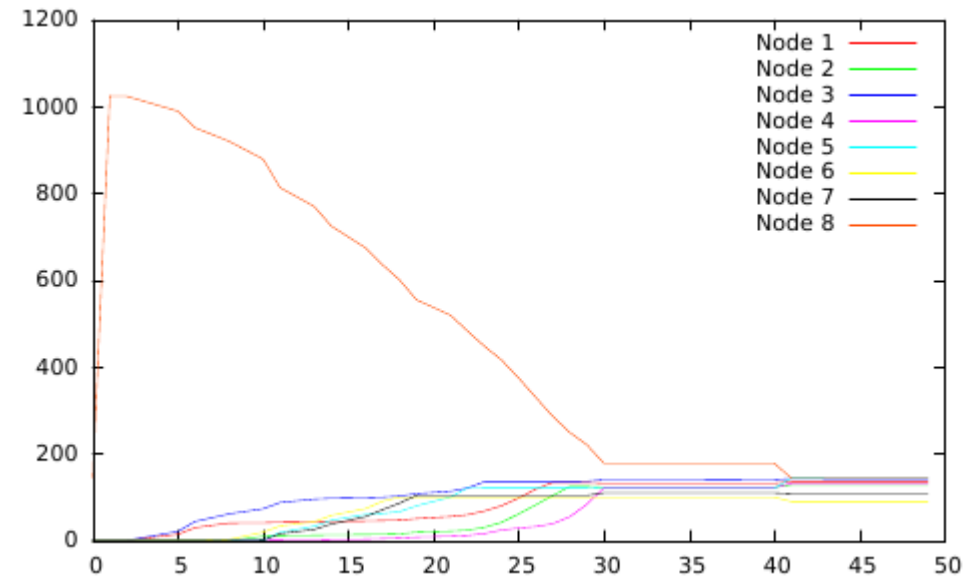
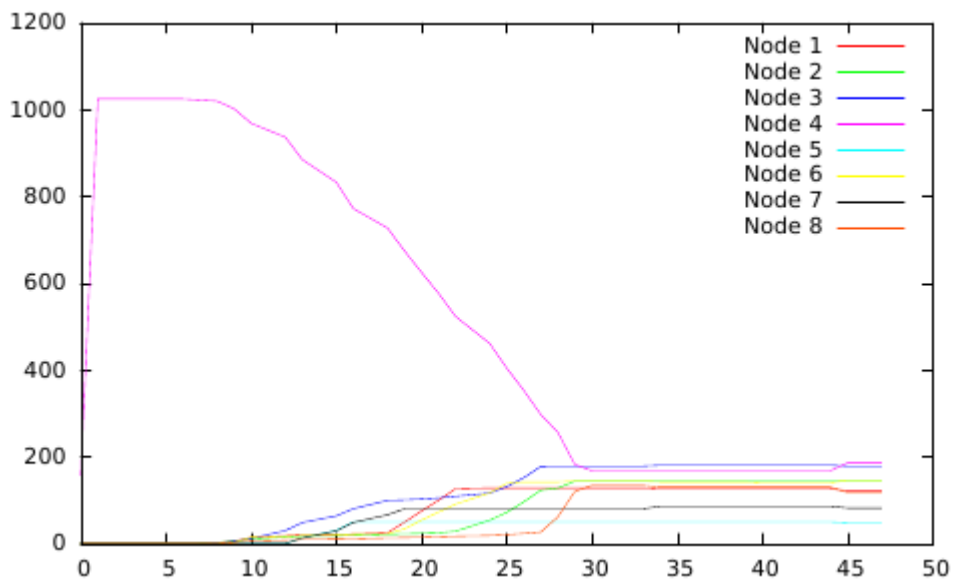
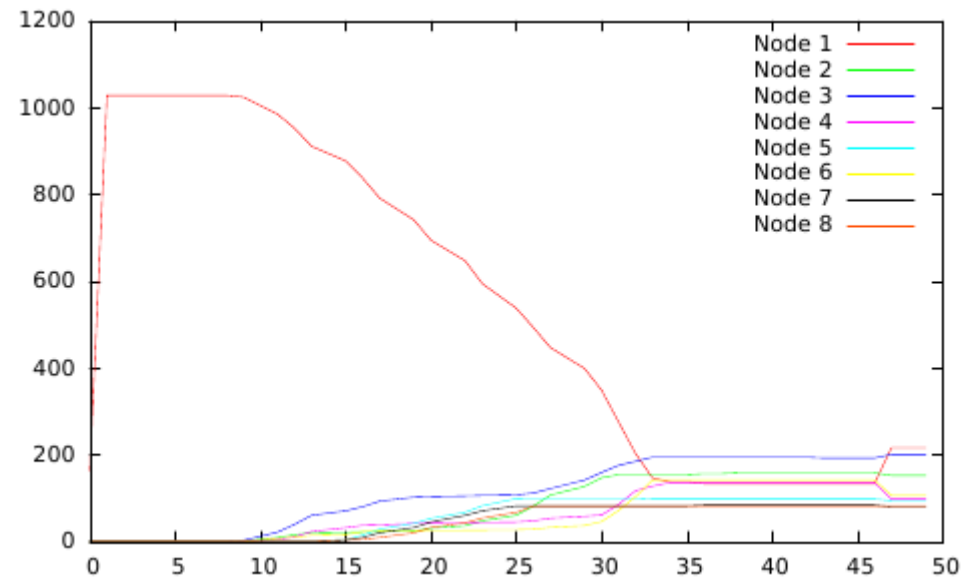
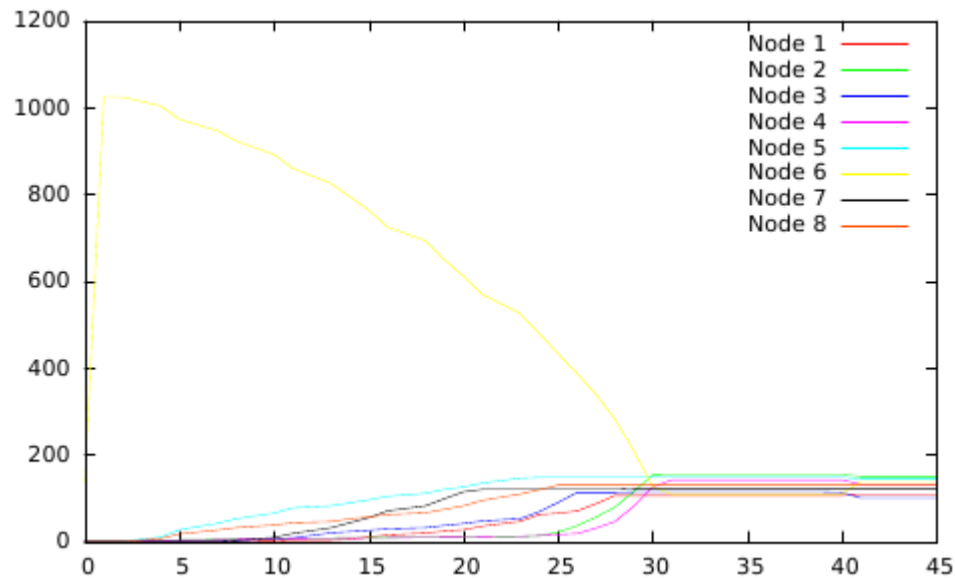


AutoNUMA23 numa02



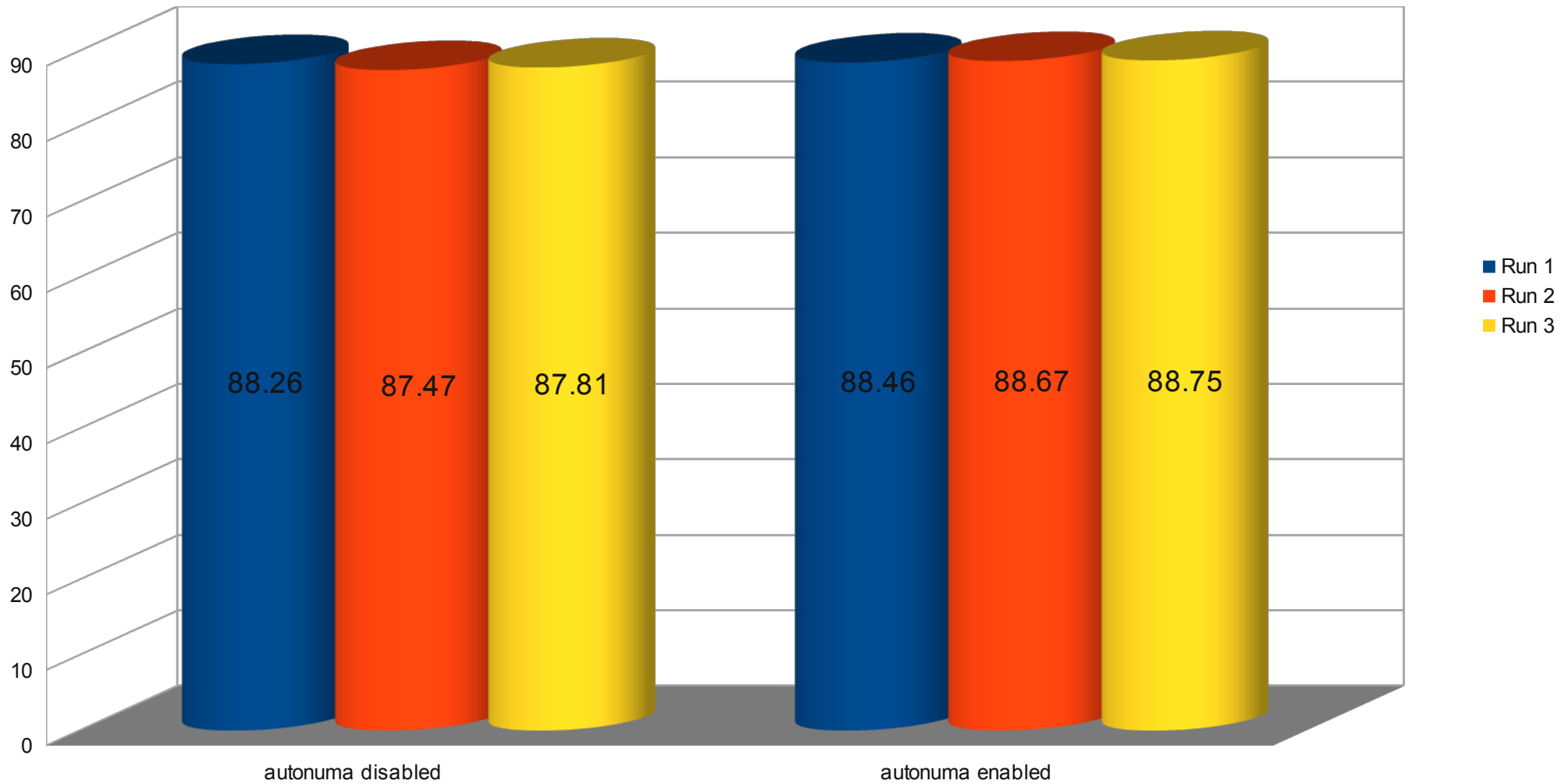
8 nodes

AutoNUMA23 numa02_SMT

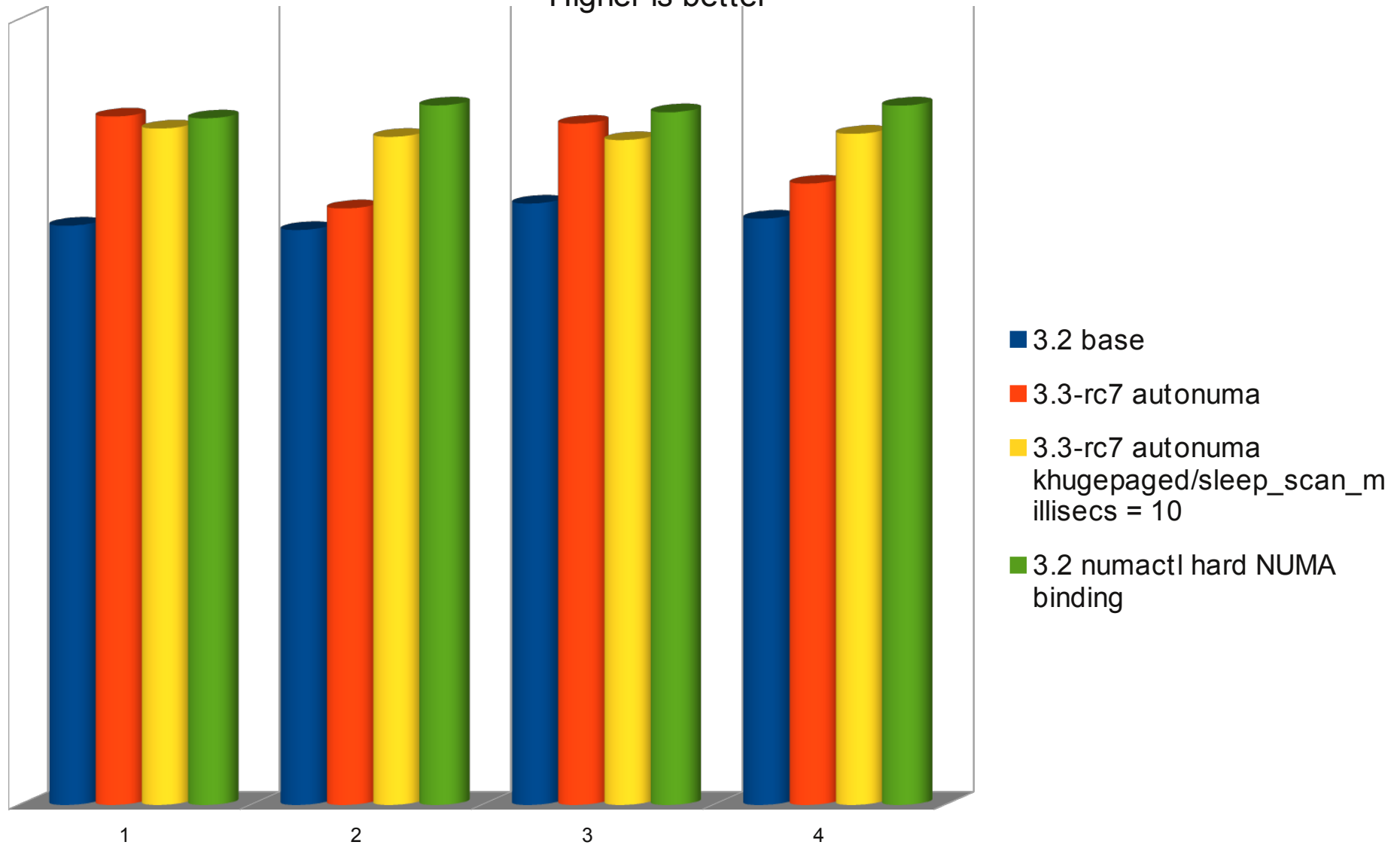


Kernel build time in seconds on tmpfs (make -j32)
Autonuma enabled includes one knuma_scand pass every 10sec
Lower is better

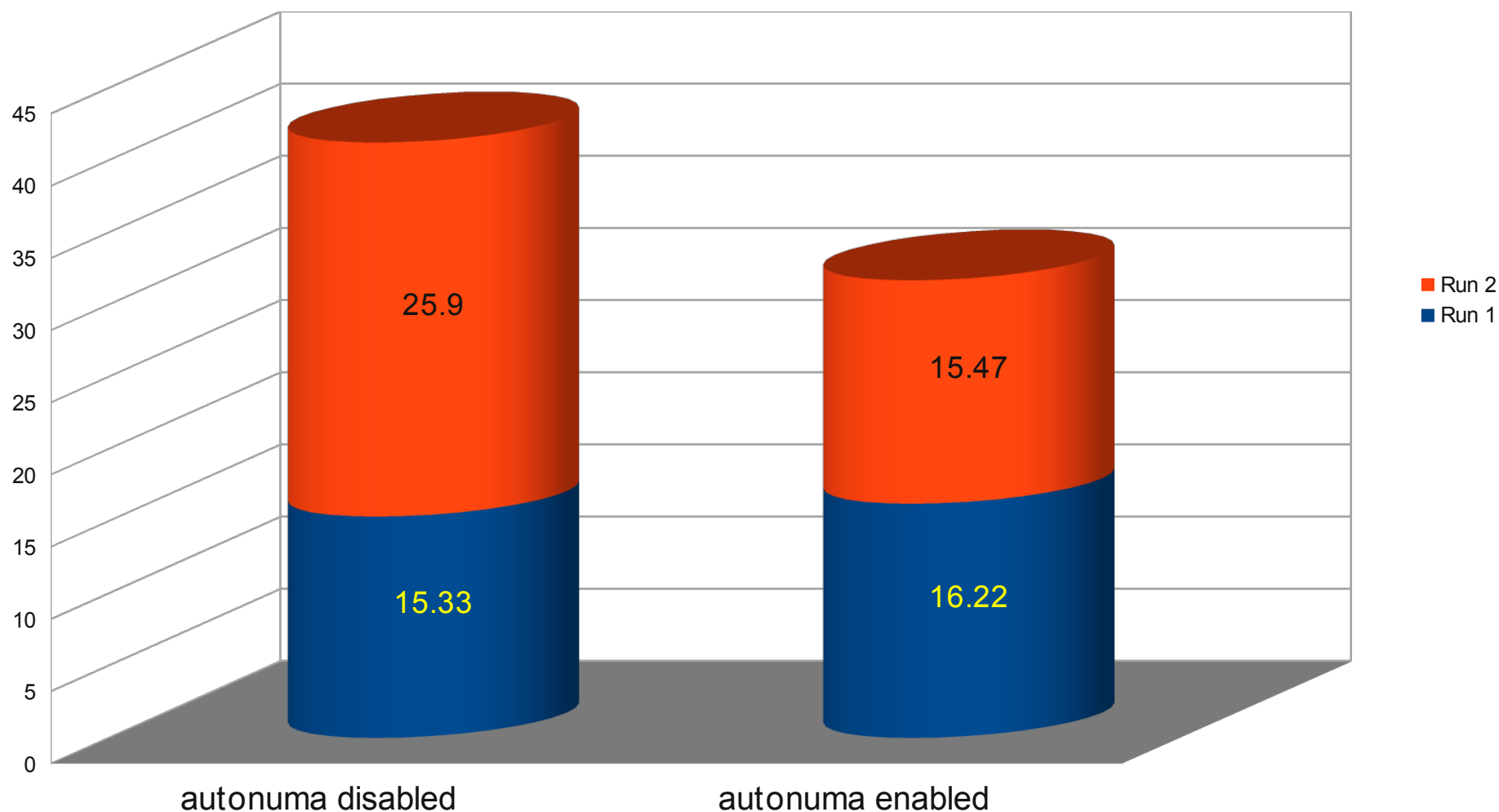
Worst possible case for AutoNUMA (gcc too short lived)
Average increase in build time 0.88%



SPECjbb results 2 NUMA nodes, 8 CPUs per node, 16 CPUs total
THP enabled, no virt
Higher is better

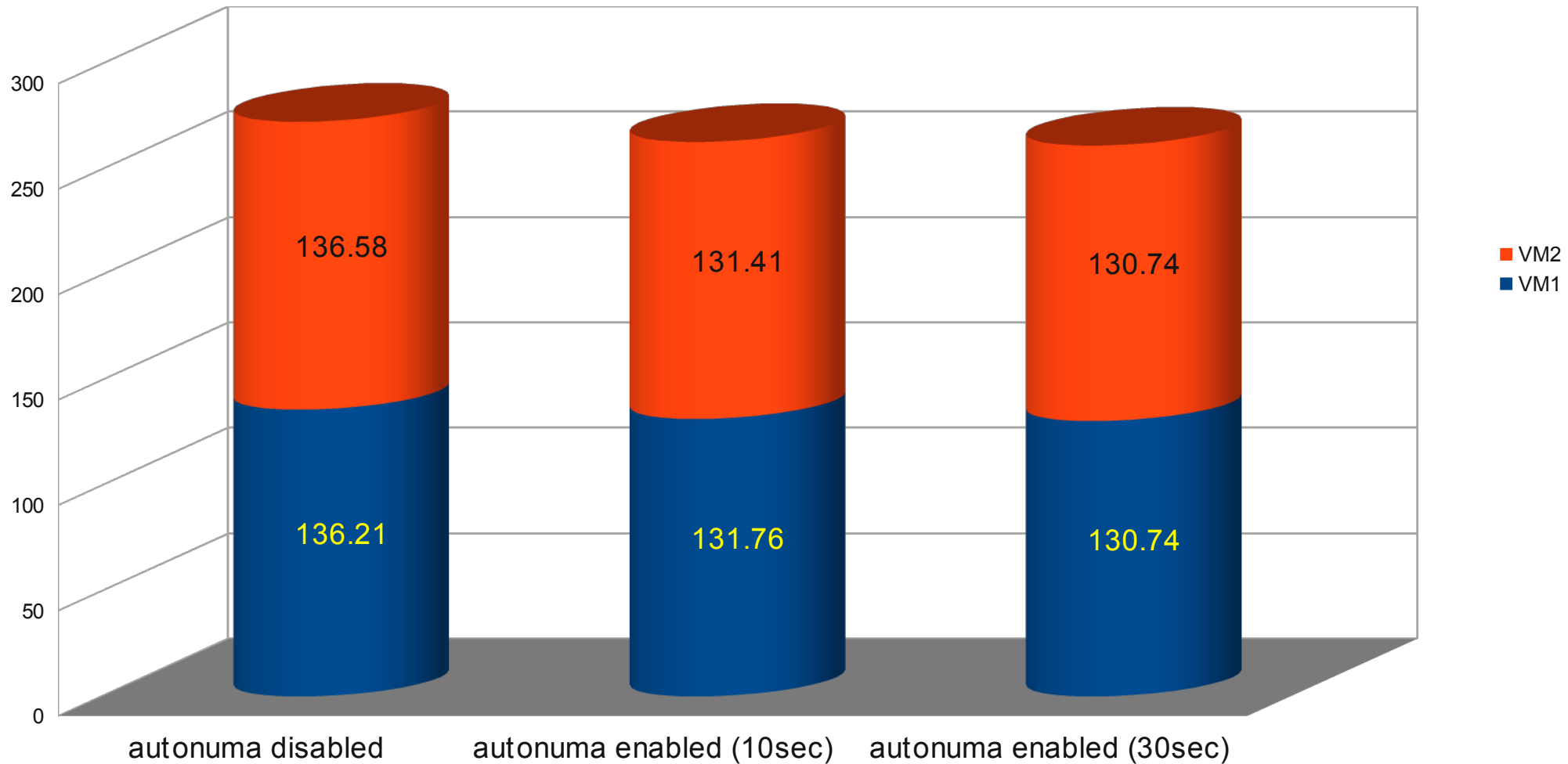


Virt guest "memhog -r100 1g" (autonuma includes 1 knuma_scand pass every 10 sec)
KVM host autonuma enabled/disabled, THP enabled
Guest VM fits in one host NUMA node
Lower is better



kernel build -j16 in parallel in 2 KVM (both in tmpfs, in a loop started in sync)
Both guest VM fits in one host NUMA node
autonuma/knuma_scand/scan_sleep_pass_millisecs = 5000 | 15000 (10sec | 30sec)
Lower is better

Host autonuma enabled/disabled, THP on, 12 vcpu per guest, 24 CPUs total on host



KVM/KSM/THP bench run by IBM

SPECjbb2005 on Linux 3.4-rc2

1. Higher is better
2. Can't compare VM1 scores with VM2/VM3 scores.
3. Can compare VM2 score with VM3 score (should be equal)
4. One Node = 12 GiB (slightly > 12 GB) with 6 cores + 6 hyperthreads

■ VM 1 (~ 1 node: 12 GB/14 GB with 12 vCPUs) ■ VM 2 (4 GB/3 GB with 6 vCPUs) ■ VM 3 (4 GB/3 GB with 6 vCPUs)

