# Running Oracle on a 32 socket server with 24T of memory

Frits Hoogland

# `whoami`

- Frits Hoogland
- Working with Oracle products since 1996

- Blog: http://fritshoogland.wordpress.com
- Twitter: @fritshoogland
- Email: frits.hoogland@accenture.com
- Oracle ACE Director
- OakTable Member

# Goals & prerequisites

- Goal: Learn about characteristics of a huge system.

- Prerequisites:
  - Basic understand of hardware architecture.
  - Basic understanding of C and Linux.
  - Basic understanding of Oracle running on Linux.

# SGI UV300

- CPUs: Intel(R) Xeon(R) CPU E7-8890 v2 @ 2.80GHz
  - (Ivy Bridge EX)
- 32 sockets
- 480 cores (15 cores/socket)
- 960 threads (intel hyper threading)

- 32s480c960t

enkitec

# Memory

- Total memory size: 24TB
- Memory is local to a socket

- $(24*1024)/32 = 768$ GB / socket

```
# numactl --hardware | grep size
node 0 size: 753624 MB
...
node 31 size: 753648 MB
```

# Memory

- Memory is DDR3 @ 1333Mhz

```
# dmidecode | grep -A13 'Memory Device'
Memory Device
    Array Handle: 0x0001
    Error Information Handle: Not Provided
    Total Width: 72 bits
    Data Width: 64 bits
    Size: 32 GB
    Form Factor: DIMM
    Set: 8
    Locator: DIMMD2
    Bank Locator: MEM8
    Type: DDR3
    Type Detail: Synchronous
    Speed: 1333 MHz
    Manufacturer: Samsung
```
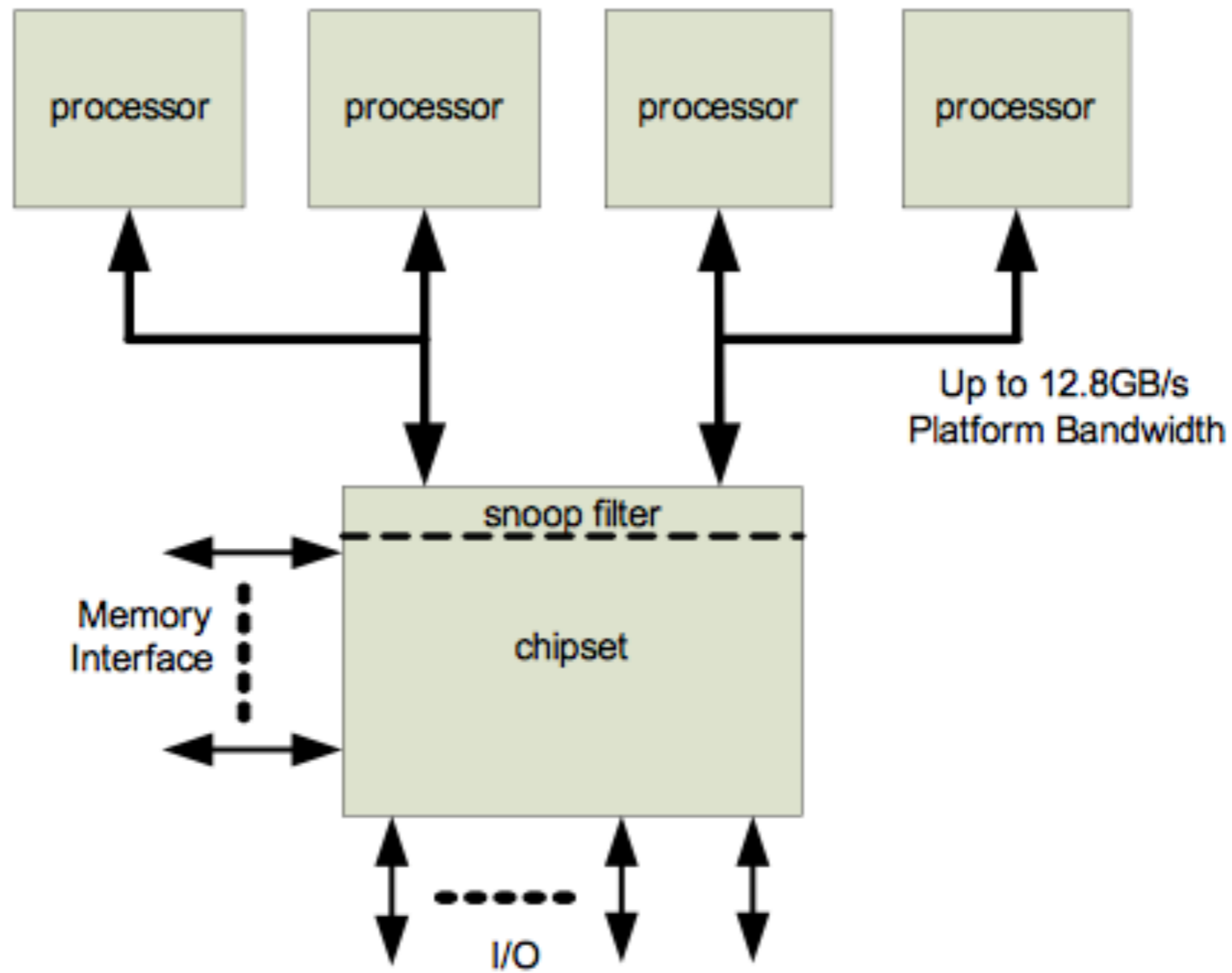
enkitec

# History: UMA

- Uniform Memory Access

- SMP in the 90s.
- Intel bus architecture: FSB.
  - Pentium Pro & Pentium II
  - Northbridge (memory controller hub)
  - Southbridge (I/O controller hub)
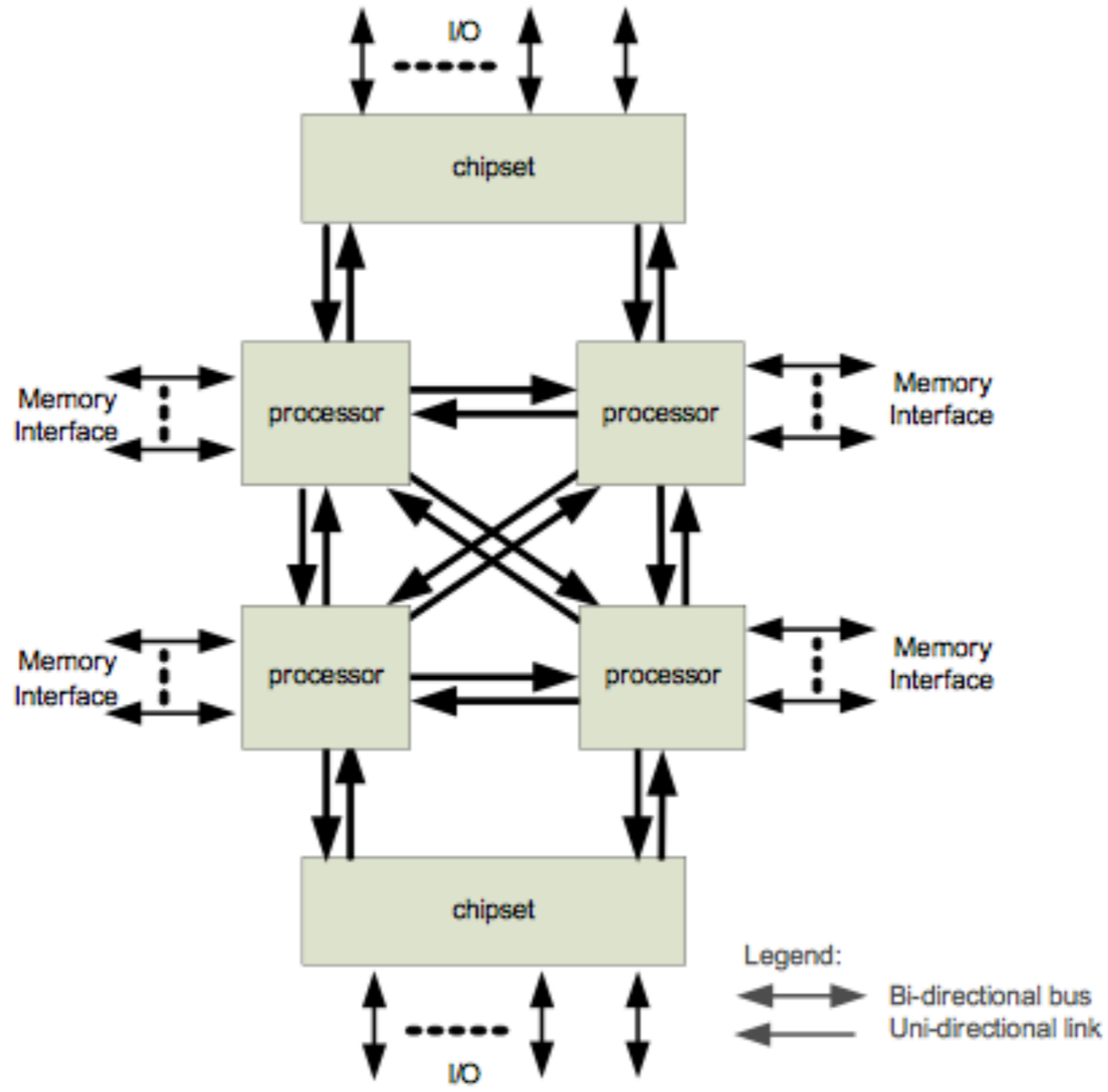  - Architecture provided limited scalability.
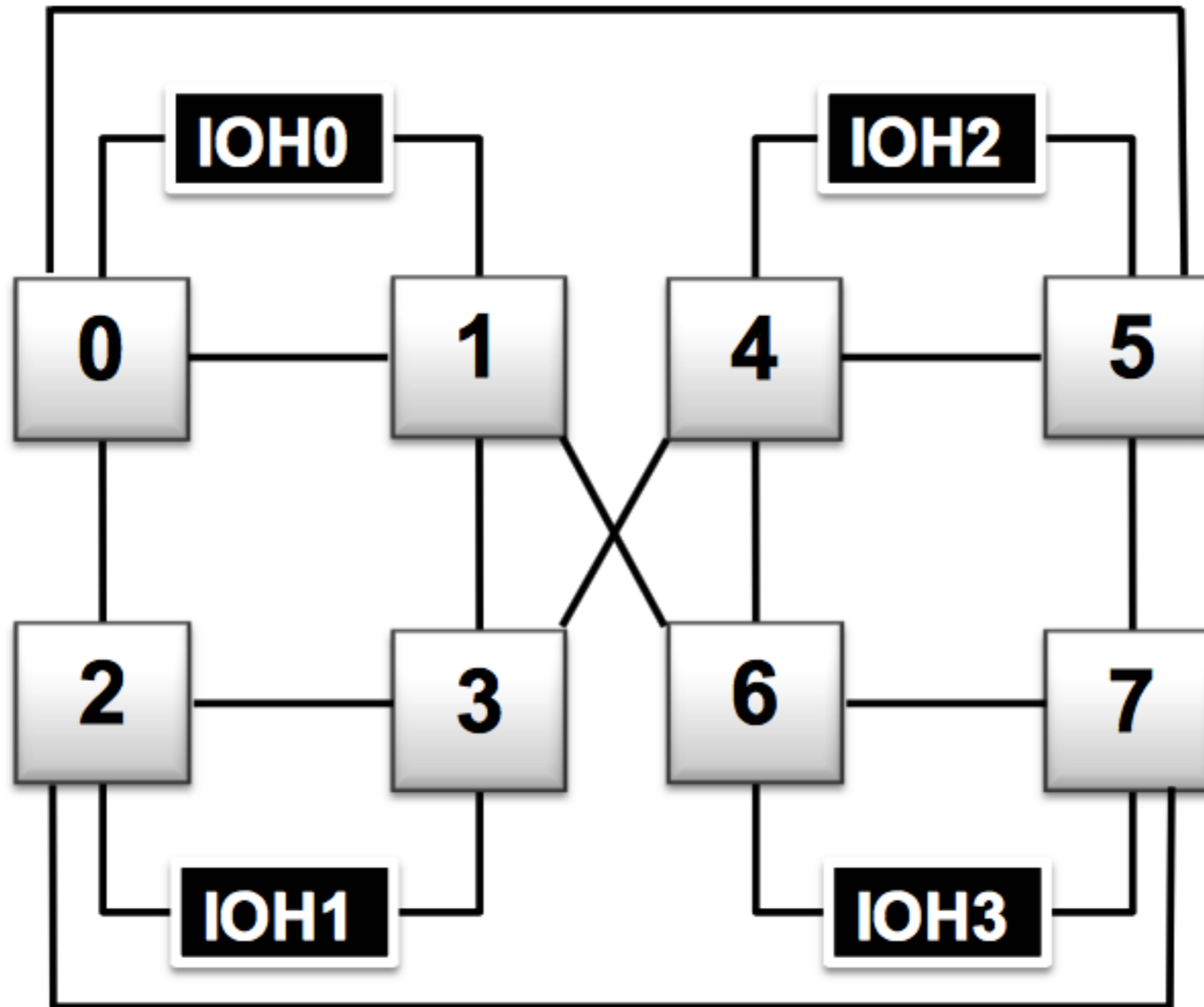
enkitec

# History: UMA

# NUMA

- Non Uniform Memory Access

- Memory local to Socket.
  - Allowing much more memory in a server.
- Each socket can also have its own IO channel.
  - Allowing higher IO rates.
- Sockets interconnected using QPI.
  - For Intel based system starting from Nehalem.

enkitec
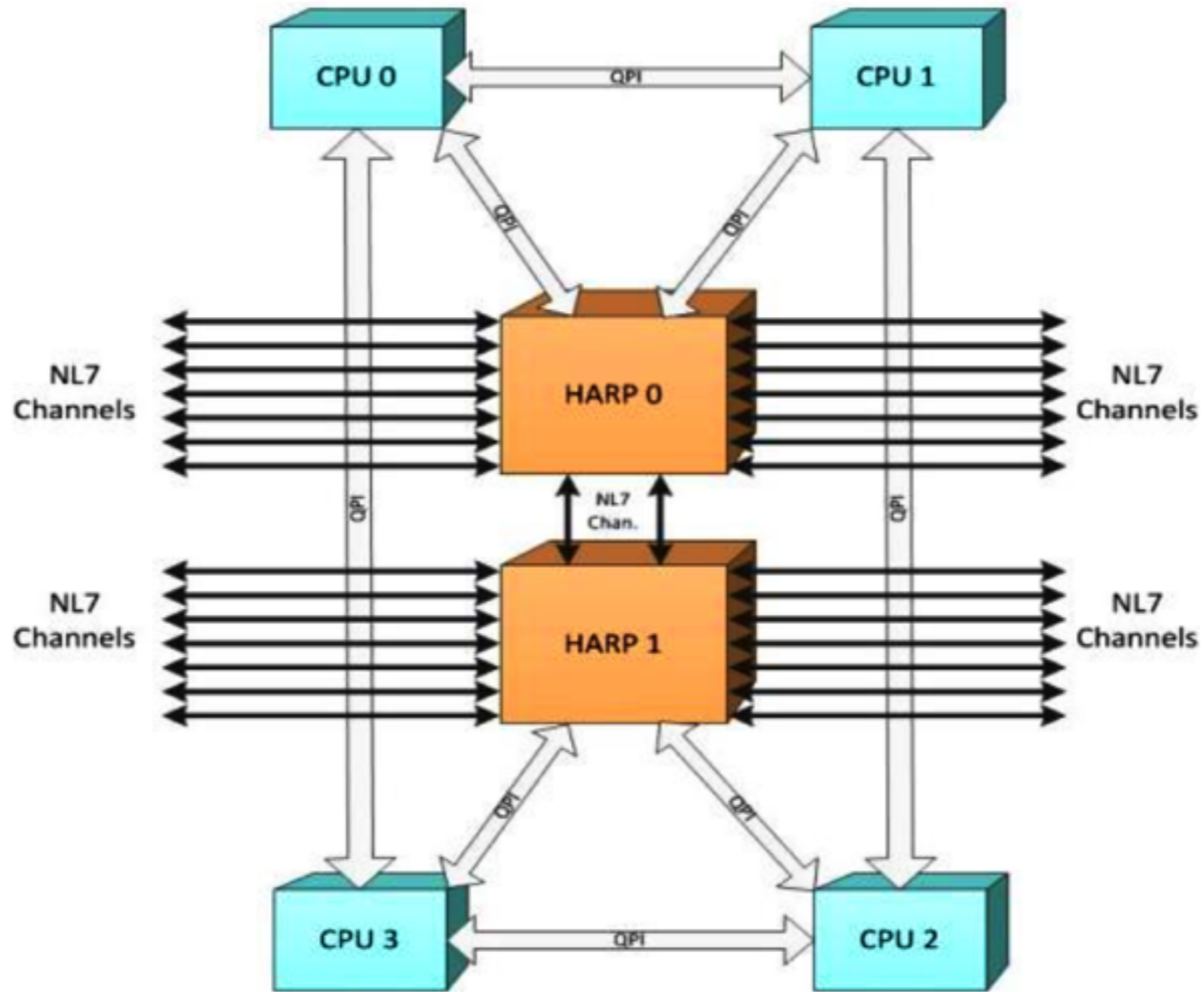
# NUMA

# NUMA - scaling up beyond 4 s.

# NUMA - SGI UV 300

- So how can the UV300 have 32 sockets?

- The sockets are grouped by 4.
- And include two "HARPs": CPU interconnects.
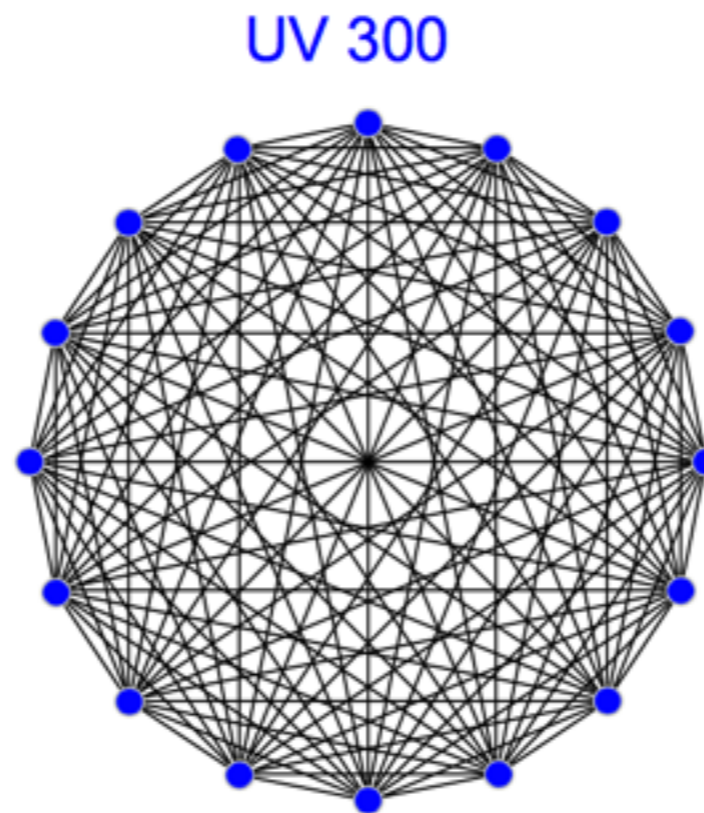- HARPs use SGIs NumaLink7 interconnect.

enkitec

# NUMA - SGI UV 300



CPU and HARP Connections
within the UV 300 Chassis

# NUMA - SGI UV 300

- The HARPs are connected to every other HARP.
- This means a socket is local, 1 or 2 hops away.

UV 300



enkitec

# Performance

- Use 'numactl --hardware' to learn the 'distance'

```
# numactl --hardware
```

**Next socket**

**Local sockets**

node distances:

| node | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | … |
|------|----|----|----|----|----|----|----|----|----|----|----|----|----|---|
| 0: | 10 | 16 | 19 | 16 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | … |
| 1: | 16 | 10 | 16 | 19 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | … |
| 2: | 19 | 16 | 10 | 16 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | |
| 3: | 16 | 19 | 16 | 10 | 50 | 50 | 50 | | | | | | | |
| 4: | 50 | 50 | 50 | 50 | 10 | 16 | 19 | | | | | | | |
| 5: | 50 | 50 | 50 | 50 | 16 | 10 | 16 | | | | | | | |

**Same distance for numalink remote sockets.**

**1 hop**

enkitec

15

# Systems performance

- How does such a system behave?

- Test memory read performance with SLB!
- http://kevinclosson.net/2010/11/17/reintroducing-slb-the-silly-little-benchmark/

- Reads anonymous memory into CPU register

# Systems performance

- What SLB does (snippets of memhammer.c)

make read go into cpu register

```
#define A_RANDOM_LINE  ((mybuffer *) (base_addr + pick_buff()))->buffer[(ops %
LINES_PER_BUFFER) * INTS_PER_LINE ]
volatile register char tmp;

  memset ((void *) base_addr, (int) 0, ((sizeof (mybuffer)) * num_buffers));

  for (i = 0, cur_buffer = base_addr; i < num_buffers; i++, cur_buffer++) {
      for (j = 0; j < INTS_PER_BUFFER; j++)  {
          cur_buffer->buffer[j] = (int) i;
      }
  }

  for (ops=0LL ; ops < 3000000000LL; ops+=1LL) {
      tmp = (char) A_RANDOM_LINE;
  }
```

randomizer
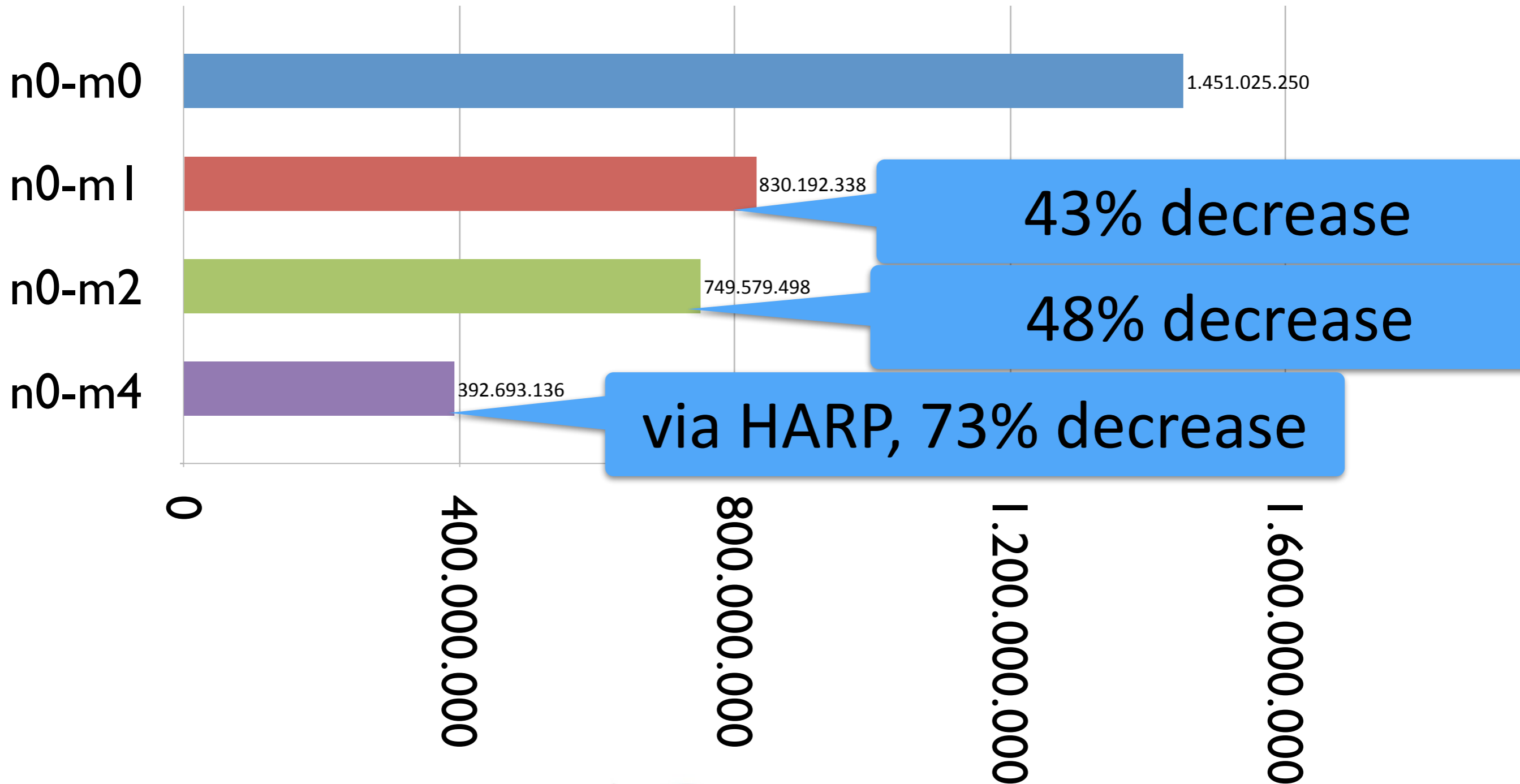
allocate memory and initialise it

read data from main memory

enkitec

# Systems performance

- Why describe this detailed?

- This has nothing to do with Oracle!
- Pure flow from main memory to CPU, 8 bits.

- Reveals memory latency very well!

# SLB - 1 reader throughput



n0-m0    1.451.025.250

n0-m1    830.192.338

43% decrease

n0-m2    749.579.498

48% decrease

n0-m4    392.693.136

via HARP, 73% decrease

0    400.000.000    800.000.000    1.200.000.000    1.600.000.000
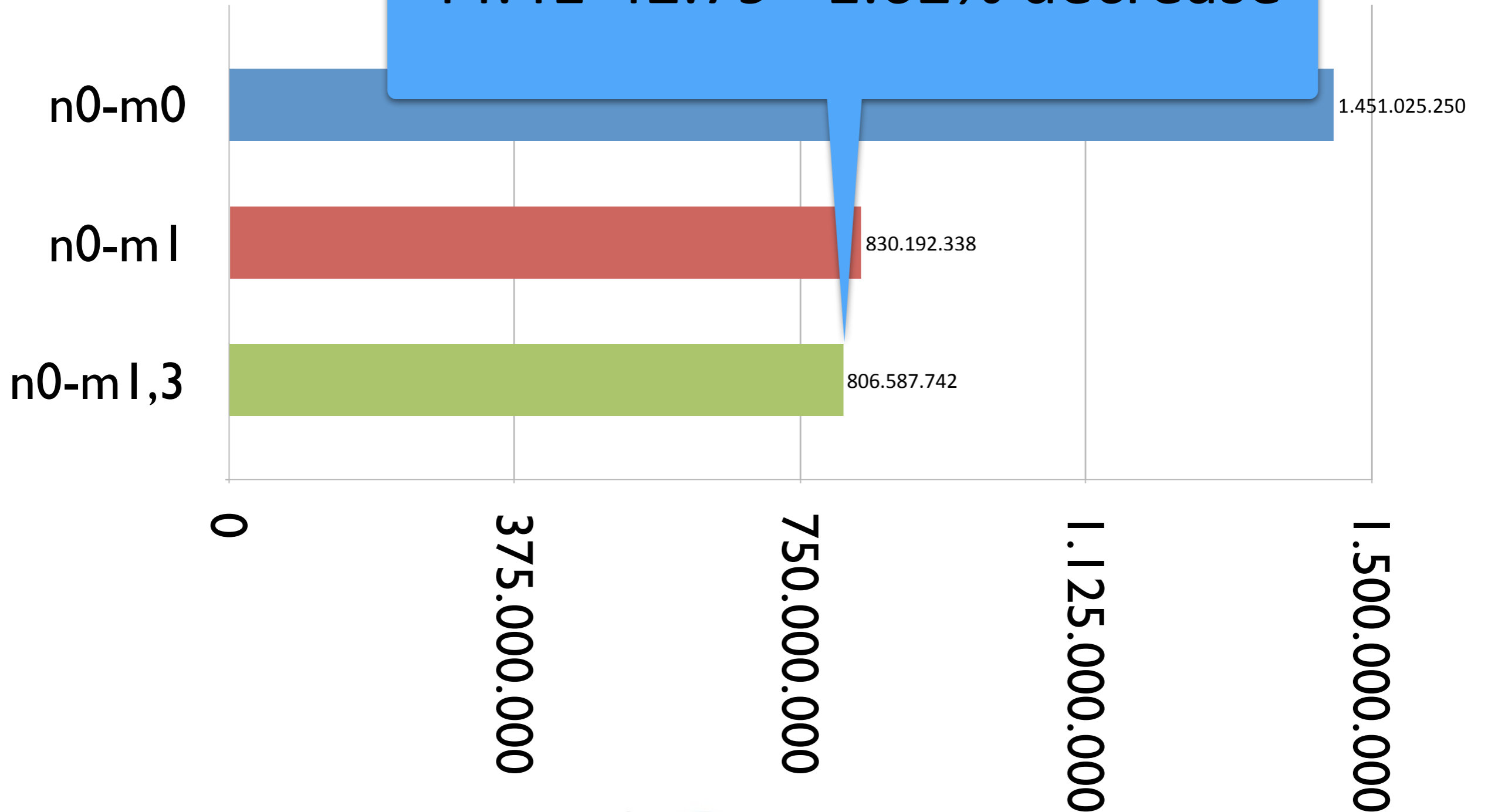
enkitec

19

# Systems performance

- Is this a problem/flaw?

- No: fact of life.
  - Further away resources means increase in latency.

# SLB - Multiple node memory

44.41-42.79= 1.62% decrease

| | |
|---|---|
| n0-m0 | 1.451.025.250 |
| n0-m1 | 830.192.338 |
| n0-m1,3 | 806.587.742 |

0    375.000.000    750.000.000    1.125.000.000    1.500.000.000

enkitec
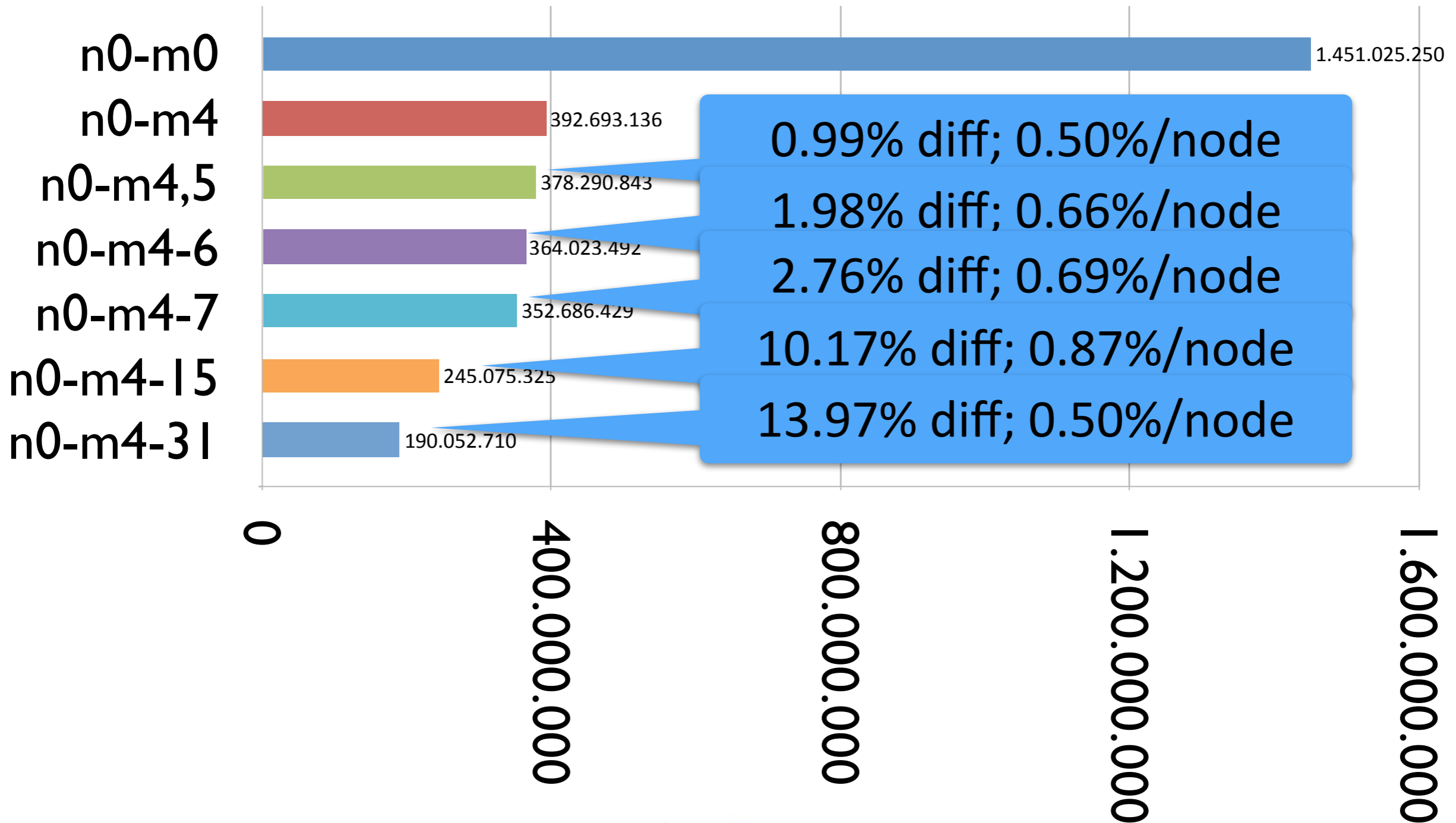
# Systems performance

- By increasing the number of sockets/nodes
  - Latency increases.
  - There's only two nodes at distance 16.

- Let's look at accessing memory via the HARP!

# SLB - Multiple node memory



| | |
|---|---|
| n0-m0 | 1.451.025.250 |
| n0-m4 | 392.693.136 |
| n0-m4,5 | 378.290.843 |
| n0-m4-6 | 364.023.492 |
| n0-m4-7 | 352.686.429 |
| n0-m4-15 | 245.075.325 |
| n0-m4-31 | 190.052.710 |

0.99% diff; 0.50%/node
1.98% diff; 0.66%/node
2.76% diff; 0.69%/node
10.17% diff; 0.87%/node
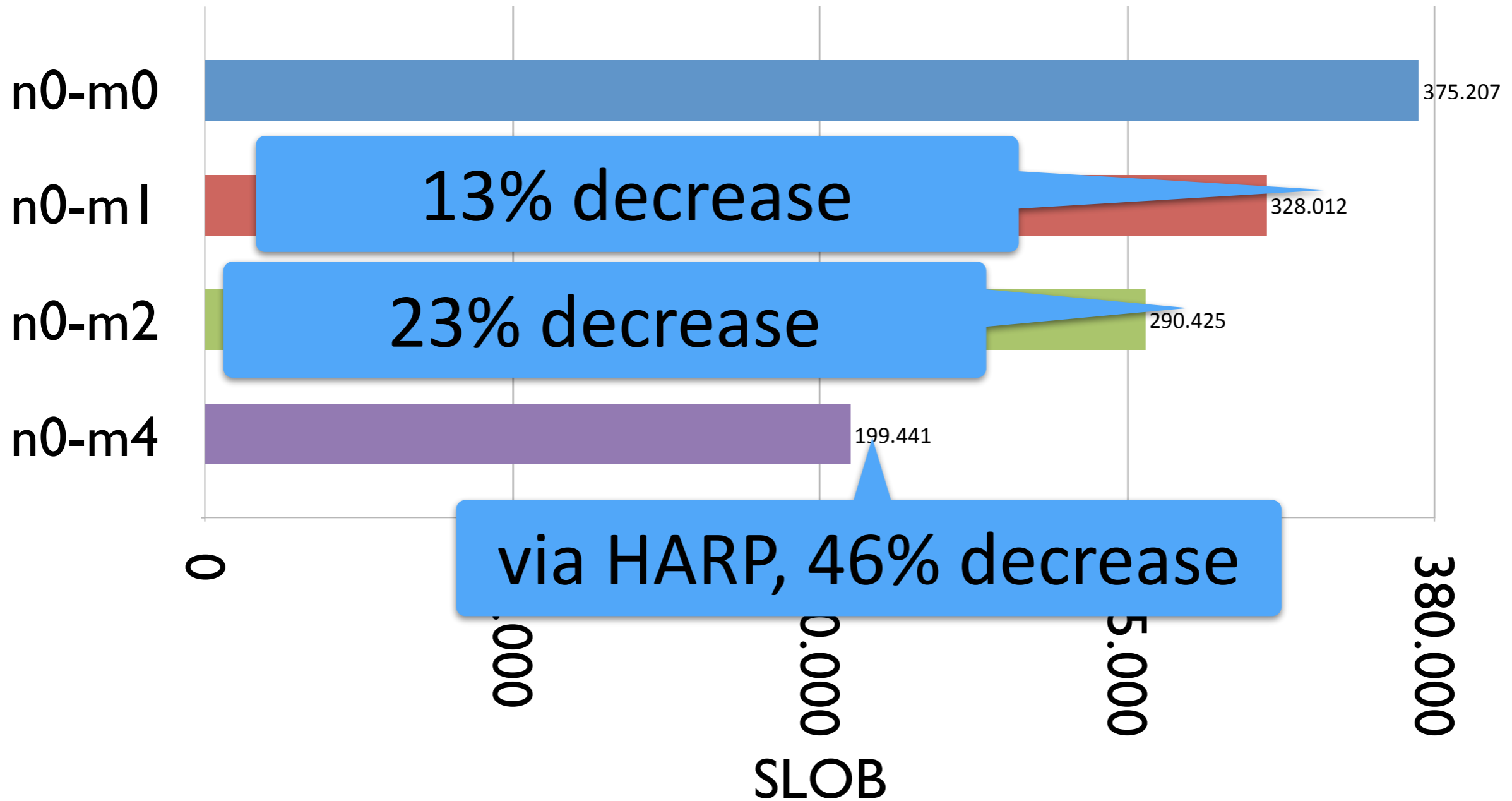13.97% diff; 0.50%/node

enkitec

23

# Systems performance

- Conclusions so far:
  - Accessing remote memory increases latency.
  - Further away memory gives higher latency.
  - Accessing multiple sockets' memory increases latency to:
    - Local sockets: 0.81%/socket (1.62/2)
    - Via HARP: between 0.50-0.87%/socket.

# Oracle performance

- Now let's measure running the Oracle database!

- http://kevinclosson.net/slob

- Test single block access in memory.
- Also known as 'LIO benchmark'.

# SLOB - 1 reader throughput



n0-m0 — 375.207

n0-m1 — **13% decrease** — 328.012

n0-m2 — **23% decrease** — 290.425

n0-m4 — 199.441 — **via HARP, 46% decrease**

0  000  000  000  380.000
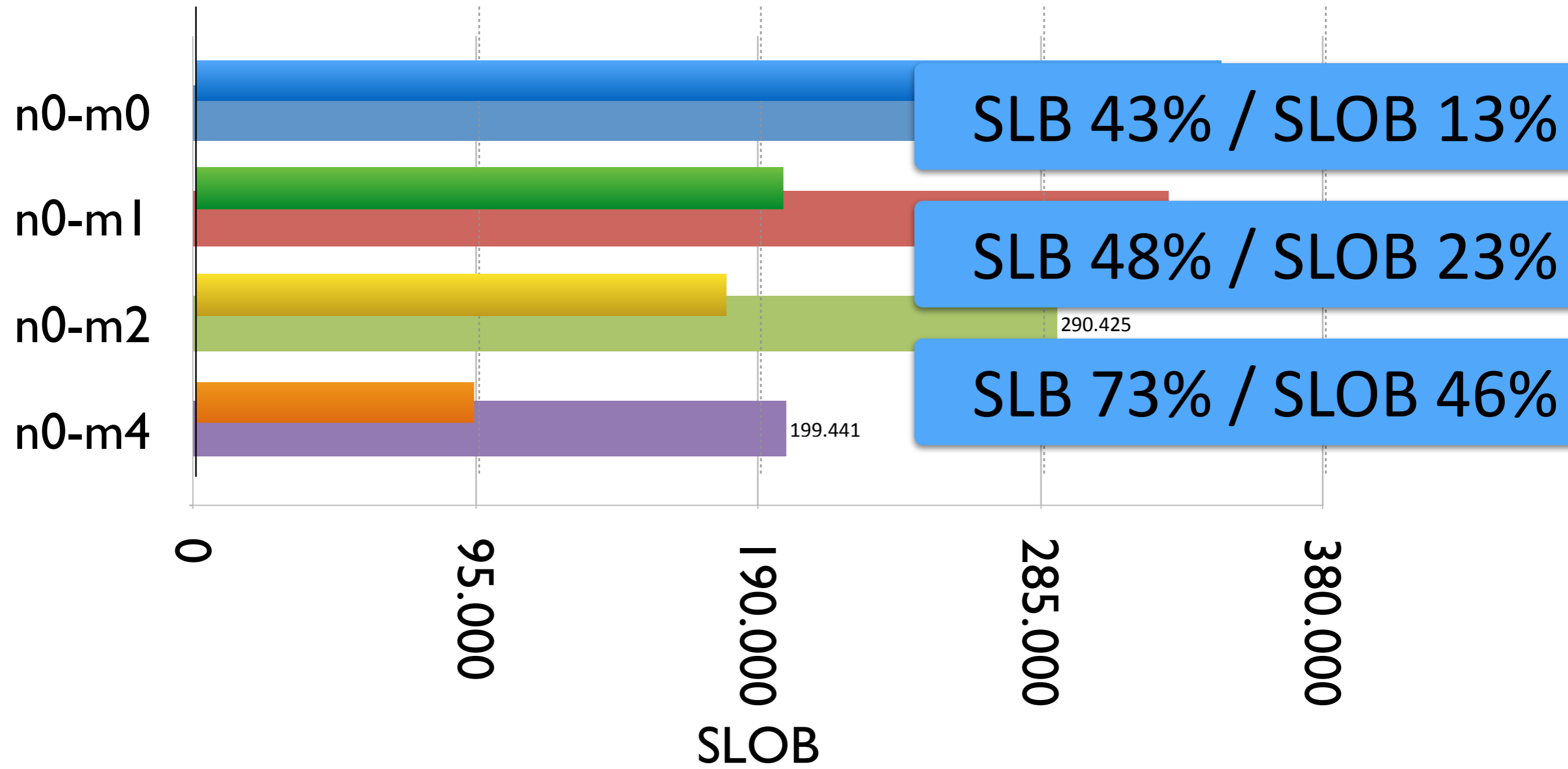
SLOB

# SLOB - 1 reader throughput

- This looks different than the memory benchmark

- Let's overlay the SLOB results with the SLB results.

enkitec

# SLOB - 1 reader throughput

SLB 43% / SLOB 13%

SLB 48% / SLOB 23%

290.425

SLB 73% / SLOB 46%

199.441

**n0-m0**

**n0-m1**

**n0-m2**

**n0-m4**

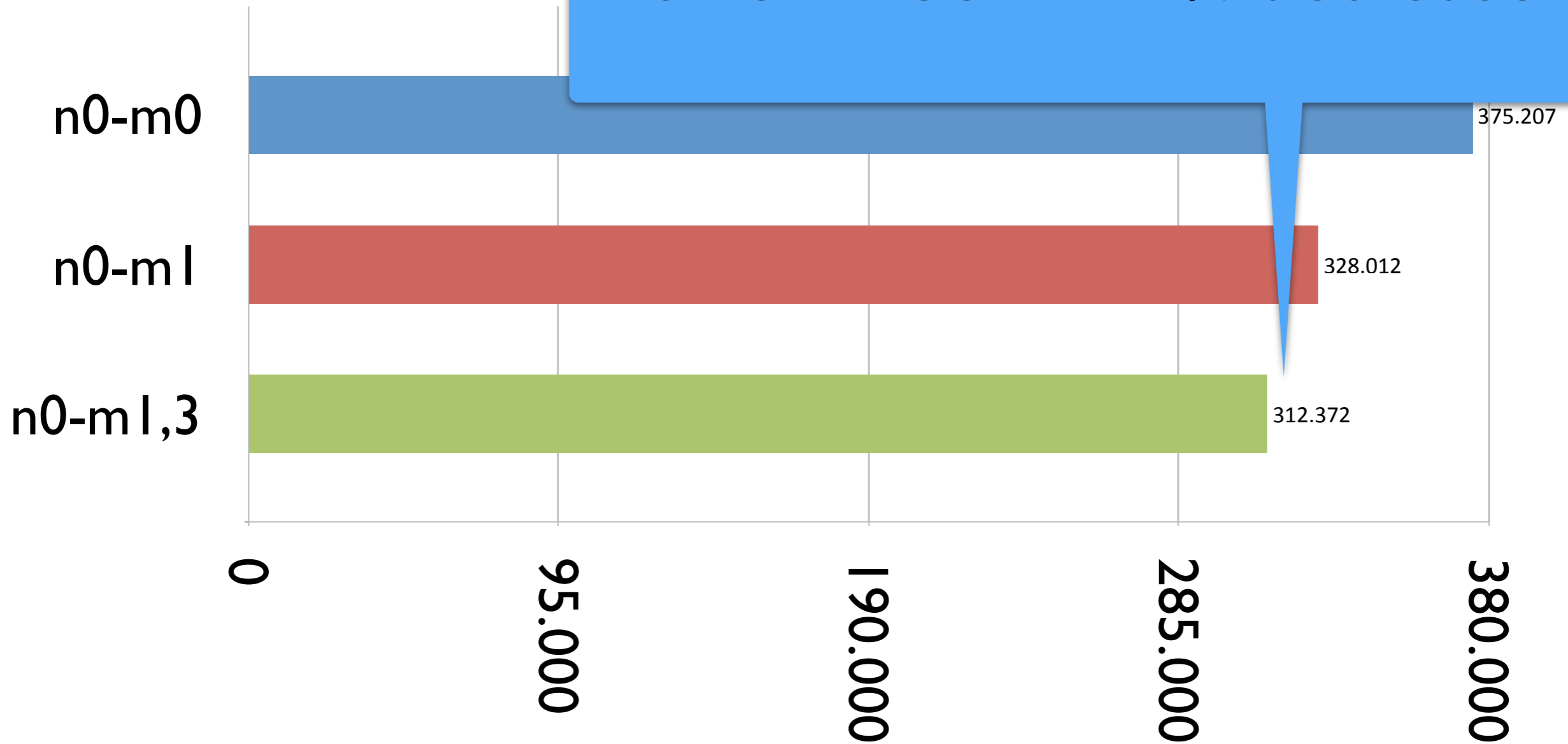0     95.000     190.000     285.000     380.000

SLOB

enkitec

# SLOB - 1 reader throughput

- The Oracle throughput does NOT decline as fast as the SLB one.
  - (this specific) Oracle load is not only accessing memory.
  - This probably means:
    - It is doing processing using L1/2/3 caches!
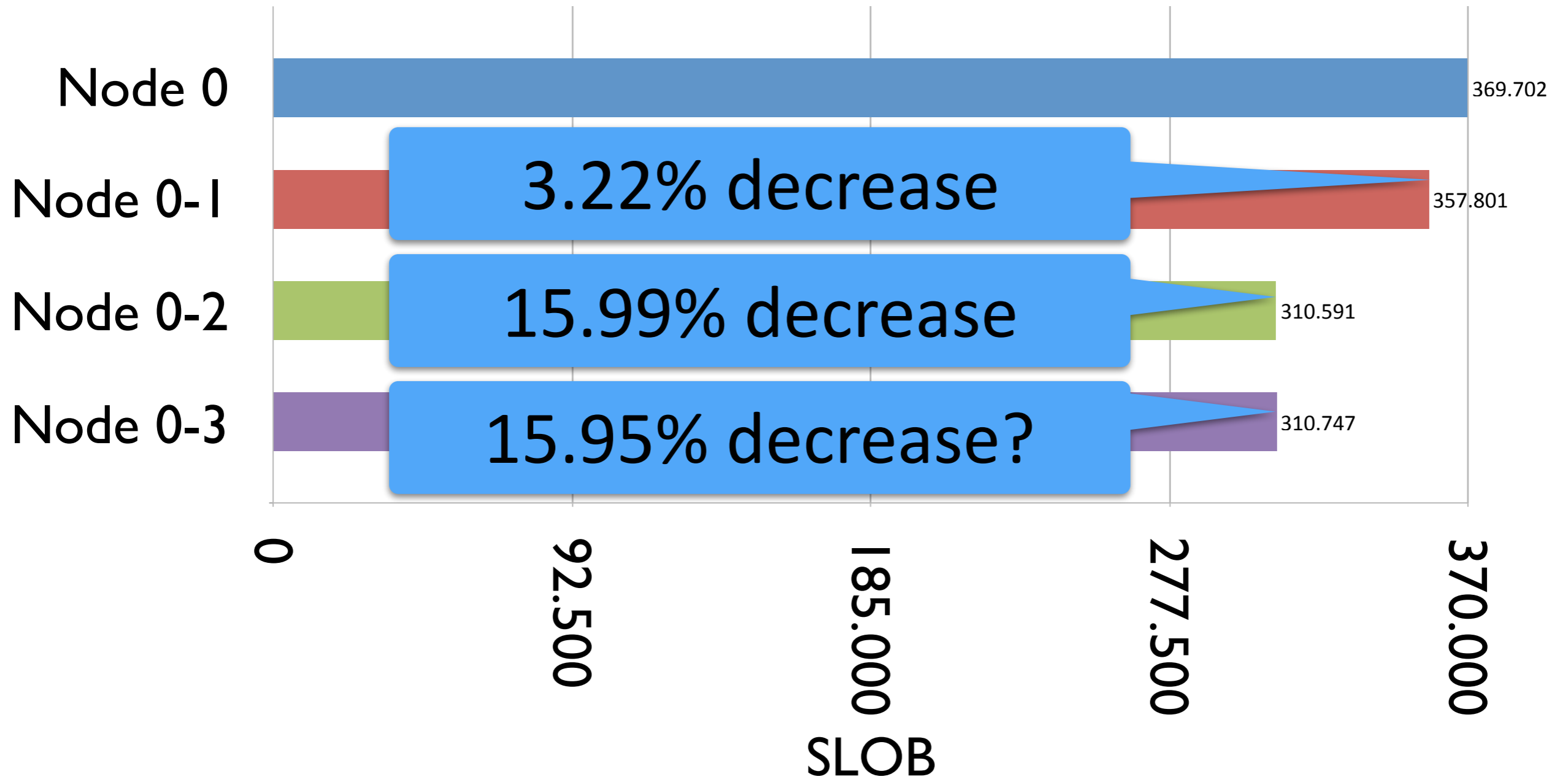    - Probably a result of many years of tuning.

enkitec

# SLOB - Multiple node memory

16.75-12.58= 4.17% decrease

| Category | Value |
|----------|-------|
| n0-m0 | 375.207 |
| n0-m1 | 328.012 |
| n0-m1,3 | 312.372 |

Axis: 0, 95.000, 190.000, 285.000, 380.000

enkitec

# SLOB - Local NUMA nodes

- Previous measurements are done to measure memory latencies.

- These **do not** show real life usage.

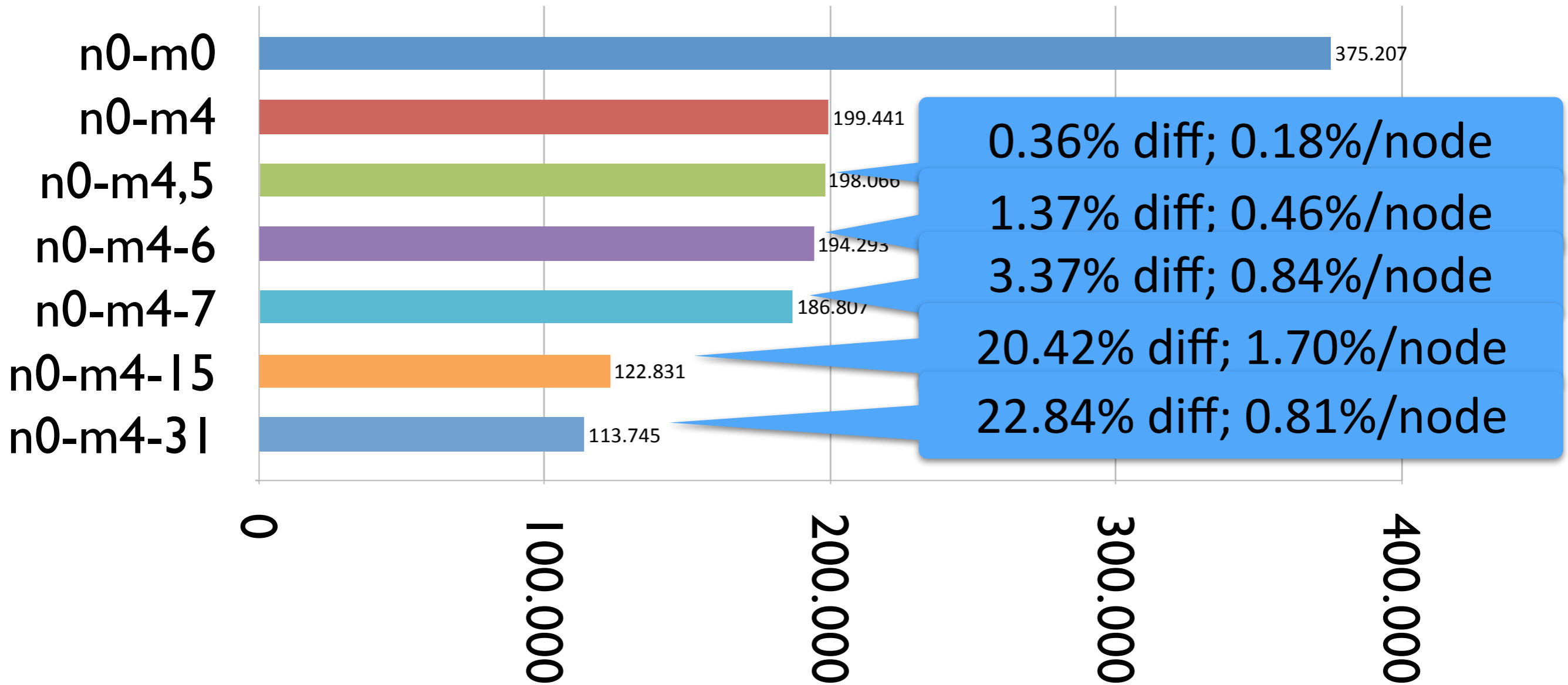- The next slide is an overview of running on one to four NUMA nodes.

enkitec

# SLOB - 1 to 4 nodes / 1 reader



| | |
|---|---|
| Node 0 | 369.702 |
| Node 0-1 | **3.22% decrease** 357.801 |
| Node 0-2 | **15.99% decrease** 310.591 |
| Node 0-3 | **15.95% decrease?** 310.747 |

SLOB

0    92.500    185.000    277.500    370.000

# SLOB - Local NUMA nodes

- Conclusion for Oracle for up to 4 sockets:
  - Latency increases moderately.
    - Way less than pure memory access.
  - For two node servers, don't enable NUMA.
  - For up to four nodes.
    - Milage varies. Probably not worth the effort.
    - Test your own load.
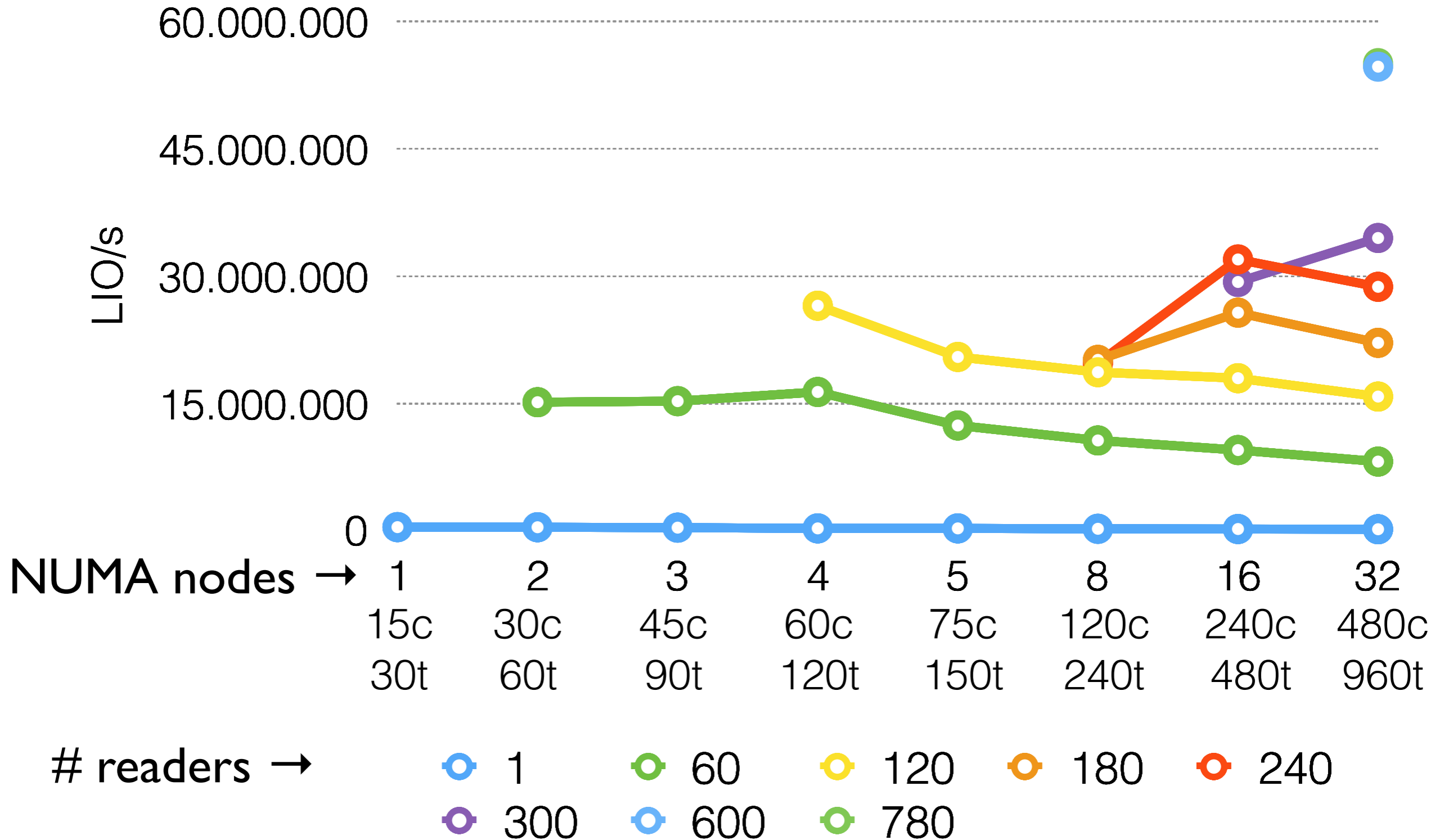
enkitec

# SLOB - Multiple node with HARP



| Label | Value |
|---|---|
| n0-m0 | 375.207 |
| n0-m4 | 199.441 |
| n0-m4,5 | 198.066 |
| n0-m4-6 | 194.293 |
| n0-m4-7 | 186.807 |
| n0-m4-15 | 122.831 |
| n0-m4-31 | 113.745 |

0.36% diff; 0.18%/node

1.37% diff; 0.46%/node

3.37% diff; 0.84%/node

20.42% diff; 1.70%/node

22.84% diff; 0.81%/node

enkitec

# Oracle/SLOB performance

- Conclusions so far:
  - Accessing multiple sockets' memory increases latency to:
    - Local sockets: 2.09%/socket (4.17/2)
    - Via HARP: between 0.18-1.70%/socket.
      - Reason for the difference Oracle <> SLB:
        - Probably L1/2/3 cache influence.

# Oracle performance

- Now let's focus on bandwidth with LIO:
  - Start independent sessions without affinity.
  - Using SLOB.
  - Every reader reads its own schema.
  - Index range scans.

- Run SLOB until PIO vanishes from AWR.
- Then measure SLOB run.
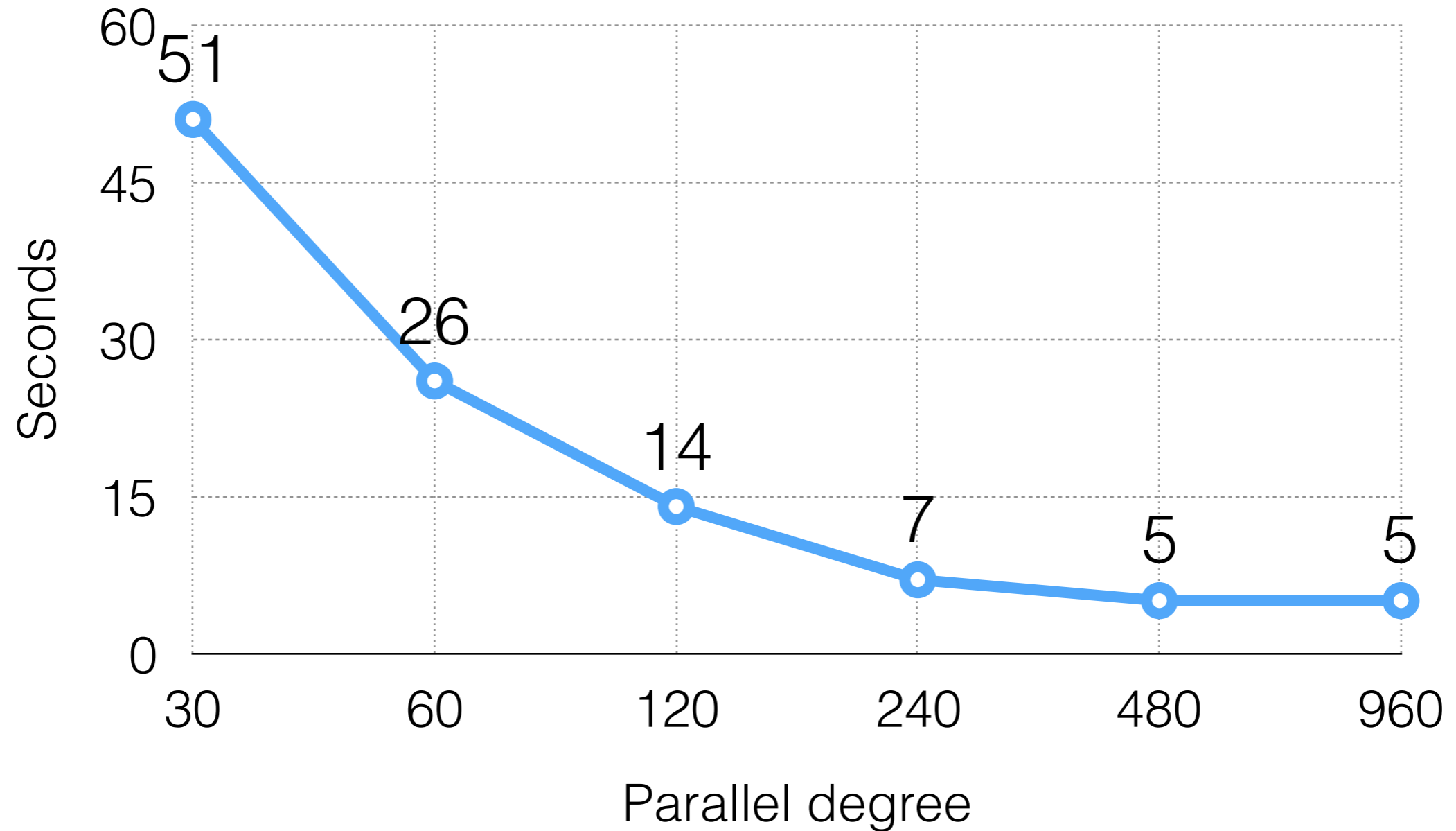
SLOB readers throughput

enkitec

37

# Scan a big table with PQ

- Created a set of tables with the TPCH kit.
  - Table H_LINEITEM is the biggest one.
  - Size: 739G / 96'864'152 blocks.
  - 5'999'989'709 rows. 6 billion rows!

  - Created with SCALE=1000

# Scan a big table with PQ

- Set my buffercache to be 10T.
- Scanned table with in memory PQ option
  - alter session set parallel_degree_policy=auto;

- Normal scanning only read 1/3rd in the buffer cache.

- Set my KEEP pool to 1T.
- Altered the table H_LINEITEM to the KEEP pool
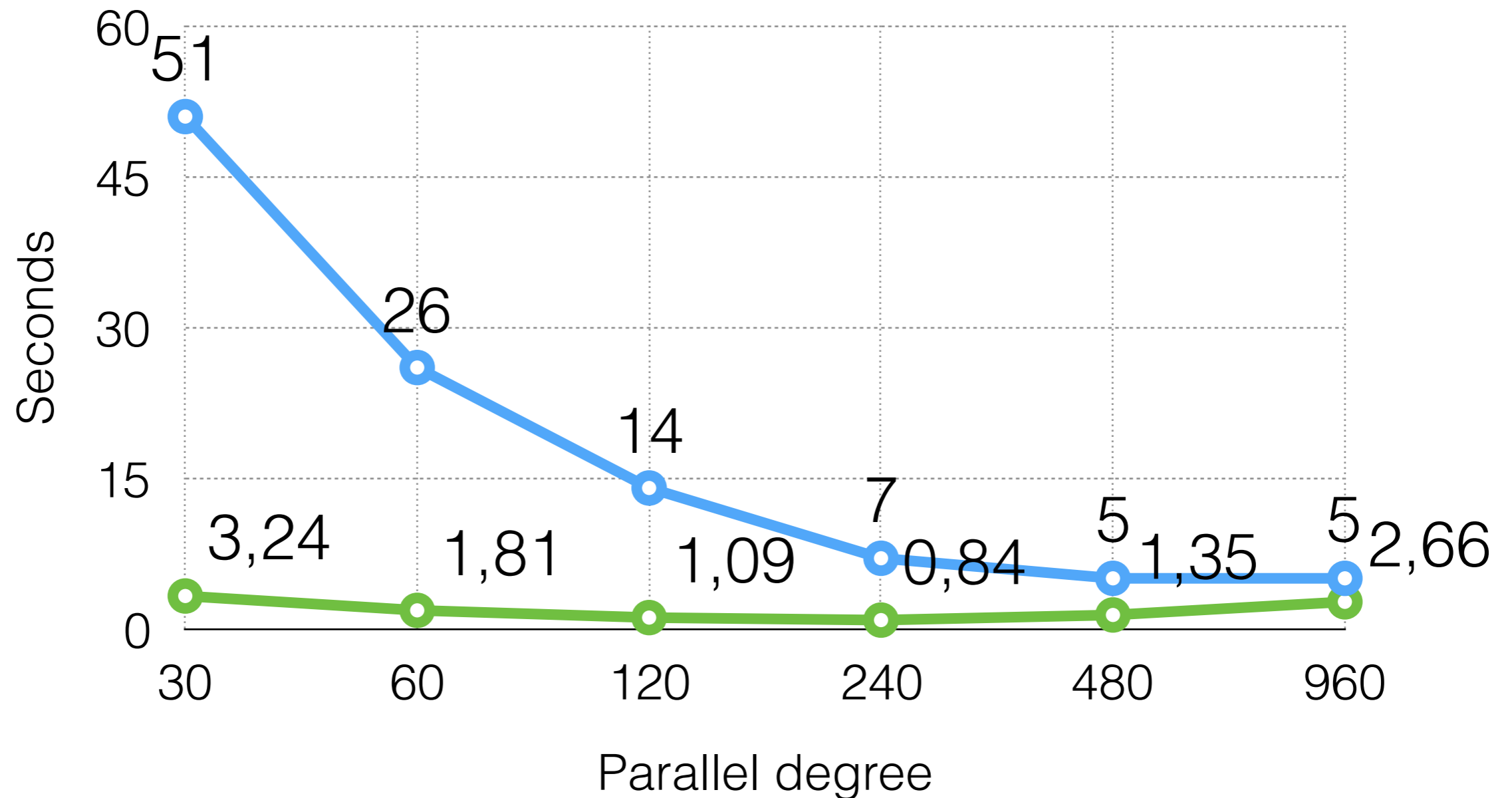
# Scan a big table with PQ

# Scan a big table with PQ

- Near linear scaling up to core count.
  - All slaves busy, no producer/consumer model.

- Let's see what the in-memory option can do!
  - Added 1T in-memory pool.
  - Copied table for in-memory (query high).
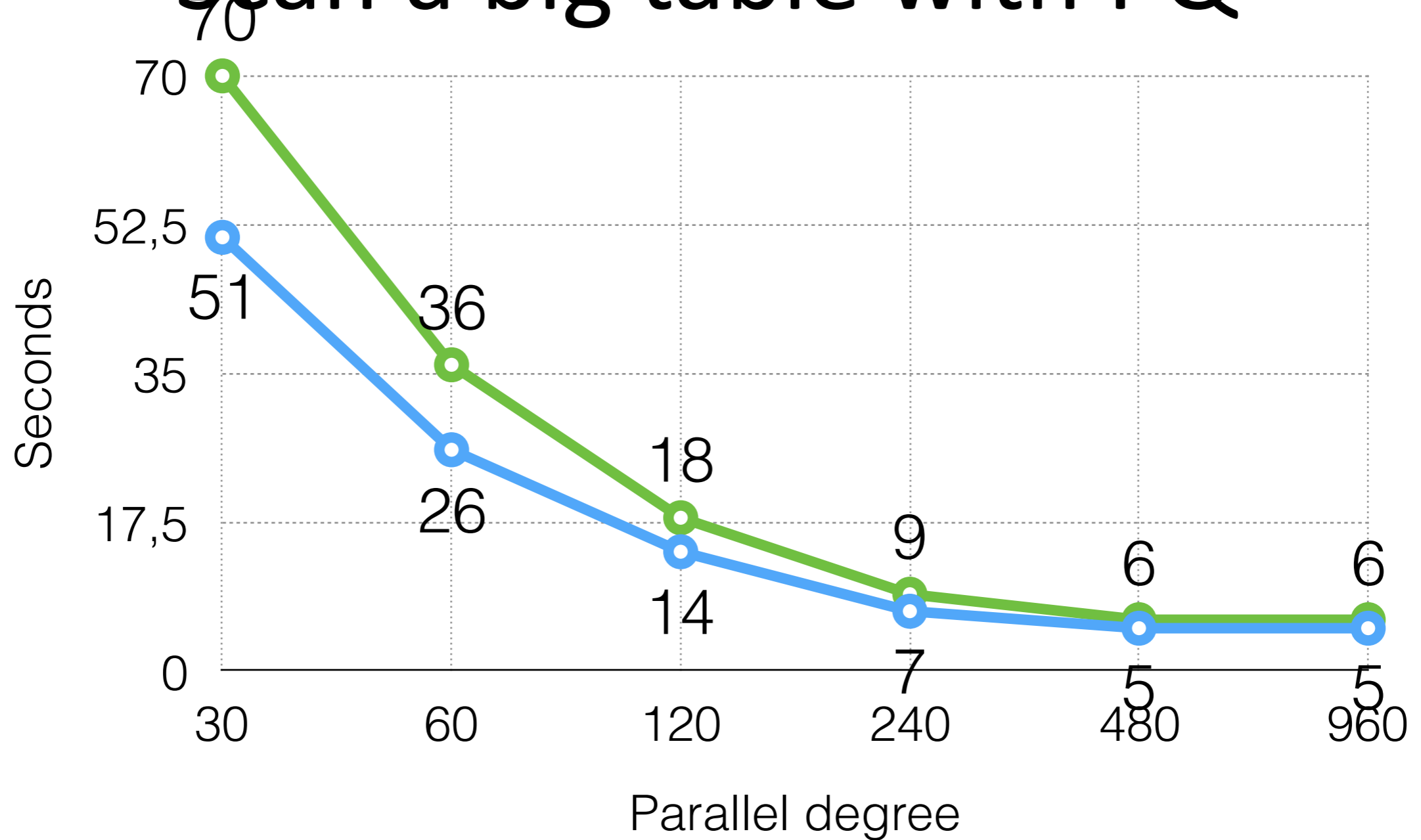
# Scan a big table with PQ

# Scan a big table with PQ

- Unbelievable performance with in-memory option.

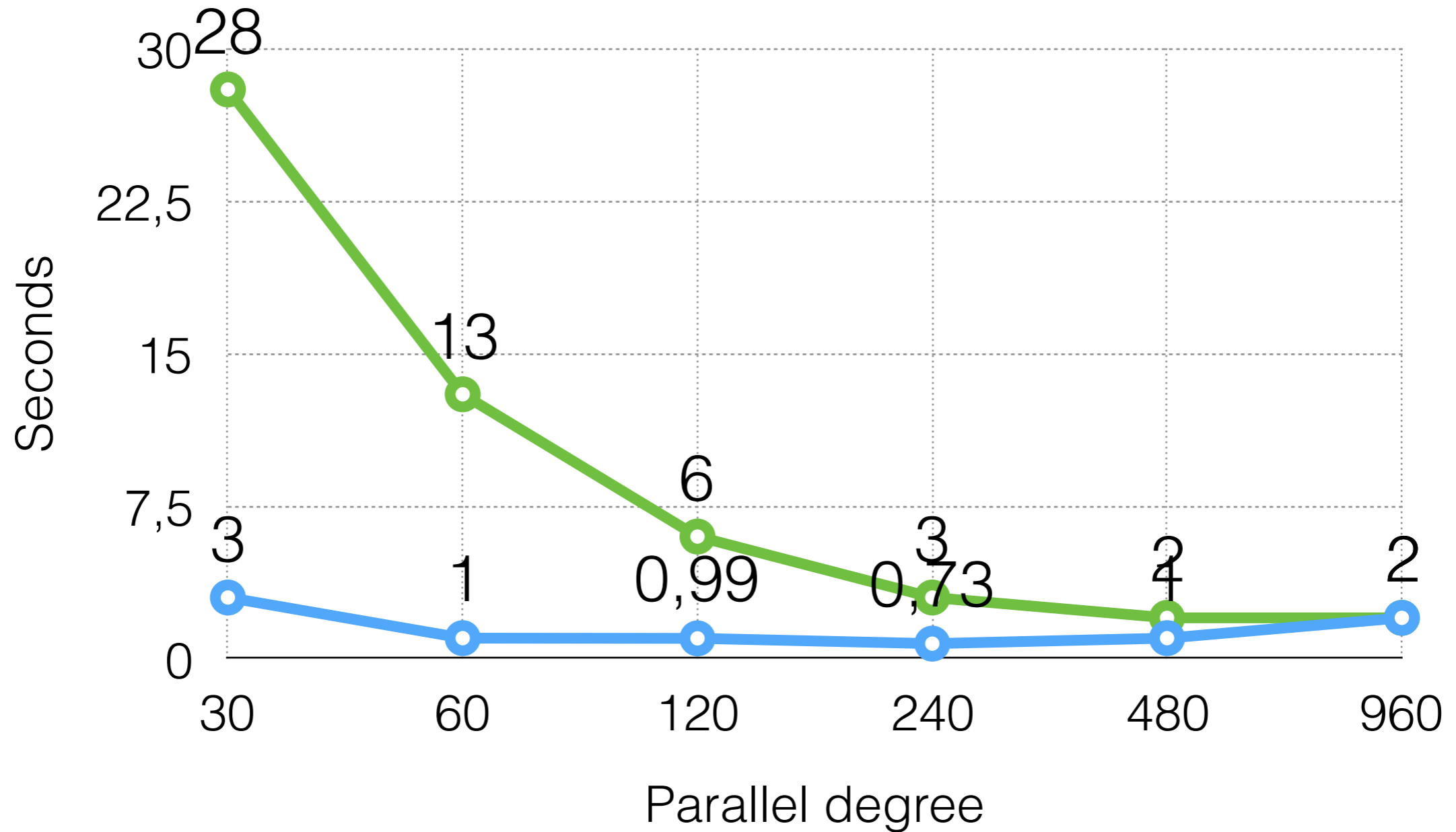- No such thing as magic.

# Scan a big table with PQ

- in-memory:
  - Suspicion of HCC count(*) optimization.

- To mitigate potential pre-computed result:
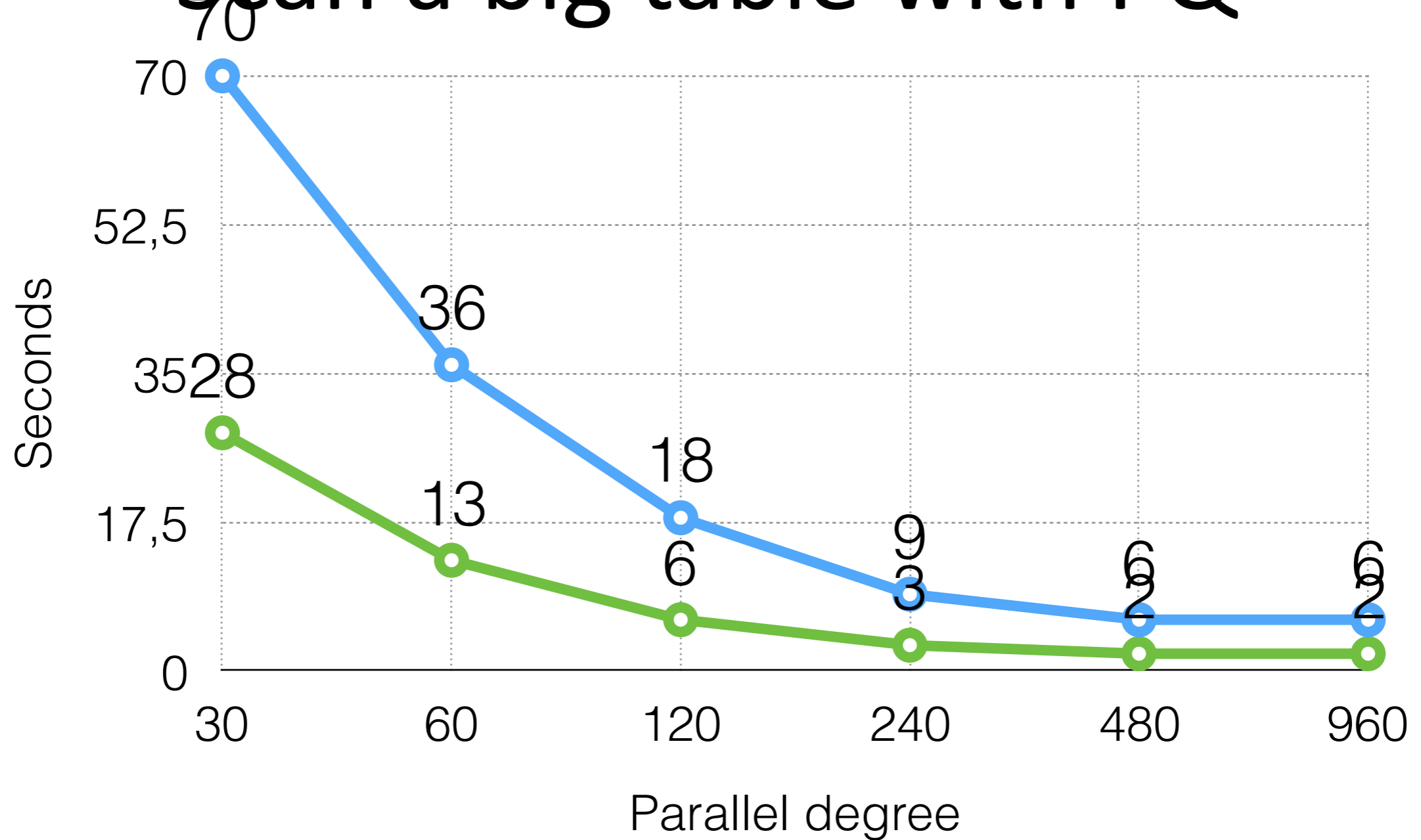  - Try different function: avg()

# Scan a big table with PQ



- ○ count(*) KEEP cache
- ○ avg(l_orderkey) KEEP cache

45

# Scan a big table with PQ



count(*) in-memory
avg(l_orderkey) in-memory

# Scan a big table with PQ



Chart: "Scan a big table with PQ"

Y-axis: Seconds (0, 17,5, 35, 52,5, 70)
X-axis: Parallel degree (30, 60, 120, 240, 480, 960)

avg(l_orderkey) KEEP cache (blue): 70, 36, 18, 9, 6, 6
avg(l_orderkey) in-memory (green): 28, 13, 6, 3, 2, 2

Legend:
- avg(l_orderkey) KEEP cache
- avg(l_orderkey) in-memory

enkitec

# Scan a big table with PQ

- in-memory option reduces 66% of time(!!)
- No vector (SIMD) processing tested.

# Conclusion

- Increase in NUMA nodes increases random memory access latency.

- However, core count increases processing capacity too.

- Memory placement is important with NUMA count > 4 - 8.

enkitec

# Conclusion

- UV 300 has constant distance via HARP.
  - This means adding NUMA nodes scales linearly.

- Oracle OLTP processing takes very efficient usage of on-die caches (L1/2/3).

- Disable NUMA on low socket count (>=4) servers for Oracle, unless you can prove it benefits.

enkitec

# Conclusion

- PQ can be turned to cached reads by:
  - Setting the NOCACHE attribute to CACHE*.
  - Moving a table in the KEEP pool.
- With no cache fixation, Oracle might restrict blocks in cache to 1/3rd of the total.

- The in-memory parallel query option scans almost all blocks into cache.

enkitec

# Conclusion

- UV 300 processing works well together with PQ.

- In-memory compression has pre-computed count(*) optimization.

- There is overhead involved in PQ processing.

enkitec

# Oracle Findings

- The SGA huge pages are initialised by a single process.
  - Initialising a 10T SGA takes a significant time.
- Shutdown normal/immediate never finishes.
  - PMON failed to acquire latch error.
  - Process shutting down the instance continuously running through /proc/stat.

# Oracle Findings

- Heavy parallel scan on 750G table in memory.
  - Took ~ 5 seconds right after startup.
  - Took ~ 20 seconds after 14 days uptime.
  - 5 seconds was restored after bounce.
  - Uncertain what the cause is…

enkitec

# Q & A