# NUMA-aware Optimization on Apache Spark*

Xiaochang Wu

Intel Big Data Team

# Legal Disclaimer

# Optimization Notice

Software and workloads used in performance tests may have been optimized for performance only on Intel® microprocessors. Performance tests, such as SYSmark* and MobileMark*, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.

For more information go to http://www.intel.com/performance.

# Agenda

- Introduction to NUMA

- Apache Spark* patch to enable NUMA support

- NUMA tuning for Apache Spark* applications

- Experiments on various workloads

- Conclusion

# What is NUMA?

- <u>N</u>on <u>U</u>niform <u>M</u>emory <u>A</u>ccess

- Address the problem of performance hit when several processors attempt to address the same memory.

- Multiple physical CPUs in a system. Each CPU has
  - Local memory: memory attached to it, fast.

  - Remote memory: memory attached not directly, slower.

- Developed commercially during the 1990s. Intel announced NUMA compatibility for its x86 servers in late 2007 with its Nehalem CPU.

- Trend: bring memory nearer to processor

(intel)
Software

# Typical 2-Sockets Intel® Xeon® NUMA Topology

Node 0

Node 1

Memory ↔ Processor ↔ Processor ↔ Memory

Interconnect

NUMA performance considerations:

• Latency: higher latency of accessing remote memory

• Bandwidth: interconnect contention

Avoid remote memory accesses!

$ numactl --hardware

available: 2 nodes (0-1)

node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65

node 0 size: 98207 MB

node 0 free: 92847 MB

node 1 cpus: 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87

node 1 size: 98304 MB

node 1 free: 94228 MB

node distances:

node   0   1

 0:  10  21

 1:  21  10

# NUMA on Linux*

- Linux* kernel's default NUMA policy allocates pages from local node of the processor which makes the request, and fall back to other nodes if no local memory available.

- When a task (run by thread) is scheduled to a node, all previously allocated pages in other nodes will not be migrated to this node, from then on, the accessing to those pages will be remote and cost much more time.

- Solution: Manual binding using numactl command

```
$ numactl --cpunodebind=$numa_node_bind --preferred=$numa_node_bind
<application command line>
```

# NUMA memcpy test on Linux*

- memcpy test of numademo* tool
  - First allocate specified size of memory based on NUMA policy and copy the first half data to the rest half

- Default policy allocates pages on the node of the CPU that triggers the allocation.

- Set cpu-bind to node 0 and preferred memory-bind to node 0. Compare the performance between "alloc on node 1" vs. "local allocation" to simulate the worst case differences.

| Data Size | 4GB | 8GB | 16GB | 32GB | 64GB | 128GB |
|-----------|-----|-----|------|------|------|-------|
| Scaling | 1.086x | 1.074x | 1.124x | 1.116x | 1.134x | 1.079x |

"local allocation" vs "alloc on node 1"

# Apache Spark Deployment Revisited

**Master Node**

Master

Driver

**Worker Node**

Worker

Executors

**Worker Node**

Worker

Executor 0:
Pin to
NUMA Node 0

Executor 1:
Pin to
NUMA Node 1

Binding Executors:
- Executors only run on binding CPU node
- Memory for executors allocated on binding CPU node preferably

Issues:
- Sometimes binding node may be too busy
- Chances to allocate memory on remote node
- Can't predict workloads on application-level

No Silver Bullet!

# Apache Spark Patch to enable prefix command and NUMA binding

- Patch: https://github.com/apache/spark/pull/16411

- This patch will support adding a prefix command to the original executor launch command line. Eg:

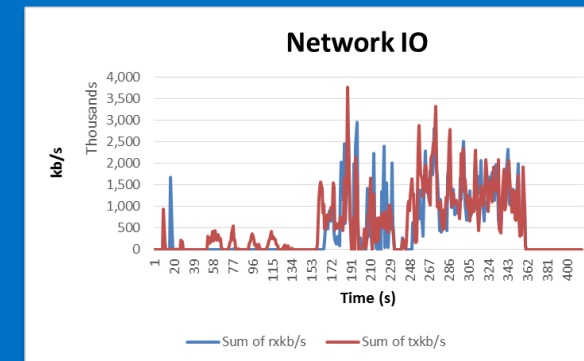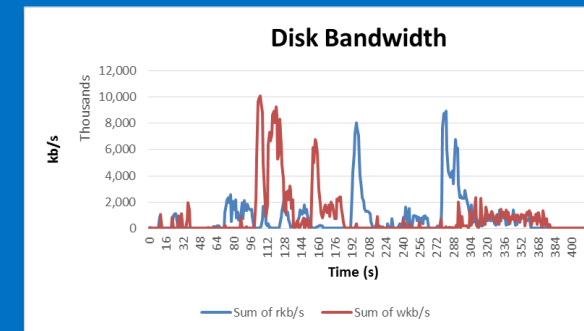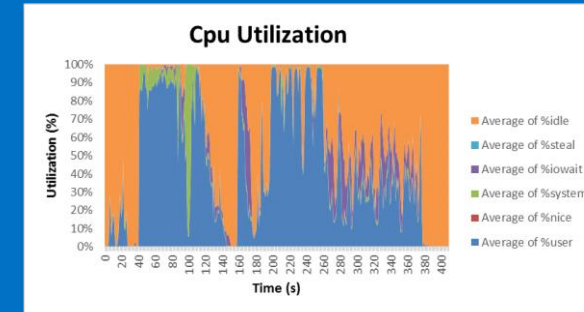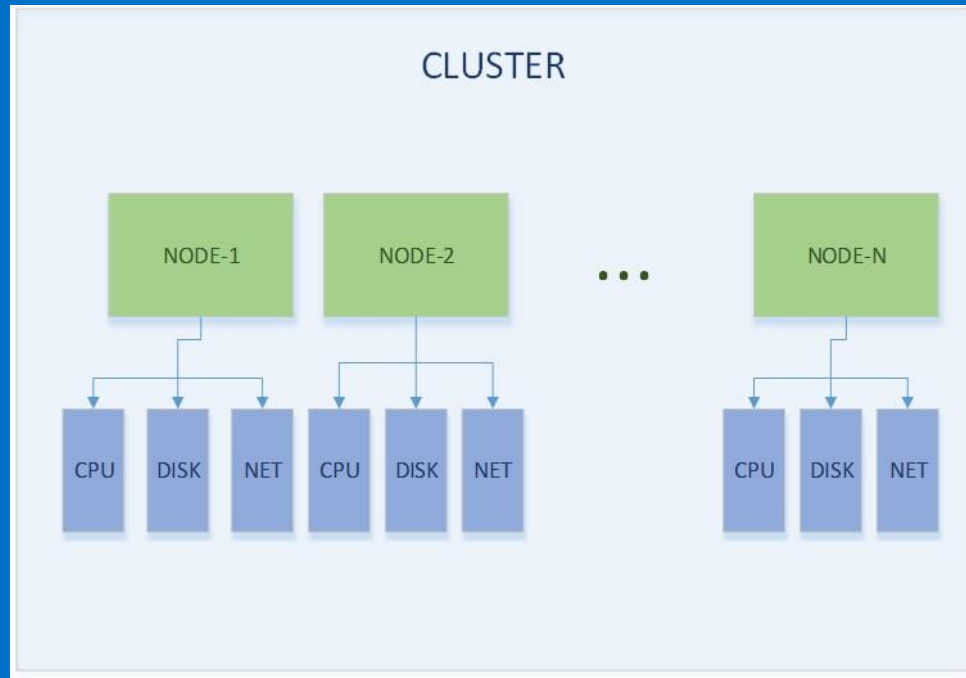| Prefix Command | Original Executor Command Line |
|---|---|
| *spark-numa.sh* | {{JAVA_HOME}}/bin/java -server -Xmx4096m -Djava.io.tmpdir={{PWD}}/tmp '-Dspark.driver.port=49187' --driver-url spark://CoarseGrainedScheduler@10.0.2.192:49187 --executor-id 26 --hostname sr593 |

- This prefix command (*spark-numa.sh* in this case) can be a user customized script which invokes *numactl* to bind each executor to a specific node according to executor id.

- In our experiments, a simple interleave binding is applied to balance task load, users need to fine tune script for more complex situation.

# NUMA Tuning Workflow

- Prepare Workloads

- Benchmark baseline results with no NUMA support

    - Use Apache Spark* Web UI to check how Spark application executes and Use PAT check CPU, Memory, Disk IO and Network IO usages and tune Spark parameters for optimization.

    - Run each workload several times until the result doesn't deviate much from previous ones. The first several runs will warm up cache.

    - Calculate average of last 3 results.

- Checking real-time NUMA metrics with numatop*

- Use Intel® VTune™ for Linux to do memory access analysis. Check remote memory access (RMA) and local memory access (LMA) ratios.

- If RMA/LMA ratio is high, apply NUMA-support patch and benchmark to see if there is any performance gain.

# Intel® PAT – Intel® Performance Analysis Tool



- Image from : https://github.com/intel-hadoop/PAT

# Checking real-time NUMA metrics with numatop

- numatop: alternative tool to check real-time NUMA stats
- NUMA access stats for TPCx-BB* q01 sample



NUMA-Unaware

Applied NUMA-aware patch

# VTune™ Memory Access Analysis Example

- Identify memory access issues with Intel® Vtune™

- Collect data with Linux command line and analyze with Windows GUI

- Capture 5 secs data on worker node:

```
$ amplxe-cl –collect memory-access --duration 5
```

| NUMA-unaware VTune™ data | | VTune™ data after applying NUMA-aware patch | |
|---|---|---|---|
| Elapsed Time: | 5.006s | Elapsed Time: | 5.007s |
| CPU Time: | 349.969s | CPU Time: | 374.452s |
| Memory Bound: | 28.5% | Memory Bound: | 25.3% |
| L1 Bound: | 14.6% | L1 Bound: | 15.6% |
| L2 Bound: | 1.3% | L2 Bound: | 1.5% |
| L3 Bound: | 6.0% | L3 Bound: | 5.3% |
| DRAM Bound: | 7.9% | DRAM Bound: | 4.1% |
| Memory Bandwidth: | 11.4% | Memory Bandwidth: | 12.7% |
| Memory Latency: | 59.3% | Memory Latency: | 61.2% |
| Remote / Local DRAM Ratio: | 0.675 | Remote / Local DRAM Ratio: | 0.000 |
| Loads: | 275,790,673,472 | Loads: | 347,496,824,592 |
| Stores: | 91,926,178,872 | Stores: | 109,500,842,488 |
| LLC Miss Count: | 208,806,264 | LLC Miss Count: | 144,804,344 |
| Average Latency (cycles): | 9 | Average Latency (cycles): | 9 |
| Total Thread Count: | 3,497 | Total Thread Count: | 3,660 |
| Paused Time: | 0s | Paused Time: | 0s |

Memory access VTune™ analysis for TPCx-BB q01 sample

# Cluster Configurations for TPC-DS* and TPCx-BB* experiments

| Hardware | |
|---|---|
| CPU | Intel® Xeon® CPU E5-2699 v4 @ 2.20GHz<br>22 cores 44 threads * 2 sockets<br>2 NUMA Nodes |
| Memory | Kinston* DDR4-2133 8G * 24 slots = 192G<br>96G for each NUMA node |
| Storage | For OS:<br>- SSD 1 X Intel® SSD DC S3510 Series（480GB,2.5in SATA 6Gb/s)<br>for Data:<br>- HDD：7 X Seagate* Constellation.2 ST9500620NS 500GB 7200 RPM 64MB Cache SATA 6.0Gb/s |
| Network | Intel® Ethernet Controller 10 Gigabit X540-AT2 |

| Software | |
|---|---|
| OS | CentOS* Linux release 7.1.1503 (Core)<br>Kernel version: 3.10.0-514.2.2.el7.x86_64 |
| Java | OpenJDK* Runtime Environment (build 1.8.0_111-b15) |
| Hadoop | Hadoop* 2.7.3 |
| Spark | Spark* 2.1 |

- 1 master node + 2 worker nodes
- Spark SQL parameters:

--master spark://sr592:7077

--driver-memory 16g

--num-executors 32

--executor-memory 11520m

--executor-cores 5

--conf "spark.sql.shuffle.partitions=480"

# Workloads for Experiments

- ## Query samples from TPC-DS
  - TPC-DS models the decision support functions of a retail product supplier. User queries convert operational facts into business intelligence
  - 1000 GB text database compressed to 277.4 GB parquet format and stored on HDFS

- ## Query samples from TPCx-BB
  - TPCx-BB has analytical queries in the context of retailers with physical and online store presence. The queries are expressed in SQL for structured data and in machine learning algorithms for semi-structured and unstructured data.
  - 1000 GB text database compressed to 198.1 GB ORC format and stored on HDFS

- ## A typical Telco SQL workload
  - SQL queries to summarize customers consumption characteristics utilizing billing data.
  - 10GB text database, compressed to 3GB parquet format and stored on HDFS.

# TPC-DS Results

- ## Query Samples:

  - Query 55: interactive query

  - Query 58: deep reporting

  - Query 73: data mining

- ## Run TPC-DS Queries:

  ```
  $ spark-sql --database $database_name -f $query_name.sql
  ```

- ## Results:

| Sample | NUMA-Unaware(s) | NUMA-Aware(s) | Scaling |
|--------|-----------------|---------------|---------|
| q55    | 24              | 23.6          | 1.017x  |
| q58    | 52              | 52            | 1.000x  |
| q73    | 28              | 27            | 1.037x  |

# TPCx-BB Results

- ## Query Samples:

  - Query 01: UDF/UDTF on structured data

  - Query 05: machine learning on semi-structured data

  - Query 27: UDF/UDTF/NLP on un-structured data

- ## Run TPCx-BB Queries

  ```
  $ ./bin/bigBench runQuery -q $query_number -U
  ```

- ## Results:

| Sample | NUMA-Unaware(s) | NUMA-Aware(s) | Scaling |
|--------|-----------------|---------------|---------|
| q01 | 80 | 74 | 1.081x |
| q05 | 249 | 228.7 | 1.089x |
| q27 | 114 | 105 | 1.086x |

# A Typical Telco SQL Workload Analysis and Results

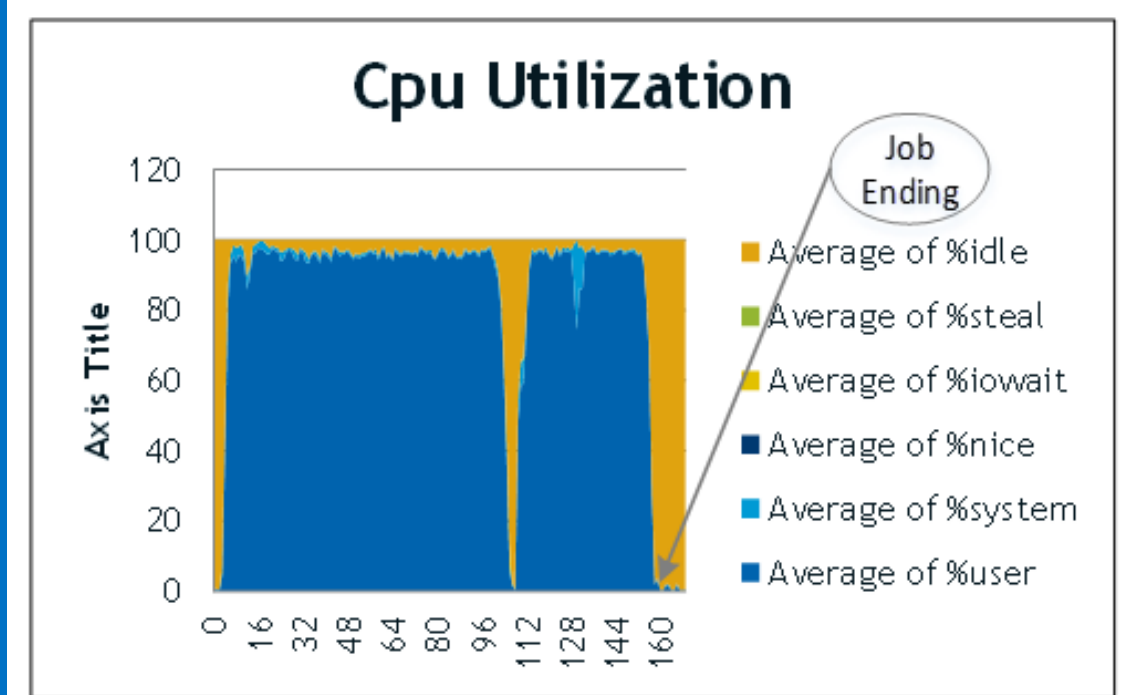


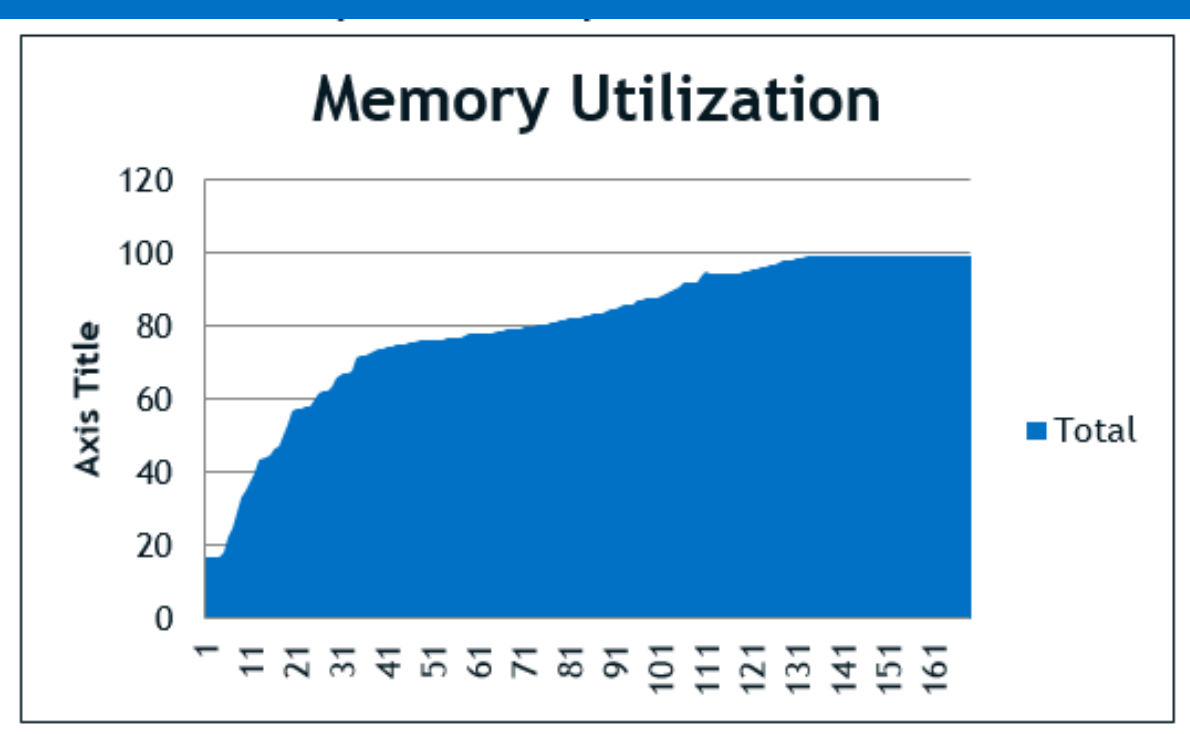


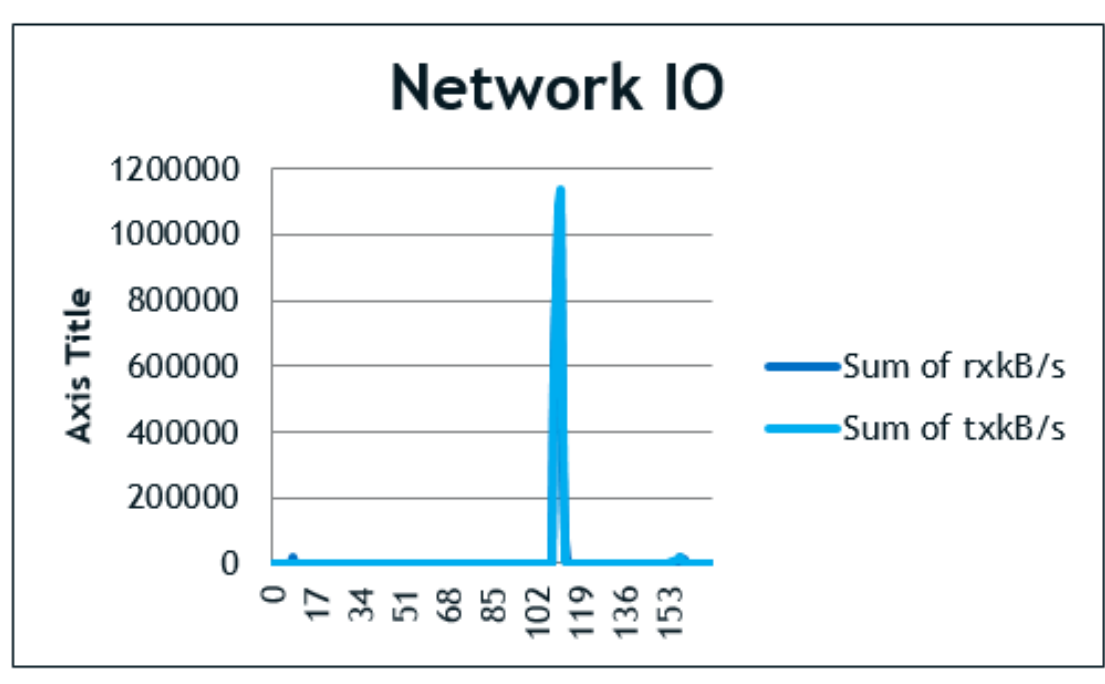


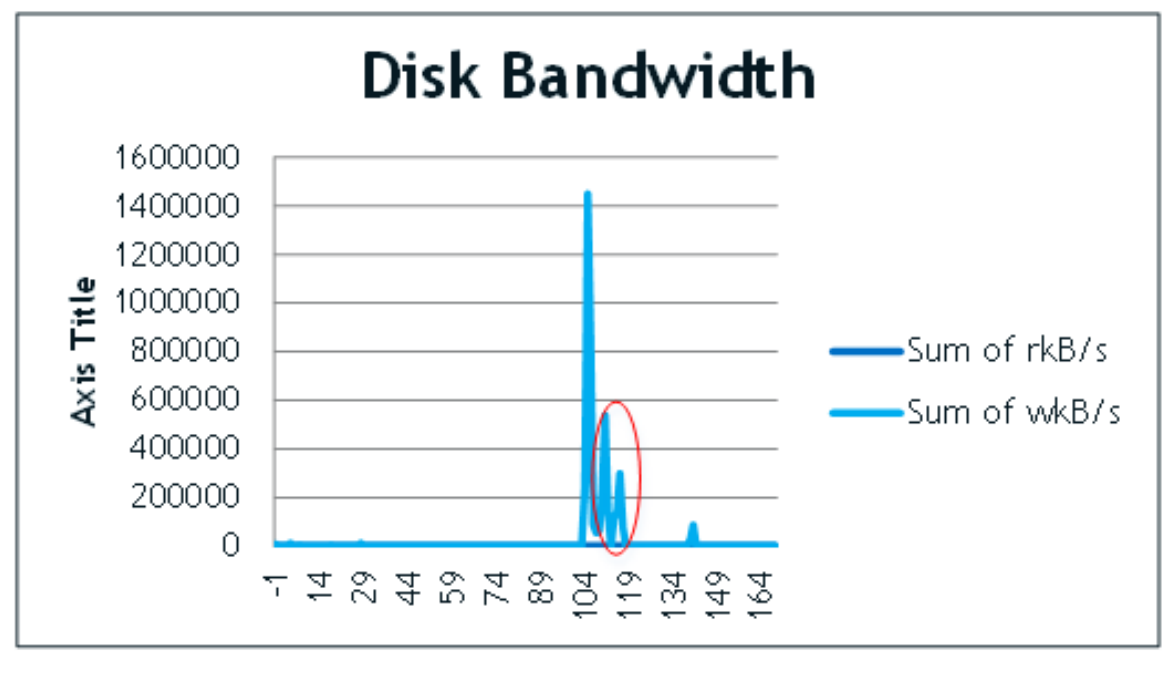


*Be noted: All slave nodes have same behaviors, we present one of them*

| Events Count | CPU_CLK_UNHALTED.THREAD | INST_RETIRED.ANY | MEM_LOAD_UOPS_L3_MISS_RETIRED.LOCAL_DRAM | MEM_LOAD_UOPS_L3_MISS_RETIRED.REMOTE_DRAM | UOPS_RETIRED.RETIRE_SLOTS |
|---|---|---|---|---|---|
| w/ NUMA-Aware Patch | 25,325,367,481,579 | 21,680,659,933,335 | 2,648,340,894 | 58,198,229 | 22,606,108,509,456 |
| w/o NUMA-aware Patch | 26,956,329,603,414 | 21,782,319,467,891 | 1,855,031,856 | 1,220,313,270 | 22,769,700,890,331 |

- Cycles Per Instruction

CPI is a little high.
NUMA-Aware: 1.2, NUMA-Unaware: 1.2

- Pipeline Slot Utilization

Utilization is low, vectorization computing may help.
NUMA-Aware: 22% , NUMA-Unaware: 21%

- L3 Miss

The L3 miss impact is minor.
NUMA-Aware: 2% Overall Cycles, NUMA-Unaware: 2.6% Overall Cycles

- Remote Node Access

Remote Memory Accessing is not the bottleneck.
NUMA-Aware: 0.07% Overall Cycles, NUMA-Unaware: 1.2% Overall Cycles

*Be noted: All slave nodes have same behaviors, we present one of them*

Profiling with PAT

Profiling with VTune™

Result:

| Sample | NUMA-Unaware(s) | NUMA-Aware(s) | Scaling |
|---|---|---|---|
| Telco SQL | 163 | 156 | 1.045x |

# Conclusion / Key Takeaways

- NUMA penalties varies as workloads change

  - Hard to identify memory access patterns

  - Not all big data workloads have NUMA issues, need to analyze case by case

- Leverage platform tools such as Intel® VTune™ / numatop to examine memory access

- Manual CPU binding helps when NUMA issues are big

# References

- NUMA on Linux: http://man7.org/linux/man-pages/man7/numa.7.html

- TPC-DS: http://www.tpc.org/tpcds/default.asp

- TPCx-BB: http://www.tpc.org/tpcx-bb/default.asp

- PAT: https://github.com/intel-hadoop/PAT

- Intel VTune for Linux: https://software.intel.com/en-us/intel-vtune-amplifier-xe/

- numatop: https://01.org/numatop