



Java and the modern CPU,
Part 1: Memory and the cache
hierarchy

Instructions, CPUs, memory,
and cache

Data size and L1, L2, and L3
caches

Data access patterns and the
perf tool

Using the Microbenchmark
Harness

Conclusion

Dig deeper

JAVA SE

Java and the modern CPU, Part 1: Memory and the cache hierarchy

You can understand application performance—and optimize your software approach—by understanding how CPUs, memory, and caches affect execution.

by *Michael Heinrichs*

January 29, 2021

[*Java Magazine* is pleased to republish a popular series on CPUs that was written in 2015. We have made only a few updates to the original articles.—*Ed.*]

Imagine an array with 67,000 integer elements. You run two loops over the array, as shown in **Listing 1**. Both loops multiply the elements of the array by three. However, while the first loop changes every element, the second loop modifies only every 16th element. How much faster will the second loop be compared to the first? Take a guess: 1/16th of the time?

Listing 1. Which loop will run faster?

```
private static final int ARRAY_SIZE = 64 * 10
public int[] array = new int[ARRAY_SIZE];
for (int i = 0, n = array.length; i < n; i++)
    array[i] *= 3;
}
for (int i = 0, n = array.length; i < n; i+=16)
    array[i] *= 3;
}
```

The perhaps surprising answer is that if the code is executed on a typical laptop, both loops take roughly the same amount of time to run. **Figure 1** shows measurements from three computers, and as you can see, the difference is negligible. The

second loop does only a fraction of the work, so how is it possible that the first loop runs at the same speed?

	SMALL STEPS EACH ELEMENT	LARGE STEPS EVERY 16TH ELEMENT
I7-4980HQ @ 2.8 GHZ, MAC OS X YOSEMITE	30.4 MS	29.7 MS
I7-3770 @ 3.4 GHZ, LINUX MINT 14	25.8 MS	26.1 MS
T7200 @ 2 GHZ, LINUX MINT 14	193.0 MS	184.2 MS

Figure 1. Comparing the performance of the two loops on three different machines

To understand this behavior, consider how the CPU and the memory system work. On the lowest level, modern computers can show surprising behavior, very much like quantum mechanics, which appears to contradict daily experience. But sometimes quantum mechanics has noticeable effects on the real world. And sometimes the effects of processes at the hardware level have a noticeable effect on programs. This article takes a look at how modern computers work at the lowest level and explores the things that can affect performance.

Instructions, CPUs, memory, and cache

If you try to imagine how the loops in **Listing 1** are executed, your initial interpretation might be that the array is stored in main memory and the CPU reads element after element, multiplies each element by three, and writes the result back, as you can see in **Figure 2**. This interpretation is useful for understanding the *functionality* of the loops, but it's not what really happens inside a computer.

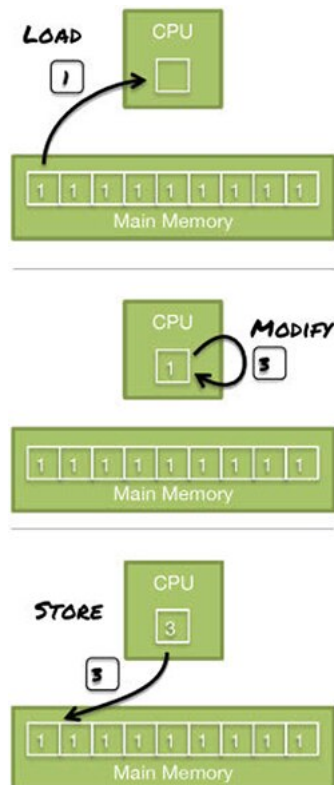


Figure 2. A functional view of the sample code

Figure 3 shows a graph of relative performance improvements that CPUs and memory went through in recent decades. Memory performance improved steadily during the whole period, but that was nothing compared to the improvements in CPU speed, especially during the 1990s. In recent years, plain CPU speed hit a limit, but do not be fooled! The scale in **Figure 3** is logarithmic. Even though it might look as if memory performance is catching up, the gap is still huge.

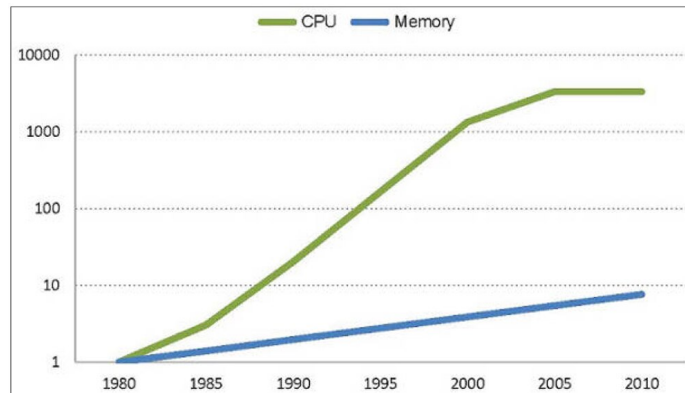


Figure 3. Relative improvements in processor and memory speed over the past three decades (this is a logarithmic scale)

The reality is that if a computer worked exactly as imagined in **Figure 2**, it would be terribly inefficient, because the superfast CPU would wait most of the time for the far-slower main memory to deliver the next element.

To overcome this bottleneck, processor designers added a small memory cache between the CPU and main memory. The cache is a much faster memory module, whose whole purpose is to mitigate the performance gap. **Figure 4** shows an improved model of the CPU and memory system.

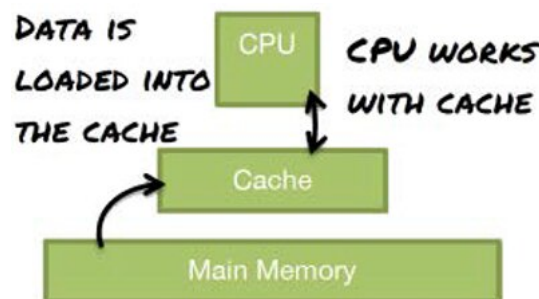


Figure 4. Adding cache into the functional diagram

Programs tend to access the same data and code several times within a short period (*temporal locality*), and memory access is often limited to small regions (*spatial locality*). This means that if you load all the data that you use into the cache, there is a high chance that you'll need it again later. Because the next time you need the data it's possibly already in the cache, the performance of your programs increases tremendously.

Now you might wonder: If faster memory can be put between the CPU and main memory, why can't all of the computer's memory be made faster? There are two main reasons:

- Main memory is a lot larger than cache, and it simply takes more time to find the right address within 16 GB (the typical size of main memory, as of this writing) than to find the right address within 8 KB (the typical size of Level 1 [L1] cache).
- The electronic components of cache memory are much more expensive than the ones used in main memory in terms of heat and space. Heat and space are limiting factors in modern chip design and, indeed, in modern computer design overall.

To exploit spatial locality, the cache doesn't work with individual bytes but uses cache lines instead. A *cache line* is an adjacent part of the memory, typically 64 bytes.

Keeping cache lines in mind, what really happens when you iterate over the large loops from **Listing 1** can be seen in **Figure 5**. The CPU loads a complete cache line from main memory into the cache and modifies the elements in the cache directly. The first loop modifies all elements in the cache line, while the second loop modifies only one element (16 integers, each 4 bytes long).

The limiting factor in this setup is loading the cache line into the cache; it almost doesn't matter how many operations you execute on each cache line. This explains why the performance of both loops is roughly the same.

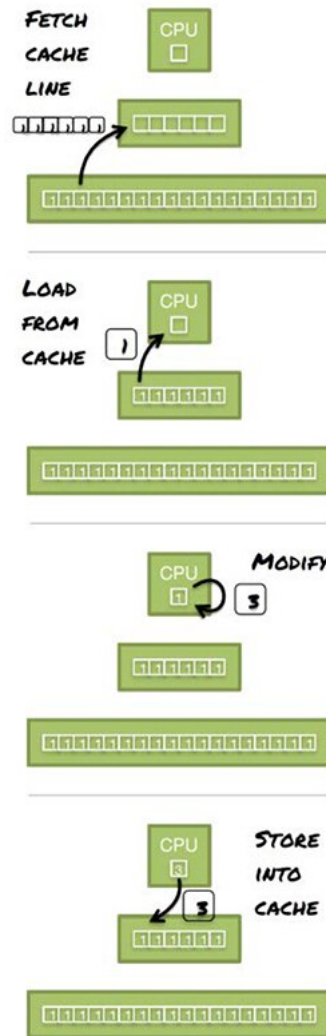


Figure 5. Adding cache lines to the functional diagrams shows why the two loops execute in almost the same amount of time.

Counting instructions to estimate the performance of an algorithm is a useful approximation, because that method is easy and usually gives a good indication. But as this example shows, you have to keep in mind that it's just an approximation. In reality, the execution times of single instructions vary widely, and you can't rely on this number only.

Data size and L1, L2, and L3 caches

Does the size of a data structure affect software's runtime performance? To answer this question, try a small experiment using the code in **Listing 2**. Take the second loop from the first code example and run it repetitively. This time, however, change the size of the array and measure the average time to run a single loop iteration.

The purpose of this experiment is to run a trivial algorithm over a data structure whose size you can control. This should show the relationship between the size of the array and the time needed to modify a single element.

Listing 2. A test to see how varying a data structure's size might affect runtime performance

```
private static final int ARRAY_CONTENT = 777;
@Param({"1024", "2048", "4096", "8192", "16384", "32768", "65536", "131072", "262144", "524288", "1048576", "2097152", "4194304", "8388608", "16777216", "33554432", "67108864", "134217728", "268435456", "536870912", "1073741824", "2147483648", "4294967296", "8589934592", "17179869184", "34359738368", "68719476736", "137438953472", "274877906944", "549755813888", "1099511627776", "2199023255552", "4398046511104", "8796093022208", "17592186044416", "35184372088832", "70368744177664", "140737488355328", "281474976710656", "562949953421312", "1125899906842624", "2251799813685248", "4503599627370496", "9007199254740992", "18014398509481984", "36028797018963968", "72057594037927936", "144115188075855872", "288230376151711744", "576460752303423488", "1152921504606846976", "2305843009213693952", "4611686018427387904", "9223372036854775808", "18446744073709551616", "36893488147419103232", "73786976294838206464", "147573952589676412928", "295147905179352825856", "590295810358705651712", "1180591620717411303424", "2361183241434822606848", "4722366482869645213696", "9444732965739290427392", "18889465931478580854784", "37778931862957161709568", "75557863725914323419136", "151115727451828646838272", "302231454903657293676544", "604462909807314587353088", "1208925819614629174706176", "2417851639229258349412352", "4835703278458516698824704", "9671406556917033397649408", "19342813113834066795298816", "38685626227668133590597632", "77371252455336267181195264", "154742504910672534362390528", "309485009821345068724781056", "618970019642690137449562112", "1237940039285380274899124224", "2475880078570760549798248448", "4951760157141521099596496896", "9903520314283042199192993792", "19807040628566084398385987584", "39614081257132168796771975168", "79228162514264337593543950336", "158456325028528675187087900672", "316912650057057350374175801344", "633825300114114700748351602688", "1267650600228229401496703205376", "2535301200456458802993406410752", "5070602400912917605986812821504", "10141204801825835211973625643008", "20282409603651670423947251286016", "40564819207303340847894502572032", "81129638414606681695789005144064", "162259276829213363391578010288128", "324518553658426726783156020576256", "649037107316853453566312041152512", "1298074214633706907132624082305024", "2596148429267413814265248164610048", "5192296858534827628530496329220096", "10384593717069655257060992658440192", "20769187434139310514121985316880384", "41538374868278621028243970633760768", "83076749736557242056487941267521536", "166153499473114484112975882535043072", "332306998946228968225951765070086144", "664613997892457936451903530140172288", "1329227995784915872903807060280344576", "2658455991569831745807614120560689152", "5316911983139663491615228241121378304", "10633823966279326983230456482242756608", "21267647932558653966460912964485513216", "42535295865117307932921825928971026432", "85070591730234615865843651857942052864", "170141183460469231731687303715884105728", "340282366920938463463374607431768211456", "680564733841876926926749214863536422912", "1361129467683753853853498429727072845824", "2722258935367507707706996859454145691648", "5444517870735015415413993718908291383296", "10889035741470030830827987437816582766592", "21778071482940061661655974875633165533184", "43556142965880123323311949751266331066368", "87112285931760246646623899502532662132736", "174224571863520493293247799005065324265472", "348449143727040986586495598010130648530944", "696898287454081973172991196020261297061888", "1393796574908163946345982392040522594123776", "2787593149816327892691964784081045188247552", "5575186299632655785383929568162090376495104", "11150372599265311570767859136324180752990208", "22300745198530623141535718272648361505980416", "44601490397061246283071436545296723011960832", "89202980794122492566142873090593446023921664", "178405961588244985132285746181186892047843328", "356811923176489970264571492362373784095686656", "713623846352979940529142984724747568191373312", "1427247692705959881058285969449495136382746624", "2854495385411919762116571938898990272765493248", "5708990770823839524233143877797980545530986496", "11417981541647679048466287755595961091061972992", "22835963083295358096932575511191922182123945984", "45671926166590716193865151022383844364247891968", "91343852333181432387730302044767688728495783936", "182687704666362864775460604089535377456991567872", "365375409332725729550921208179070754913983135744", "730750818665451459101842416358141509827966271488", "1461501637330902918203684832716283019655932542976", "2923003274661805836407369665432566039311865085952", "5846006549323611672814739330865132078623730171904", "11692013098647223345629478661730264157247460343808", "23384026197294446691258957323460528314494920687616", "46768052394588893382517914646921056628989841375232", "93536104789177786765035829293842113257979682750464", "187072209578355573530071658587684226515959365500928", "374144419156711147060143317175368453031918731001856", "748288838313422294120286634350736906063837462003712", "1496577676626844588240573268701473812127674924007424", "2993155353253689176481146537402947624255349848014848", "5986310706507378352962293074805895248510699696029696", "11972621413014756705924586149611790497021399392059392", "23945242826029513411849172299223580994042798784118784", "47890485652059026823698344598447161988085597568237568", "95780971304118053647396689196894323976171195136475136", "191561942608236107294793378393788647952342390272950272", "383123885216472214589586756787577295904684780545900544", "766247770432944429179173513575154591809369561091801088", "1532495540865888858358347027150309183618739122183602176", "3064991081731777716716694054300618367237478244367204352", "6129982163463555433433388108601236734474956488734408704", "12259964326927110866866776217202473468949912977468817408", "24519928653854221733733552434404946937899825954937634816", "49039857307708443467467104868809893875799651909875269632", "98079714615416886934934209737619787751599303819750539264", "196159429230833773869868419475239575503198607639501078528", "392318858461667547739736838950479151006397215279002157056", "784637716923335095479473677900958302012794430558004314112", "1569275433846670190958947355801916604025588861116008628224", "3138550867693340381917894711603833208051177722232017256448", "6277101735386680763835789423207666416102355444464034512896", "12554203470773361527671578846415332832204710888928069025792", "25108406941546723055343157692830665664409421777856138051584", "50216813883093446110686315385661331328818843555712276103168", "100433627766186892221372630771322662657637687111424552206336", "200867255532373784442745261542645325315275374222849104412672", "401734511064747568885490523085290650630550748445698208825344", "803469022129495137770981046170581301261101496891396417650688", "1606938044258990275541962092341162602522202993782792835301376", "3213876088517980551083924184682325205044405987565585670602752", "6427752177035961102167848369364650410088811975131171341205504", "12855504354071922204335696738729300820177623950262342682411008", "25711008708143844408671393477458601640355247900524685364822016", "51422017416287688817342786954917203280710495801049370729644032", "102844034832575377634685573909834406561420991602098741459288064", "205688069665150755269371147819668813122841983204197482918576128", "411376139330301510538742295639337626245683966408394965837152256", "822752278660603021077484591278675252491367932816789931674304512", "1645504557321206042154969182557350504982735865633579863348609024", "3291009114642412084309938365114701009965471731267159726697218048", "6582018229284824168619876730229402019930943462534319453394436096", "13164036458569648337239753460458804039861886925068638906788872192", "26328072917139296674479506920917608079723773850137277813577744384", "52656145834278593348959013841835216159447547700274555627155488768", "105312291668557186697918027683670432318895095400549111254310977536", "210624583337114373395836055367340864637790190801098222508621955072", "421249166674228746791672110734681729275580381602196445017243910144", "842498333348457493583344221469363458551160763204392890034487820288", "1684996666696914987166688442938726917102321526408785780068975640576", "3369993333393829974333376885877453834204643052817571560137951281152", "6739986666787659948666753771754907668409286105635143120275902562304", "13479973333575319897333507543509815336818572211270286240551805124608", "26959946667150639794667015087019630673637144422540572481103610249216", "53919893334301279589334030174039261347274288845081144962207220498432", "107839786668602559178668060348078522694548577690162289924414440996864", "215679573337205118357336120696157045389097155380324579848828881993728", "431359146674410236714672241392314090778194310760649159697657763987456", "862718293348820473429344482784628181556388621521298319395315527974912", "1725436586697640946858688965569256363112777243042596638790631055949824", "3450873173395281893717377931138512726225554486085193277581262111899648", "6901746346790563787434755862277025452451108972170386555162524223799296", "13803492693581127574869511724554050904902217944340773110325048447598592", "27606985387162255149739023449108101809804435888681546220650096895197184", "55213970774324510299478046898216203619608871777363092441300193790394368", "110427941548649020598956093796432407239217743554726184882600387580788736", "220855883097298041197912187592864814478435487109452369765200775161577472", "441711766194596082395824375185729628956870974218904739530401550323154944", "883423532389192164791648750371459257913741948437809479060803100646309888", "1766847064778384329583297500742918515827483896875618958121606201292619776", "3533694129556768659166595001485837031654967793751237916243212402585239552", "7067388259113537318333190002971674063309935587502475832486424805170479104", "14134776518227074636666380005943348126619871175004951664972849610340958208", "28269553036454149273332760011886696253239742350009903329945699220681916416", "56539106072908298546665520023773392506479484700019806659891398441363832832", "113078212145816597093331040047546785012958969400039613319782796882727665664", "226156424291633194186662080095093570025917938800079226639565593765455331328", "452312848583266388373324160190187140051835877600158453279131187530910662656", "904625697166532776746648320380374280103671755200316906558262375061821325312", "1809251394333065553493296640760748560207343510400633813116524750123642650624", "3618502788666131106986593281521497120414687020801267626233049500247285301248", "7237005577332262213973186563042994240829374041602535252466099000494570602496", "14474011154664524427946373126085988481658748083205070504932198000989141204992", "28948022309329048855892746252171976963317496166410141009864396001978282409984", "57896044618658097711785492504343953926634992332820282019728792003956564819968", "115792089237316195423570985008687907853269984665640564039457584007913129639936", "231584178474632390847141970017375815706539969331281128078915168015826259279872", "463168356949264781694283940034751631413079938662562256157830336031652518559744", "926336713898529563388567880069503262826159877325124512315660672063305037119488", "1852673427797059126777135760139006525652319754650249024631321344126610074238976", "3705346855594118253554271520278013051304639509300498049262642688253220148477952", "7410693711188236507108543040556026102609279018600996098525285376506440296955904", "14821387422376473014217086081112052205218558037201992197050570753012880593911808", "29642774844752946028434172162224104410437116074403984394101141506025761187823616", "59285549689505892056868344324448208820874232148807968788202283012051522375647232", "118571099379011784113736688648896417641748464297615937576404566024103044751294464", "237142198758023568227473377297792835283496928595231875152809132048206089502588928", "474284397516047136454946754595585670566993857190463750305618264096412179005177856", "948568795032094272909893509191171341133987714380927500611236528192824358010355712", "1897137590064188545819787018382342682267975428761855001222473056385648716020711424", "3794275180128377091639574036764685364535950857523710002444946112771297432041422848", "7588550360256754183279148073529370729071901715047420004889892225542594864082845696", "15177100720513508366558296147058741458143803430094840009779784451085189728165691392", "30354201441027016733116592294117482916287606860189680019559568902170379456331382784", "60708402882054033466233184588234965832575213720379360039119137804340758912662765568", "121416805764108066932466369176469931665150427440758720078238275608681517825325531136", "242833611528216133864932738352939863330300854881517440156476551217363035650651062272", "48566722305643226772986547670587972666
```

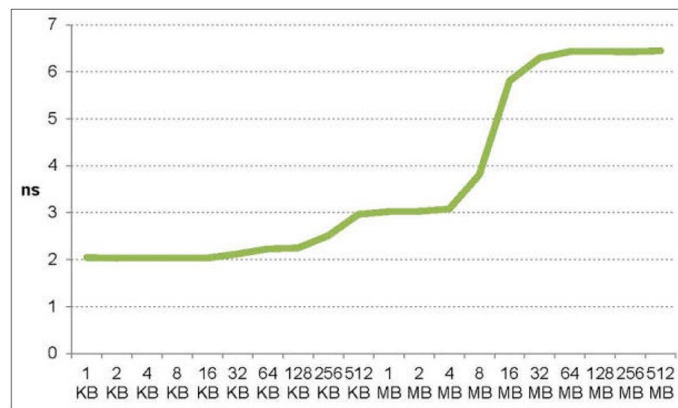


Figure 6. The dependency between array size and access time

The cache is usually not a single unit but instead consists of several hierarchical levels with different sizes and access times. As shown in **Figure 7**, the L1 cache is the smallest and fastest. L2 is larger and much slower. L3 is even larger and slower still—but still significantly faster than main memory.

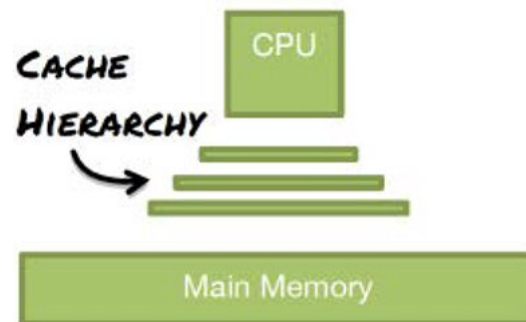


Figure 7. The cache hierarchy, going from the CPU to L1 to L2 to L3 to main memory

How large are the performance gaps between the different cache levels? To explain this in a form that is more accessible to human beings, my former colleague Richard Thompson came up with the beer cache hierarchy. Imagine that you're sitting in front of your TV watching your favorite team and you're thirsty.

- L1 cache is the bottle of beer in your hand. Access time is almost immediate (< 1 ns), but the quantity is extremely limited (for example, 32 KB on my computer).
- L2 cache is the cooler next to your sofa. Access time is still pretty low (7 ns), and the quantity is significantly larger (256 KB, which is equivalent to 8 bottles of beer).
- L3 cache is the fridge in the kitchen. Access time is noticeably larger (25 ns), but the size is so large that the analogy falls apart (8 MB, which is equivalent to 256 bottles of beer).
- Main memory is the corner store. Access time is huge (100 ns), but the quantity of beer is probably more than enough for a lifetime (16 GB, which is equivalent to more than half a million bottles of beer).

Looking at these numbers, it becomes quite obvious why both loops in the initial example took the same amount of time, because it doesn't really matter how many sips of beer you drink if you have to run to the corner store for each bottle.

With the cache level hierarchy in mind, look back at the graph in **Figure 6**. Each plateau in the graph corresponds to a level of the cache hierarchy. As long as the array fits into the L1 and L2 caches, access time is very low. But as soon as the array becomes too large and has to be read from the L3 cache, access time increases noticeably. And the same happens again as soon as the array does not fit into the L3 cache and has to be read from main memory. If you look closely, you can even see the small jump between the L1 and L2 cache.

Size does matter. Even though memory is cheaper than ever before, try to avoid wasting it. The smaller the amount of data you use for a particular data structure, the higher the chance that it will fit into the cache, which can lead to significantly better performance.

Data access patterns and the perf tool

The *size* of data influences performance. Does the *order* in which you access your data—the data access pattern—have an influence, too? You can change the previous experiment slightly to find an answer. Instead of simply running an index through the array, create a second array that stores the access order. Access the array sequentially as before for one time, and then access it randomly and measure the difference. You can see the code for both experiments in **Listing 3**.

Listing 3. Changing the order of data access

```
public int[] rndIndex;
@Setup(Level.Iteration)
public void setUp() {
    ...
    rndIndex = new int[indexes];
    final List<Integer> list = new ArrayList<>(indexes);
    for (int i=0; i<indexes; i++) {
        list.add(16 * i);
    }
    Collections.shuffle(list);
    for (int i=0; i<indexes; i++) {
        rndIndex[i] = list.get(i);
    }
}
```

If you run the experiment with different array sizes and plot the result in a graph, you get two curves, as shown in **Figure 8**. Not surprisingly, you can see the already familiar staircase pattern. Both curves show similar access times on the lower two levels, which correlate to the L1 and L2 caches.

But on the third level, the performance of the sequential access pattern is noticeably better. The fourth level shows a significant difference. Why is the access order insignificant for small arrays but significant for large arrays?

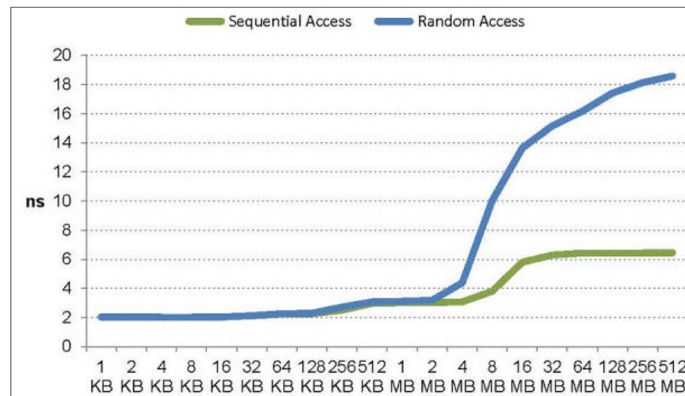


Figure 8. Sequential access and random access performance

To get a better understanding of what is going on inside the computer, you can use the Linux profiler tool `perf`, which collects and prints out events generated by the CPU and the memory system while a program is executed. The command-line interface for `perf` is similar to that of `git`. You call `perf` with the command you want to execute, for example:

```
perf COMMAND [ARGS]
```

To get a list of all commands, use `perf --help`. To get help for a specific command, you can run this:

```
perf help COMMAND
```

The most useful command is `stat`, which allows `perf` to run another program and tracks hardware events during execution, for example:

```
perf stat [ARGS] PROGRAM
```

Without any arguments, this command starts `PROGRAM`, tracks some general events, and prints out the statistics as soon as the program ends. You can see a typical output in **Figure 9**.

```
michael@desktop ~/repos/quantum/target/classes $ perf stat -p 3195 -B -r 40 sleep 5
Performance counter stats for process id '3195' (40 runs):
      4997,277608 task-clock                #    0,999 CPUs utilized          ( +- 0,00% ) [100,00%]
           529 context-switches           #    0,106 K/sec                  ( +- 0,04% ) [100,00%]
             0 CPU-migrations              #    0,000 K/sec                  [100,00%]
             0 page-faults                 #    0,000 K/sec                  ( +- 35,51% )
    19,429,967,708 cycles                   #    3,888 GHz                    ( +- 0,00% ) [83,28%]
    19,006,043,778 stalled-cycles-frontend  #    97,53% frontend cycles idle   ( +- 0,00% ) [83,31%]
    18,296,349,545 stalled-cycles-backend  #    94,17% backend cycles idle    ( +- 0,01% ) [66,74%]
    1,386,651,837 instructions              #    0,07 insns per cycle          ( +- 0,05% ) [83,37%]
           23,934,698 branches              #    13,71 stalled cycles per insn ( +- 0,07% ) [83,37%]
             23,139 branch-misses           #    0,10% of all branches         ( +- 11,00% ) [83,30%]

5,00586757 seconds time elapsed            ( +- 0,00% )
```

Figure 9. A typical output from the `perf` tool

You can specify which events should be tracked by using the `-e` option. To get a list of supported events, run the command `perf list`.

In Java, you usually don't want to track the whole program, because this would include bootstrapping the VM, just-in-time (JIT) compilation, and so on. Fortunately, with `perf` you can hook into a running process with the `-p` option. You have to specify a program that will be executed, and the measurement ends once this program ends. You can use the `sleep` command to specify the duration of the test.

First, measure both your loops using the default settings of `perf`. This provides a good overview. The most-interesting results can be seen in **Figure 10**.

	SEQUENTIAL ACCESS		RANDOM ACCESS	
CYCLES	19,430,435,800		19,429,967,708	
STALLED FRONT-END CYCLES	7,217,361,632	(37.14%)	19,006,043,778	(97.82%)
STALLED BACK-END CYCLES	843,462,646	(4.34%)	18,296,349,545	(94.17%)

Figure 10. Running the `perf` test with the default settings

The number of stalled front-end and back-end cycles differs significantly between both runs. A stalled cycle means the CPU is idle and waiting for something. One of the most likely causes of a stalled front-end cycle is a cache miss, which results in the CPU waiting for data to arrive from main memory or a slower cache. To validate this assumption, you can run `perf` again, but this time run it to specifically measure cache loads and cache misses. You can see the result in **Figure 11**.

	SEQUENTIAL ACCESS		RANDOM ACCESS	
L1 CACHE LOADS	5,758,001,370		170,655,221	
L1 CACHE MISSES	360,757,378	(6.27%)	365,959,699	(214.44%)

Figure 11. The `perf` test measuring cache loads and cache misses

The ratio between successful cache loads and cache misses differs tremendously. When the array is accessed in sequential order, only about 6% of all memory loads result in a cache miss. But when the array is accessed in random order, two out of three loads result in a cache miss. The high number of cache misses is expected, because the array doesn't fit into the cache, and you have to load everything from main memory. But why are there almost no cache misses when you access the array sequentially?

The time it takes to load data from main memory into the cache hierarchy is often a major bottleneck. For this reason, the CPU tries to help by guessing which data you'll use next and loading it into the cache in the background, as you can see in **Figure 12**. While you modify the elements of a cache line, functionality called the *prefetcher* loads the next cache line into the cache. Thus, when you need the data, it's *already* available in the cache.

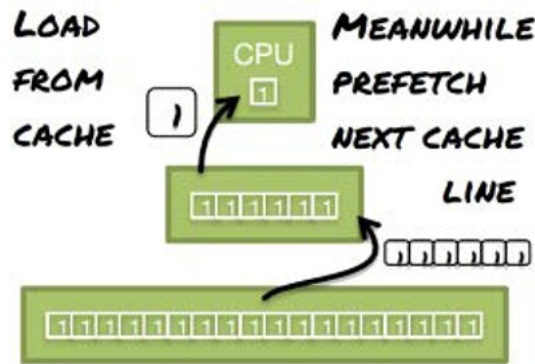


Figure 12. The prefetcher loads information from main memory into the cache line.

The prefetcher is not particularly smart. It can guess the next memory location correctly only if the memory loads follow a regular pattern, for example:

- In the first case, when you went through the array sequentially, guessing the next memory location was easy, and the prefetcher could mitigate a substantial part of the performance loss by prefetching the next memory location. Most of the time, that's a win.
- But when you accessed the array randomly, guessing the next memory location correctly was impossible, and the algorithm had to wait until data was loaded from main memory. Still, there was no downside to the prefetching.

The access pattern has a significant influence on the performance of an algorithm, but the chances to use this knowledge within Java are limited. You have close to no control over how your data is arranged in memory.

Using the Microbenchmark Harness

Microbenchmarking is another valuable tool for getting more insight into your programs. Probably the most important rule of microbenchmarking is to always use a tool that helps you avoid some of the many pitfalls, such as the insufficient warmup of the VM, dead code elimination, and loop unrolling.

The [Java Microbenchmark Harness \(JMH\)](#) is probably the best harness available right now. JMH is a Java harness for building, running, and analyzing microbenchmarks written in Java and other languages targeting the JVM.

Tests written for JMH are similar to JUnit tests. The code you want to benchmark needs to be in a single method, which must be annotated with `@Benchmark`. The test can be configured with annotations at the class level. **Figure 13** shows the most-important annotations and their meaning. You should experiment with this tool to see how it works.

ANNOTATION	DESCRIPTION
@BENCHMARKMODE	SPECIFIES WHAT SHOULD BE MEASURED—FOR EXAMPLE, THROUGHPUT OR AVERAGE TIME.
@OUTPUTTIMEUNIT	SPECIFIES THE TIME UNIT USED IN THE OUTPUT—FOR EXAMPLE, <code>TimeUnit.MILLISECONDS</code> .
@WARMUP	SPECIFIES THE WARMUP PHASE. YOU CAN SET THE NUMBER OF ITERATIONS AND THE DURATION OF A SINGLE ITERATION.
@MEASUREMENT	SPECIFIES THE MEASUREMENT PHASE AND IS SIMILAR TO @WARMUP.
@FORK	SPECIFIES HOW OFTEN YOU WANT TO FORK THE JAVA VIRTUAL MACHINE (JVM) AND RUN THE TESTS. YOU SHOULD ALWAYS DO RUNS IN SEVERAL FORKS.

Figure 13. Important annotations for the Java Microbenchmark Harness

Conclusion

Most of the time, processes at the hardware level have no significant effect on programs, but sometimes they do. Therefore, it's useful to have a rough understanding of what goes on at the hardware level.

This article focused on memory and, specifically, the cache hierarchy. In the second part of this series, I extend the investigation to the issues of multithreaded access. In the third and final part, I look at the internal workings of the CPU itself.

Dig deeper

- [How L1 and L2 CPU caches work, and why they're an essential part of modern chips](#)
- [Java Microbenchmark Harness](#)
- [Comparing workload performance](#)
- [Book review: Optimizing Java](#)



Michael Heinrichs

Michael Heinrichs is co-founder of Karakun AG in Freiburg, Germany. He's a Java Champion as well as founder and leader of the Java user group in Freiburg. Follow him on Twitter at [@net0pyr](#).

Share this Page



Contact

US Sales: +1.800.633.0738
Global Contacts
Support Directory
Subscribe to Emails

About Us

Careers
Communities
Company Information
Social Responsibility Emails

Downloads and Trials

Java for Developers
Java Runtime Download
Software Downloads
Try Oracle Cloud

News and Events

Acquisitions
Blogs
Events
Newsroom

