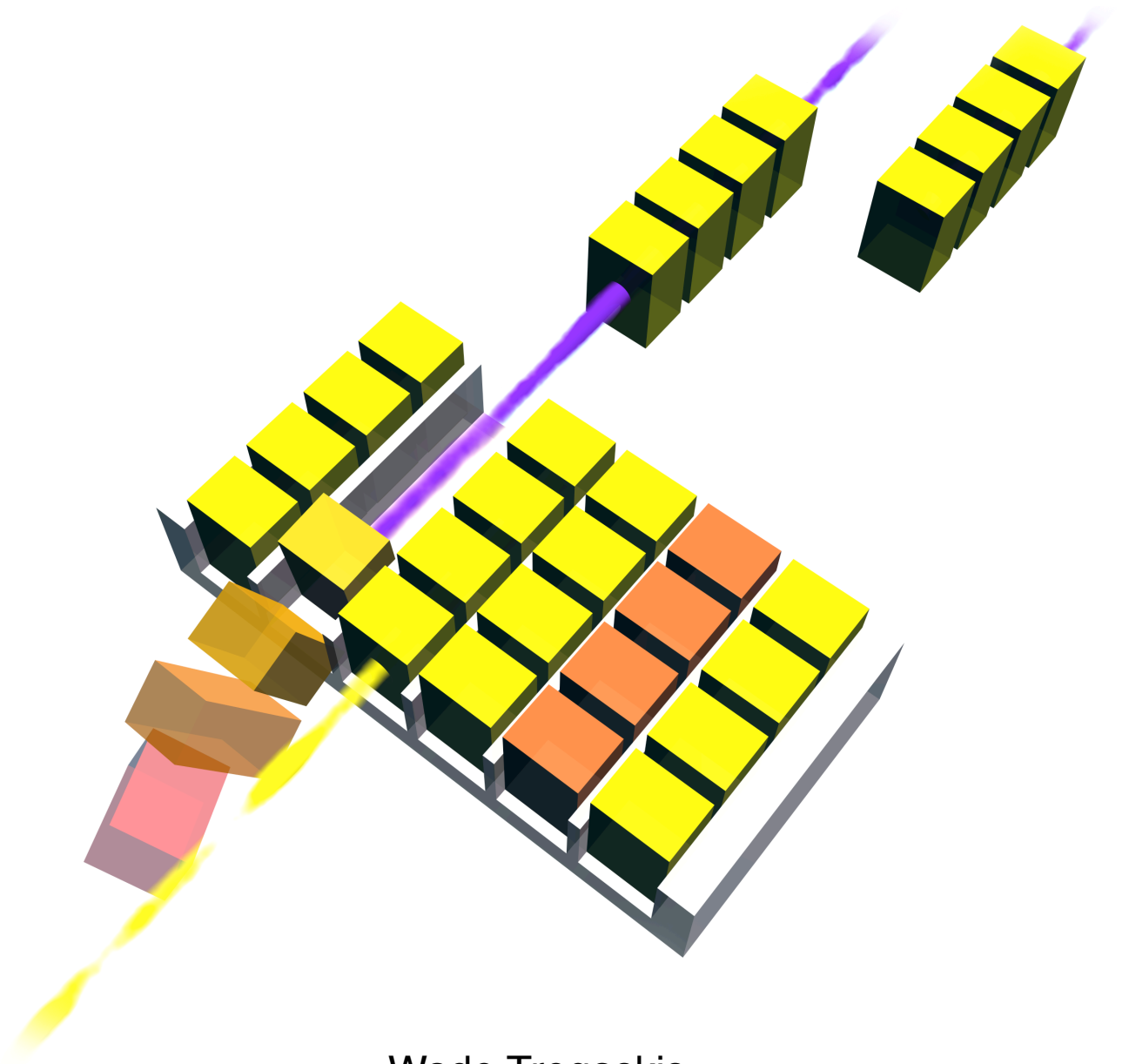


CSE32ARC

SimpleCacheSim



Wade Tregaskis

02557793

Monday, 10 October 2005

Contents

Introduction	1
Overview	1
Usage	1
<i>Program Parameters</i>	2
<i>Output Formats</i>	3
<i>Additional Notes</i>	4
Design	4
Design Wins	5
Design Losses	6
Sample Simulation	7
Method	7
Data	7
Performance as related to cache size	7
<i>Total Miss Rates</i>	8
Discussion	11
<i>Miss Breakdown</i>	12
Discussion	15
Performance as related to block size	15
<i>Total Miss Rates</i>	16
Discussion	19
<i>Miss Breakdown</i>	19
Discussion	23
Ideal Cache	23

Introduction

Overview

SimpleCacheSim is a small program for simulating the operation of memory caches. It collects statistics about the general operation of the cache under simulation - such as how many requests to the cache hit, how many missed, etc. It breaks the statistics down by operation type and exact miss type, producing a thorough 6x8 matrix of results.

It is capable of simulating several replacement policies, and caches of any size or reasonable configuration. It allows specification of block (a.k.a. line) sizes and the system word size, which then allows it to detect misaligned requests - and optionally handle them, invoking multiple cache accesses from a single request.

It also features a batch-orientated interface and mode of operation, allowing a range of parameters to be specified. The program then simulates every valid permutation of these parameters, outputting the results en-mass. It has been used for batches involving over 14,000 simulations at a time. Its output can be configured to be in numerous forms, including tabular and SQL-like, allowing results to be imported into a database (e.g. PostgreSQL) with ease.

What it does not do is simulate multiple levels of cache, nor non-unified caches. For such simulation you may wish to try the free `dinerolV` cache simulator.

Additionally, it does not perform any timing simulation. As such, all accesses are sequential. In reality many modern hardware memory caches allow multiple concurrent requests, and many also support pipelining. There are no immediate plans, at this point in time, to add such functionality.

Usage

The program has a typical *nix command line interface. It is invoked from a terminal with various parameters defining the simulation(s) to perform. It requires a list of cache requests to simulate, which can be provided via standard input or from an input file specified as a program parameter. The program is invoked as such:

```
SimpleCacheSim [OPTIONS] [INPUT]
```

All parameters are optional - without any other direction, the program assumes a generic cache configuration and awaits input. If a file is named in the parameters, that file is opened and read from instead of standard input. This has the same effect as piping the contents of the file to the program.

The options that can be specified to the program are provided on the following pages.

Program Parameters

<code>--addressbase/-A</code>	The numerical base of addresses provided as input. Useful if the addresses are hexadecimal but without a preceding "0x", which would otherwise confuse the input parser. Valid values are 2 through 36 inclusive, or the special value 0, which allows any base by trying to detect the base in use (e.g. any input starting with "0x" is hexadecimal, any starting with "0" is assumed to be octal, and if all else fails the input is assumed to be decimal. Defaults to 0.
<code>--misalignment/-M</code>	The types of misalignment to permit, if any. Valid values are "none", "word", "block" (which implies allowing word misalignment as well) or "forced" (where all addresses are forced to align). Defaults to None.
<code>--SQL/-S</code>	SQL functionality. Takes a string as an argument, which may be either "create" or "insert". The "create" directive instructs the program to output the SQL CREATE statements that generate a table and summary views suitable for holding the tabular output from this program. These statements can be executed in a database client - e.g. psql for PostgreSQL. The "insert" directive instructs the program to output all tabular results as SQL INSERT statements. The output can then be executed directly in a database client to import data into the table(s) generated by the "create" directive. Note: if your database supports a COPY command, it is usually much faster to use this on the normal tabular output than to use INSERT statements. The "insert" directive is only meaningful if a negative verbosity is specified.
<code>--associativity/-a</code>	The associativity of the cache, which may be either "direct", "full" or a numeric value indicating the number of ways mapped. Multiple associativities may be separated by commas, and/or provided as ranges distinguished by hyphens. Defaults to Full.
<code>--blocksize/-b</code>	The size (in bytes) of each cache block. This is the smallest unit that the cache will be capable of holding (not to be confused with the smallest unit which can be transferred, which is the word size). Must be an integer multiple of the word size. Multiple block sizes may be separated by commas, and/or provided as ranges distinguished by hyphens. Defaults to 0.
<code>--cachesize/-c</code>	The size (in bytes) of the cache. Must be an integer multiple of block size and associativity. Multiple cache sizes may be separated by commas, and/or provided as ranges distinguished by hyphens. Defaults to 16384.
<code>--help/-h</code>	Prints this usage information and exits.
<code>--memorysize/-m</code>	The size (in bytes) of the memory behind the cache. Optional, used only to validate input values. A value of 0 indicates no defined limit. Note that only one value is permitted, unlike other parameters, as this is purely for input validation. Defaults to 0.

<code>--replacement/-r</code>	The replacement policy in N-way associative caches ($N > 0$). May be any of "random", "LRU" (Least Recently Used), "FIFO" (First In First Out), "LRR" (Least Recently Read) or "LRW" (Least Recently Written). Multiple policies may be separated by commas. Defaults to LRU.
<code>--verbosity/-v</code>	Sets the verbosity of the output, where 0 corresponds to silent operation, 1 prints out the simulation parameters and other key information, 2 also prints out information about every cache operation as it is simulated, and 3 also prints out even more details about what's happening with each and every cache access. A negative verbosity indicates tabular output; -1 outputs only the raw data, while -2 or lower also outputs a header on the first row, indicating what value each column represents. Use high verboisities with caution - the volume of output can be tremendous. Defaults to 0.
<code>--wordsize/-w</code>	The size (in bytes) of one word, which is the smallest unit of data that can be transferred to or from cache. Multiple word sizes may be separated by commas, and/or provided as ranges distinguished by hyphens. Defaults to 4.

Output Formats

The final output of the program has three forms. The first is "standard" output, which is a form designed to be succinct and human-readable. This is the form used for a verbosity setting of 0 or above. An example is shown below. If a non-essential value is zero, it is omitted from the output for brevity.

```

Hits: 91922
      Read: 12998
      Write: 11213
      Instruction: 67711
Capacity Misses: 0
Compulsory Misses: 12
      Read: 3
      Write: 2
      Instruction: 7
Conflict Misses: 8066
      Read: 2690
      Write: 1343
      Instruction: 4033
Misalignments straddling blocks: 0
Misalignments within block: 0

Total misses: 8078
Total cache accesses: 100000

Miss rate: 0.080780
      Compulsory: 0.000120
      Capacity: 0.000000
      Conflict: 0.080660

```

Another form is for a verbosity setting of -1 or lower, without specifying the SQL INSERT format (see the --SQL/-S parameter in the parameter list, previous). This is standard tabular output, which is well suited for import into a database or spreadsheet for further processing. It is also explicit in all values - that is, it outputs all values regardless of whether they are zero or even directly relevant to the exact simulation performed. This consistency is essential for ease of use in automated systems. An example of tabular output is shown below.

100000	Direct	LRU	64	512	1	12998	11213	67711	0	0
0	0	0	0	0	0	00	0	0	0	3
2	7	0	0	0	0	0	2690	1343	4033	0
0	00	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

If the verbosity is set to -2 or lower, a heading is also printed on the first line, identifying the meaning of each column.

The last format is using the SQL INSERT functionality, by providing a negative verbosity setting as well as invoking the --SQL/-S parameter with an argument of "insert". The output in this case looks like the example shown below.

```
INSERT INTO SimpleCacheSim VALUES (100000, 'Direct', 'LRU', 64, 512, 1, 12998, 11213,
67711, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 3, 2, 7, 0, 0, 0, 0, 0, 0, 2690, 1343,
4033, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0);
```

In addition to these final outputs, a high positive verbosity setting also causes information to be output as the simulation is run - such as what each access is, notification of block movement in and out of the cache, etc. The higher verbosity settings (2 and above) are not recommended for normal use - their primary purpose is as debugging aids or teaching tools for very small simulations.

Additional Notes

Program performance is generally very good - for reference, on a 1.5 GHz Powerbook a simulation with 100,000 cache requests takes [on average] approximately 100ms. Note however that some simulations will take substantially longer than others - the worst offenders in this regard are simulations of fully associative cache, particularly large such caches.

Additionally, larger block and word sizes typically improve performance (as larger block sizes reduce the total number of blocks in the simulation).

Design

See Appendix III for the source code.

The program is written in C++, but favours a procedural approach. There is one core function - performSimulation() - which does the body of the work. main() parses input parameters and drives the simulation(s), and various other utility functions aid it in that endeavour. Execution flow is fairly trivial, and what isn't intuitive is covered in the code by appropriate commenting and inline documentation.

In terms of data structures and flow, again, the design is fairly simple. There are three primary structures - `CacheAccess`, `CacheEntry` and `CacheSet`. The first simply stores a single cache request, which involves recording the address, the access type and any other miscellaneous information necessary (such as which line the access came from in the original input). A single vector of such `CacheAccess`'s are created in `main()` by reading from the appropriate input stream, and this vector is passed to `performSimulation()`. It is reused across all calls to `performSimulation()` for performance reasons.

The `CacheEntry` structure represents a single block that is or has been in cache. It records the base address of the block, temporal information regarding when it was last loaded/used/written/read/etc, flags indicating whether it is dirty, whether it was prefetched, etc, and - importantly - a field containing the eviction reason for the last time it was thrown out of cache. The reason for this field's importance will be detailed shortly.

The `CacheSet` structure is trivial - it contains two vectors, one of which stores all the `CacheEntry`'s currently "in" the cache, while the other stores every `CacheEntry` that is or ever has been in the cache. The latter vector is used, obviously, primarily for determining why a block is not already in cache (i.e. if it does not exist in the vector, it is a compulsory miss, otherwise the reason it is not in cache can be read from the eviction reason field).

The `CacheSet` structure also contains an 'isFull' flag, indicating whether the cache has reached capacity. For performance reasons the vector of entries "in" cache is fixed at the number of possible entries, with empty spots being represented by NULL pointers. Thus `size()` or similar methods of the vector are not useful for indicating it's state. The 'isFull' flag is used as a common-case optimisation, as the capacity status of the cache is checked with every new block loaded into the set.

Finally, `performSimulation()` has a local variable - an array of `CacheSet`'s, representing the cache as a whole. A vanilla C array was used for simplicity and performance - the array is static for the life of the `performSimulation()` function, and need only be indexed; never resized or otherwise modified.

The following two subsections contrast the good and bad parts of this design, both in an absolute sense and also with an eye towards `dinerolV`, a "competing" cache simulator.

Design Wins

The design is simple and elegant. The regular approach and use of data - regardless of replacement algorithms or any other parameters - makes the implementation single-minded and straightforward. The use of vectors instead of linked lists, especially given adding elements is rare (and removals never) is a big performance win (in contrast to `dinerolV`, primarily).

In terms of features, there are some minor ones worth mentioning, as they are relatively rare in cache simulators. The first is the flexibility of input parameters. For example, it is perfectly valid to use a word size of 3, provided all other values (such as block size, cache size, etc) are multiples of 3. No other simulator known to the author supports this, and while it may seem gimmicky, it is testimony to the elegant and un-assuming design.

Additionally, several replacement policies are available that are usually overlooked, because they are not used in the real world. These are LRR (Least-Recently-Read) and LRW (Least-Recently-Written). While not representative of any real world caches, they are useful for academic purposes.

Arguably the biggest win, however, is the method for distinguishing between capacity and conflict misses. dinerolV and many other simulators have a serious design flaw in this regard - they distinguish between the two based on the state of the cache at the time of a miss. Unfortunately, the cache will likely have changed by this time, perhaps significantly. For example, at the time the block was actually evicted from cache, it may have been because of a conflict, even though the cache was not full. If the cache does become full before a miss occurs on that block, dinerolV and others will register the miss as a capacity miss, not a conflict miss. This is seen clearly in some of the comparative results in Appendix I, where dinerolV and SimpleCacheSim both recognise exactly the same total misses, and compulsory misses, but dinerolV overstates the portion of those which are capacity misses.

SimpleCacheSim stores the eviction reason in the CacheEntry as it is removed from cache, as previously noted. Thus, the eviction reason is based on the state of the cache at the actual time of the eviction, and consequently is always accurate.

Design Losses

One big area of concern is performance. While SimpleCacheSim performs admirably in most cases, for some borderline scenarios it's performance rapidly deteriorates. For example, high associativities are expensive, as are access patterns which hit a very large number of unique addresses only a few times each, causing a large number of CacheEntry's to be produced and stored, increasing the search time for each subsequent cache miss.

Another problem is scalability of the design going forward, when looking to introduce new functionality. A more object-orientated implementation would be more flexible in this regard. For example, the current implementation has no support for multiple caches, and cannot be easily configured as such.

Other small flaws include the inability for the program to issue more than one additional cache operation per access. This relates to how the program handles misalignment, or more generally multi-block references. Ideally the block size could be completely independent of word size, and the cache could generate as many block accesses as necessary for any given request. The current implementation does not support this, although to add such functionality would be relatively trivial.

In terms of implementation details, there are several instances where expensive iteration is used unnecessarily. For example, to check if the cache is full requires iteration through the entire array of CacheSets, checking their 'isFull' flag. A much better way would be to count how many are full, which requires only a single integer comparison for the same effect.

Another notable set of "problems" (really just possible avenues for future work) is the missing features that would enrich the program. Multiple levels of cache is one clear such suggestion, as well as support for non-unified caches. Furthermore, support for multiport access, pipelining and prefetching would be desirable. Also of use might be statistics regarding the bandwidth used to and from the cache to service all the requests (as dinerolV provides, albeit in a limited capacity).

And beyond that, on the completely wishful-thinking list, is pervasive support for timing. This includes calculating statistical data such as cycles spent waiting on misses, average request latency, etc. Such functionality is not trivial to add, as the timing of a real world cache is far from a simple "1 cycle hits, 50 cycle misses". Indeed, a much more modular implementation would be required, as the timing simulation itself could easily swamp the core cache simulation.

Sample Simulation

To demonstrate the use of the simulator, the following section documents the output from a simple test set of data. The sample access pattern contains 100,000 requests, which are a mixture of instruction reads, data reads and data writes. There are no advanced requests (such as cache line invalidates/flushes, or prefetches). Due to its size, it is not included in the printed form of this document - it should however be bundled with the electronic copy.

To be more specific, the sample access pattern used contains 7,261 unique addresses, which with 64-byte blocks (as used extensively in this document) map to 129 unique blocks. There are 15,691 data read requests, 12,558 data write requests, and 71,751 instruction read requests.

Method

To gather data for this analysis, the program was run with the command:

```
./SimpleCacheSim -A 16 -M block -a 0-16 -b 1,2,4,8,16,32,64,128,256,512,1024 -c
1,2,4,8,16,32,64,128,256,512,1024,2048,4096,8192,16384,32768,65536,131072 -r ran-
dom,LRU,FIFO,LRR,LRW -v -1 -w 1-16 ../../smalltex.din > results
```

That is - all associativities from direct to 16-way and then full, as well as all the first 10 powers of two for block sizes, cache sizes in powers of two up to 128kiB, every replacement policy available and word sizes from 1 to 16 bytes.

This set of parameters generates 14,100 simulations, which are output to the “results” file. This was then copied into a PostgreSQL database (as defined using the --SQL/-S create parameter) so that it could be easily manipulated.

Tiny subsets of the results were selectively copied into Excel, in order to generate the simplified tables and diagrams shown throughout this document. Due to the functional limitations of Excel, the charts are not optimal, but do suffice for the most part.

Data

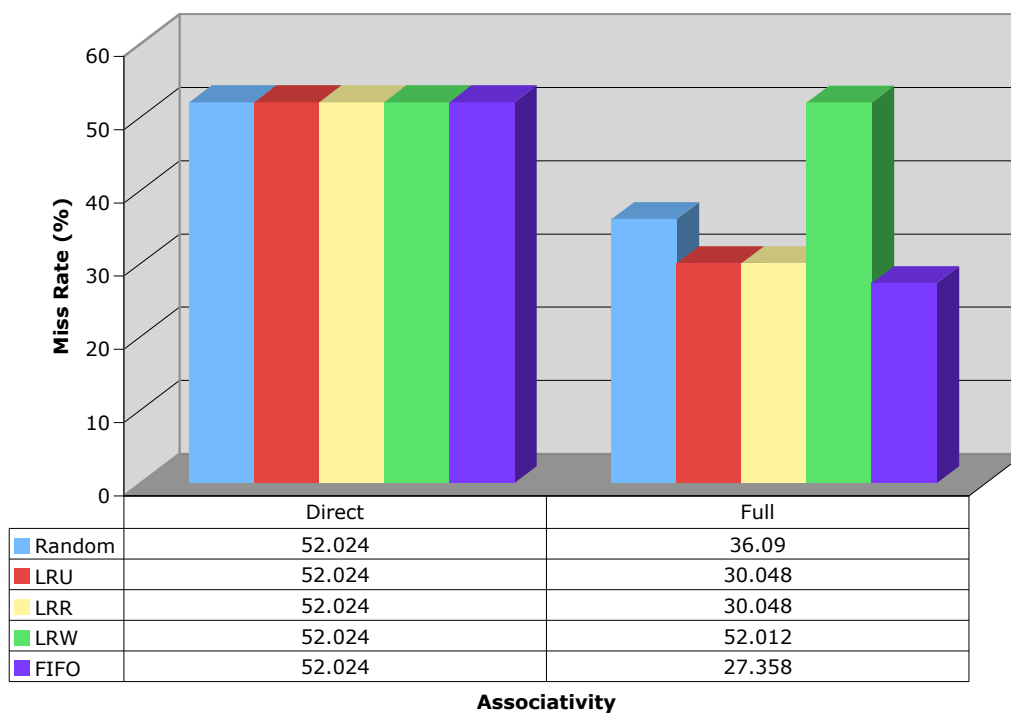
Appendix I contains the raw data used in this section, in tabular form. Italic entries indicate results from dinerolV. Entries coloured light yellow indicate discrepancies as a result of the randomness inherent in random replacement. Entries coloured blue indicate discrepancies as a result of errors in dinerolV.

Performance as related to cache size

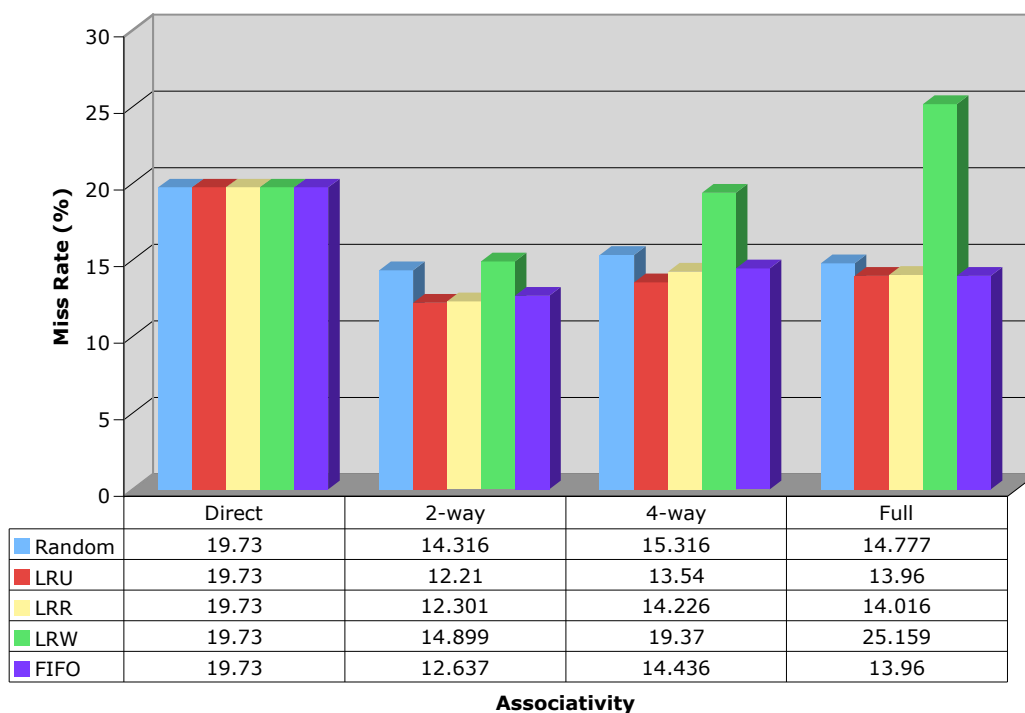
This section looks at the performance of caches of various sizes, using a fixed block size of 64 bytes (and a word size of one byte). The miss rates are shown in the following six charts. Keep in mind that the scales are not consistent between the charts, as the range of values varies significantly between them.

Total Miss Rates

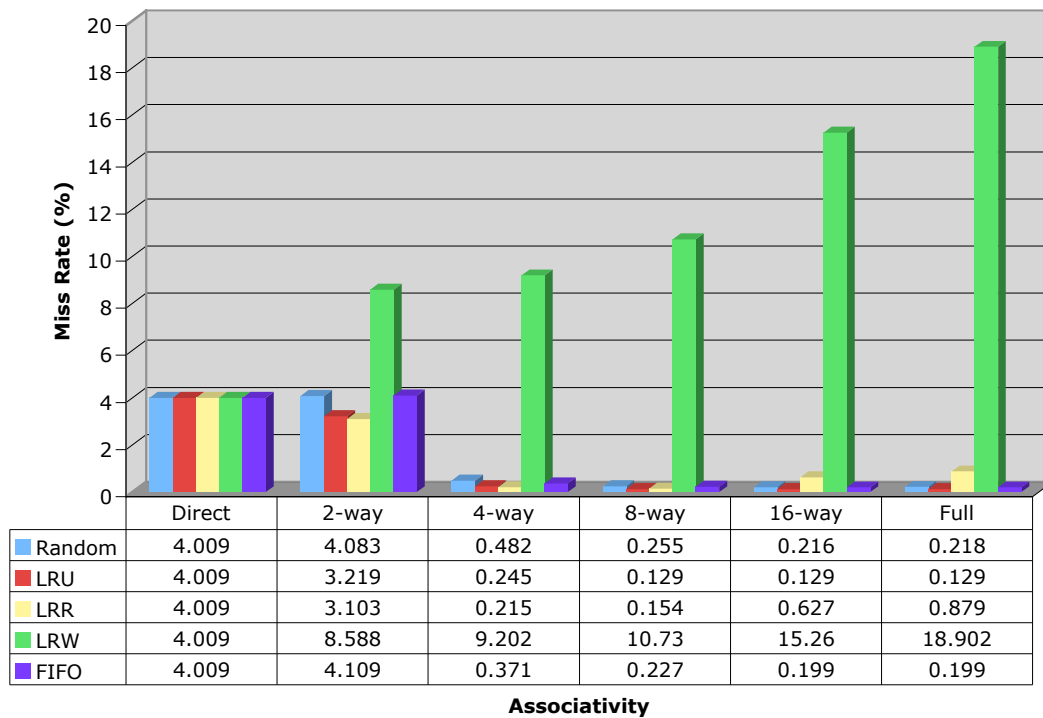
**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 128-byte cache, 1-byte words)**



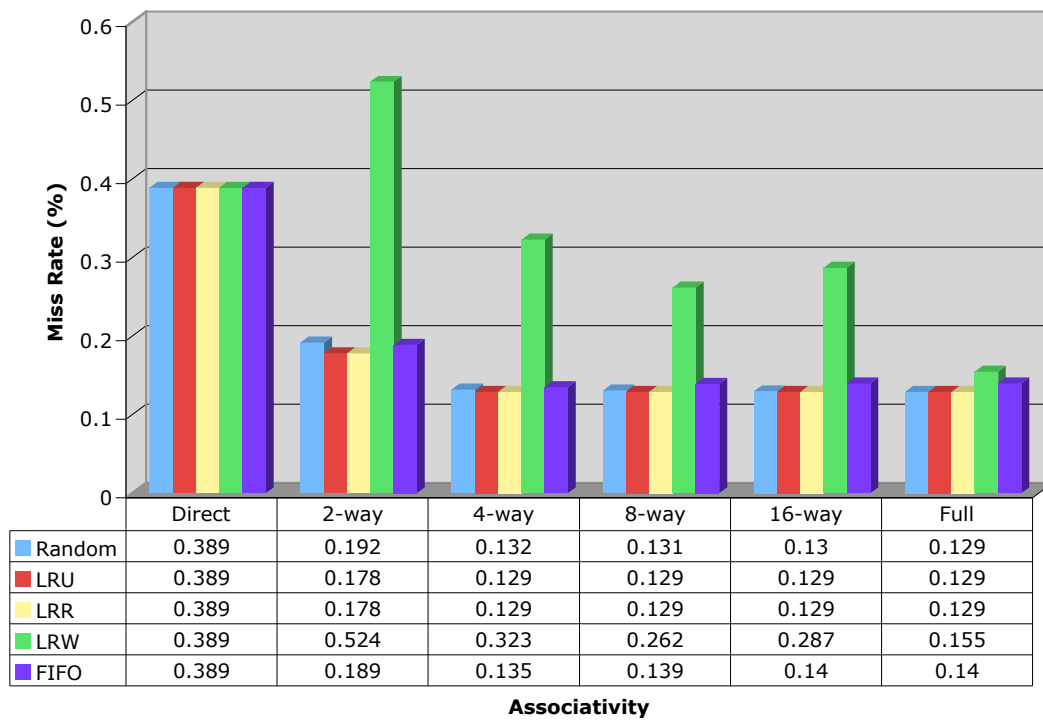
**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 512-byte cache, 1-byte words)**



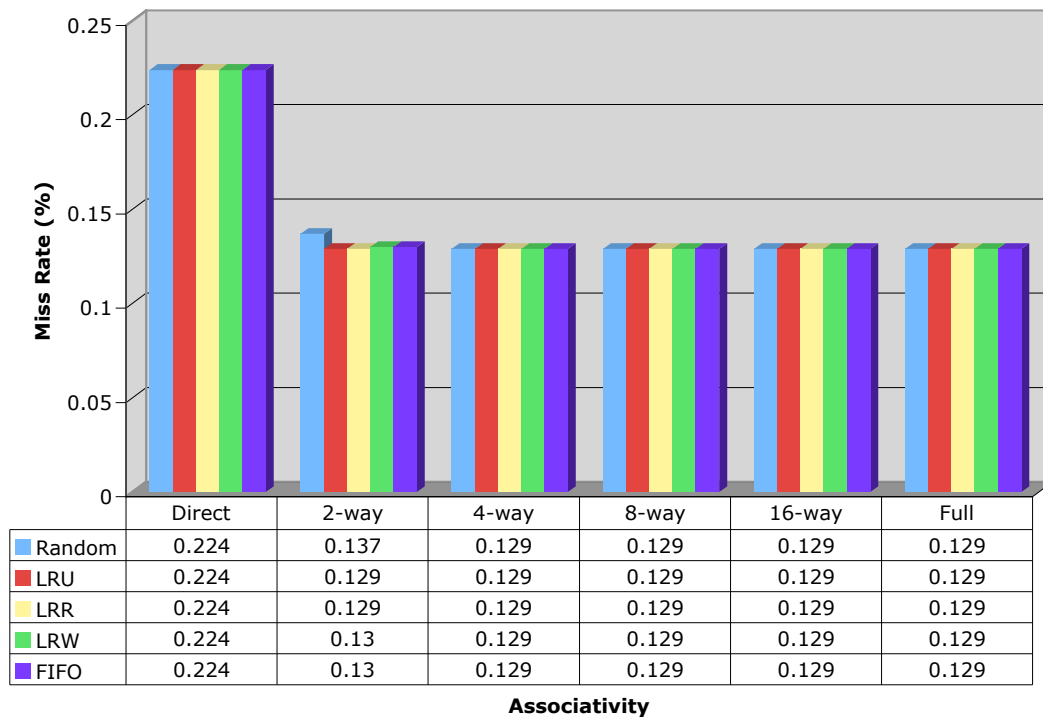
**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 2048-byte cache, 1-byte words)**



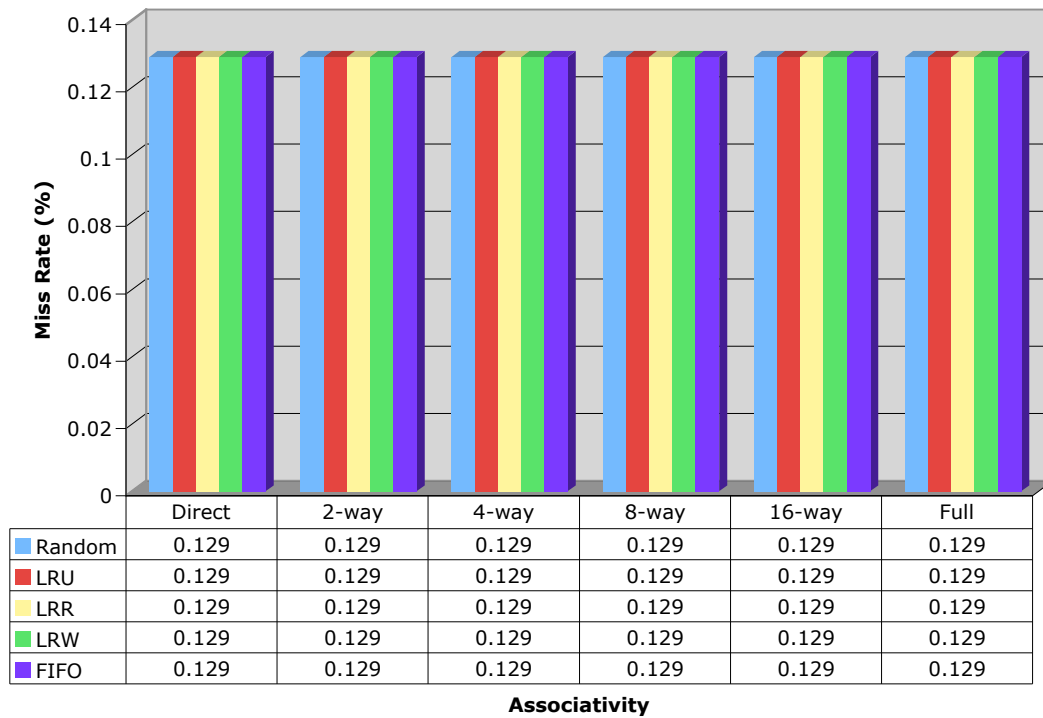
**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 8192-byte cache, 1-byte words)**



**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 16384-byte cache, 1-byte words)**



**Miss Rate over Associativity and Replacement Algorithm
(64-byte blocks, 32768-byte cache, 1-byte words)**



Note that at and above 32kiB (shown in the last chart, above) the results do not change at all - the sample access pattern is too small to have any further effect. All misses at this point are compulsory, with no evictions from cache ever occurring. Obviously a “serious” real program would perform several orders of magnitude more requests than we have simulated, which may require a far larger cache to reach optimal efficiency.

Discussion

There are several obvious conclusions to draw from these charts. The first is that, as noted previously, our sample set of requests requires only 32kiB of cache to achieve maximum efficiency. This obvious proof is intuitive, as with sufficient size the cache no longer acts as cache, but rather just a flat memory.

Another note is that even with a tiny cache - fitting just two blocks - there is still a marked improvement over no cache, with only slightly more than 50% of requests missing. More optimistically, this means nearly 50% of requests hit. Given the minimal latency a real cache would introduce, relative to the cost of going to the next level of memory, a 50% hit rate delivers fantastic performance gains.

Of interest is the varying performance of more esoteric replacement policies - particularly LRW (Least-Recently-Written). This policy performed terribly with small caches, particularly the 2kiB cache, where its miss rate was an order of magnitude greater than every other policy. However, for larger caches the performance improved again - although in the tests above never equalling the other policies. Clearly LRW is not a good policy for the sample access pattern used - intuitively this is true, as memory is more often read than written, and programs do not generally overwrite the same memory frequently, or in quick succession.

Overall the best policy was clearly LRU (Least-Recently-Used), which in nearly every case at least equalled the other policies. In only a handful of cases was it slightly outperformed by LRR (Least-Recently-Read), but these are so infrequent and the differences so small as to be insubstantial. But LRR was still a reasonable policy, and for some data sets could hypothetically be significantly better than LRU - e.g. when iteratively reading a small set of addresses, and writing a large set of varying addresses. In this case LRU might tend to throw out from the small set of read addresses, whereas LRR would favour keeping them, flushing instead the writes - which in this scenario would be unlikely to hit again frequently, if at all.

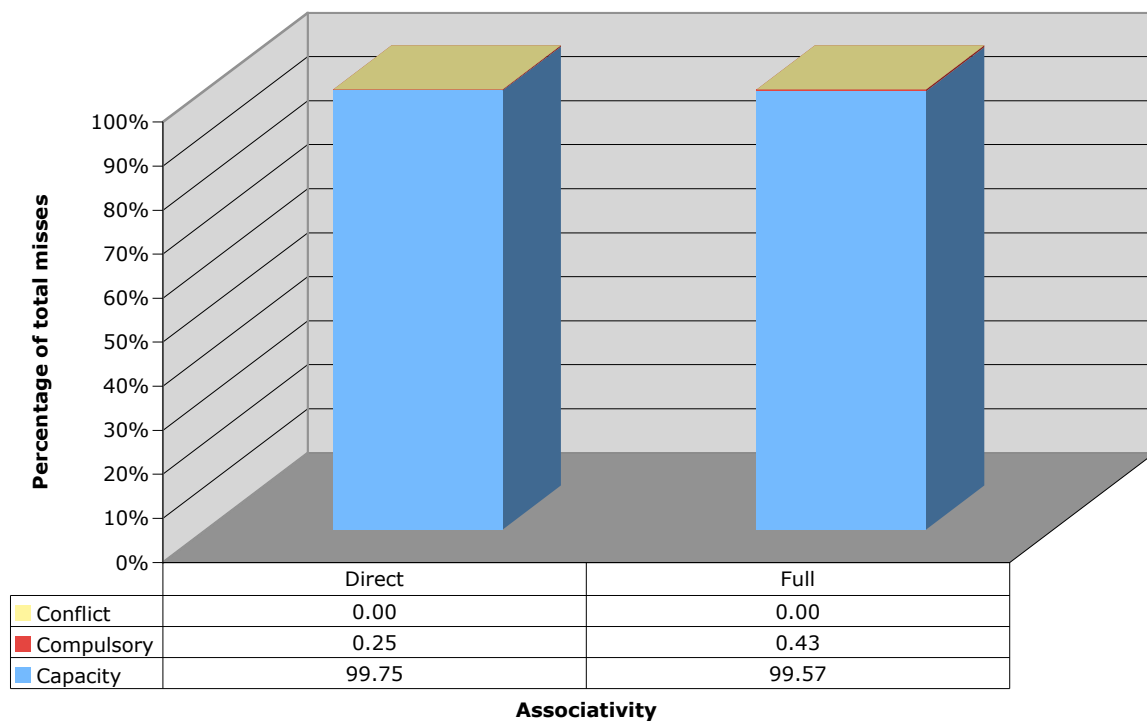
In terms of associativity, it is clear that higher associativities are generally superior. In nearly every case direct mapping resulted in twice as many (or more) misses as associative configurations. It is interesting to note, however, that in some cases 2-way associativity was superior to higher associativities - such as with the 512-byte cache. Yet in others (e.g. 2048 and 8192) it was clearly inferior.

By and large associativities of four and greater were equivalent. It is generally accepted that higher associativities represent strongly diminishing returns in real world caches - 16-way associativity is generally considered equivalent to full associativity, and 8-way associative caches are by far the most common in the real world.

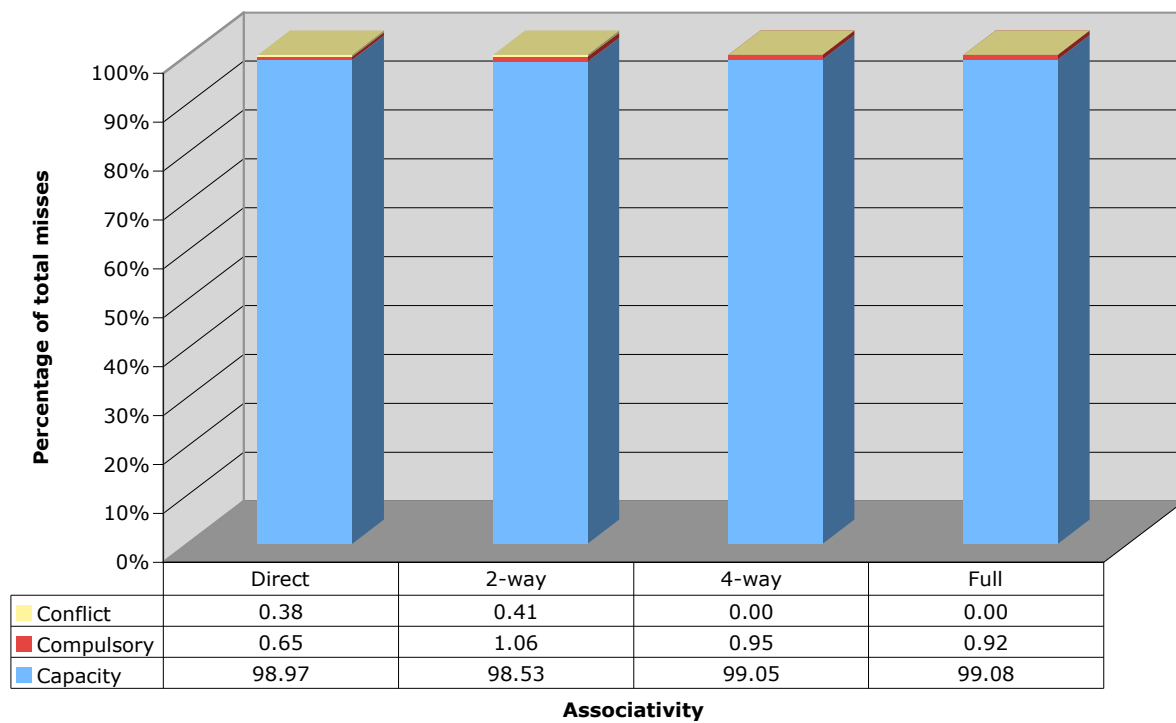
To be complete, it should finally be noted that random replacement, while potentially equivalent to any other policy - as well as potentially better or worse - did not perform particularly well in any of the tests. It was never substantially worse than LRU, but it was very rarely equivalent [or better]. Obviously rerunning the exact same simulations could produce different results, so the results reported above [for random replacement] are not conclusive.

Miss Breakdown

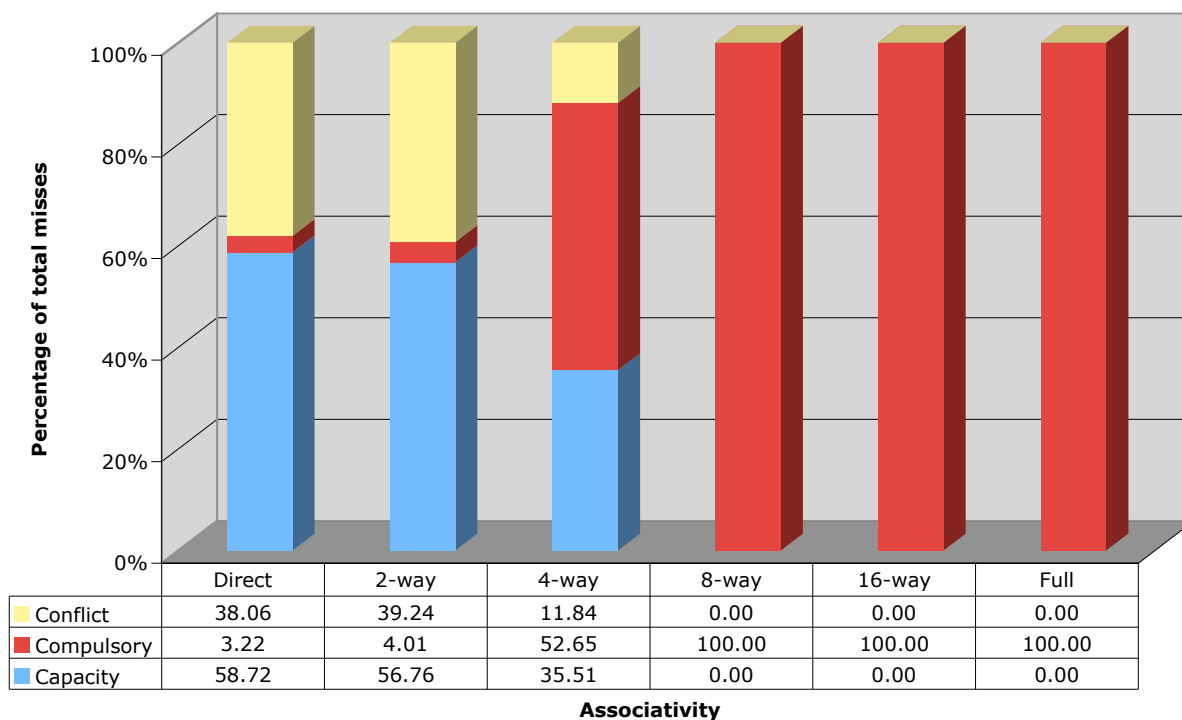
**Breakdown of Miss Types for LRU
(64-byte blocks, 128-byte cache, 1-byte words)**



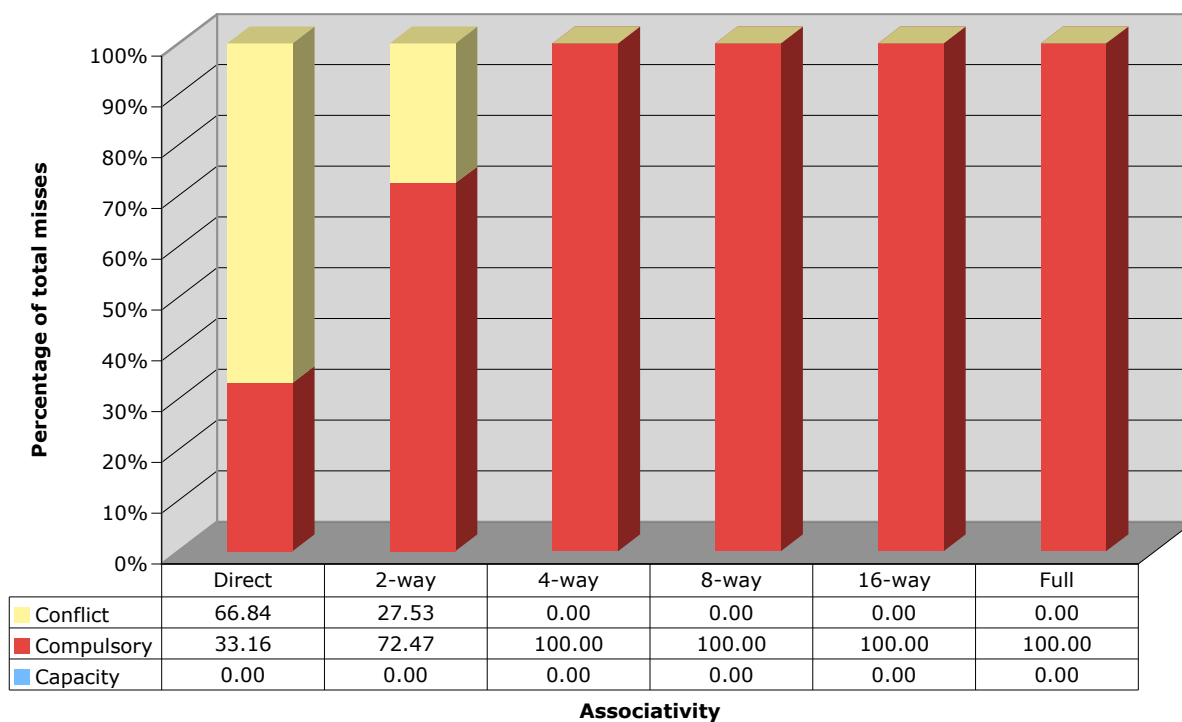
**Breakdown of Miss Types for LRU
(64-byte blocks, 512-byte cache, 1-byte words)**



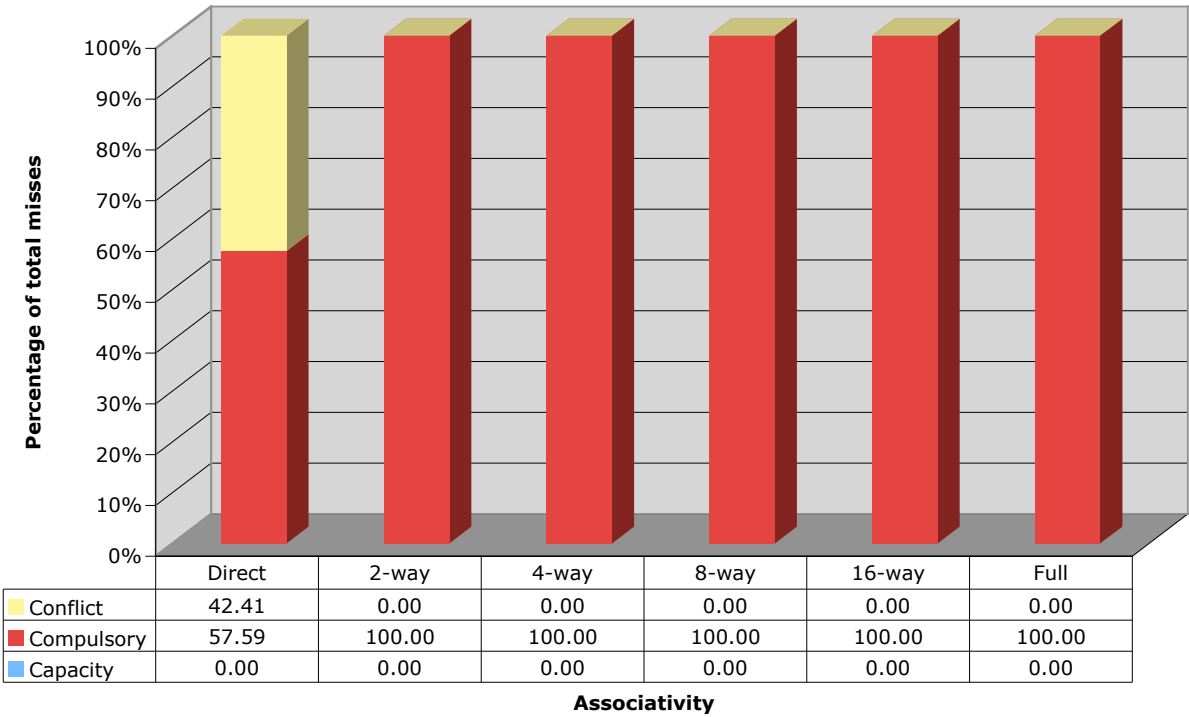
**Breakdown of Miss Types for LRU
(64-byte blocks, 2048-byte cache, 1-byte words)**



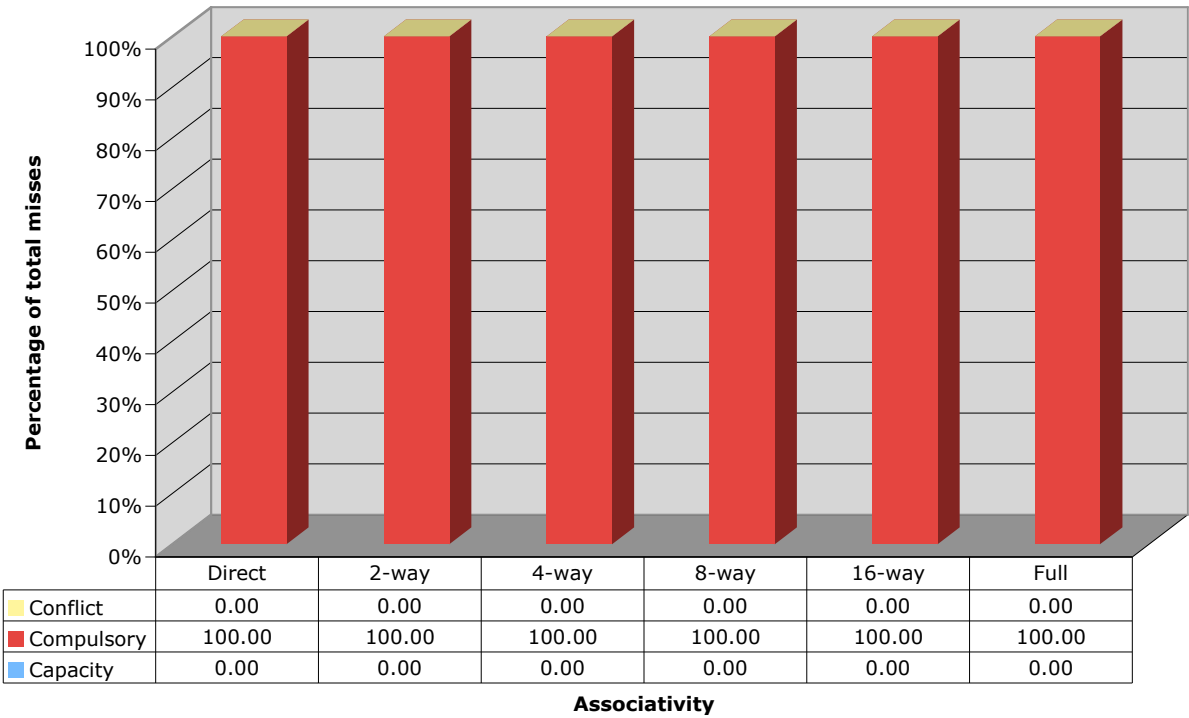
**Breakdown of Miss Types for LRU
(64-byte blocks, 8192-byte cache, 1-byte words)**



Breakdown of Miss Types for LRU
(64-byte blocks, 16384-byte cache, 1-byte)



Breakdown of Miss Types for LRU
(64-byte blocks, 32768-byte cache, 1-byte)



Discussion

Keep in mind when reviewing the charts that there are, in every case, 129 compulsory misses - since there are only 129 blocks (for 64-byte blocks) touched by the sample access pattern. Thus the size of the red segment can be used as a yardstick.

The data here is more simplistic, but shows strong conclusions. The most obvious is the transition from capacity misses as the dominant type, to compulsory misses, as cache size increases. Similarly, for caches “in the middle”, size-wise, conflict and capacity misses dominate only for the low associativities.

Overall it is clear that conflict misses are the dominant problem - both by nature and by portion of misses. A conflict miss indicates that there is room in the cache to store a new block, but the cache structure does not permit the exact block in question to be stored in any of the available locations. Such misses can occur for a cache of any size, given an appropriate access pattern. While capacity misses are relatively trivial to reduce - simply increase the cache size - conflict misses require an increase in associativity, which is a far less trivial modification.

From this we could hypothesise that policies favouring efficiency at lower associativities would be the most critical for real world caches. LRU and LRR both represent this in the results demonstrated so far.

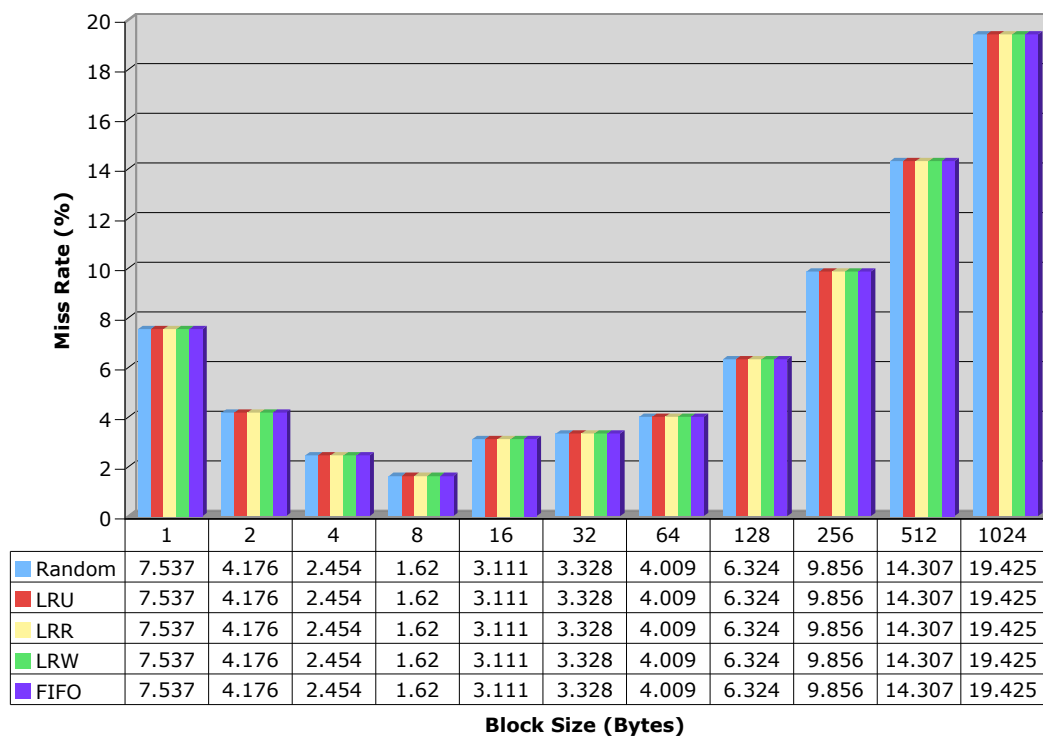
Performance as related to block size

This section looks at the performance of various replacement policies and block sizes, for a 2048-byte cache, with 1-byte words and considering associativities from direct to 16-way and then full. As in the preceding section, keep in mind that scales on the charts are not necessarily consistent.

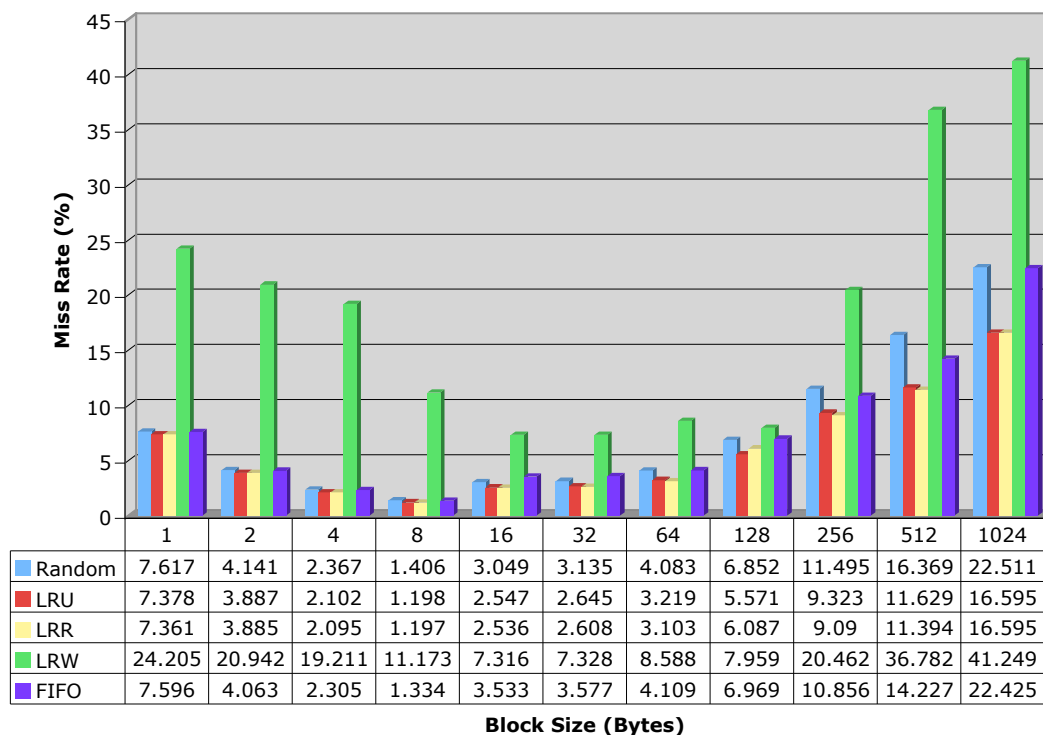
A cache size of 2048-bytes was chosen because in the previous results it seemed to represent a key point in the results, below which the performance of the cache was poor regardless of replacement policy, and above which it quickly became close to ideal for all policies. Thus 2048-bytes was chosen to ensure some interesting contention.

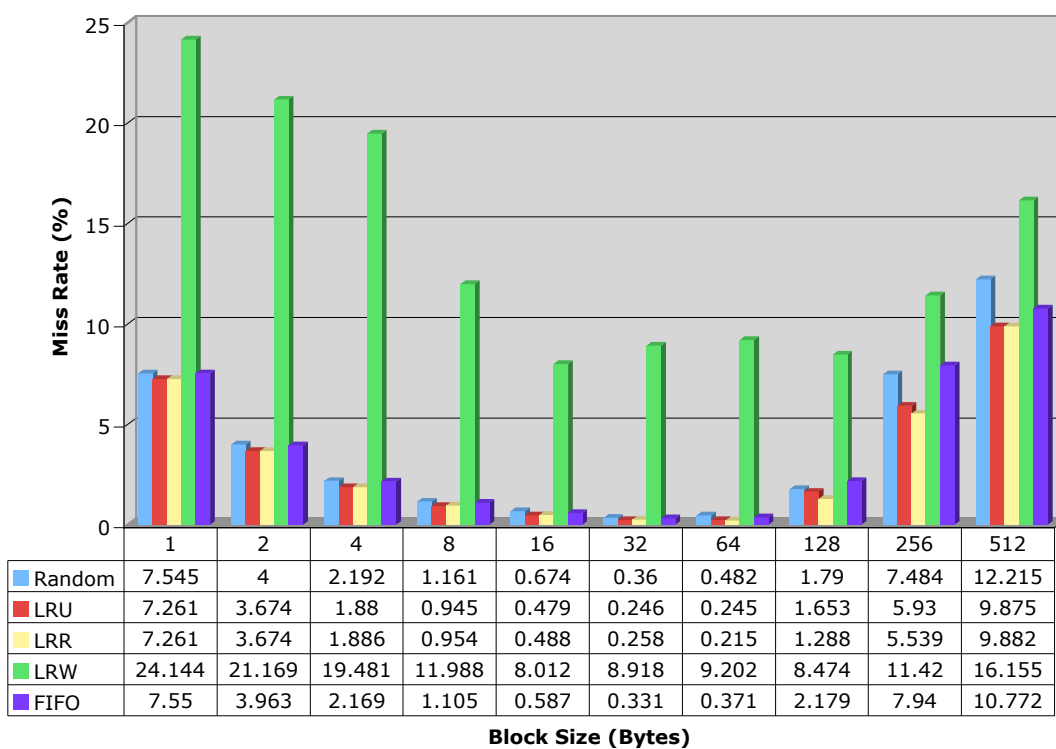
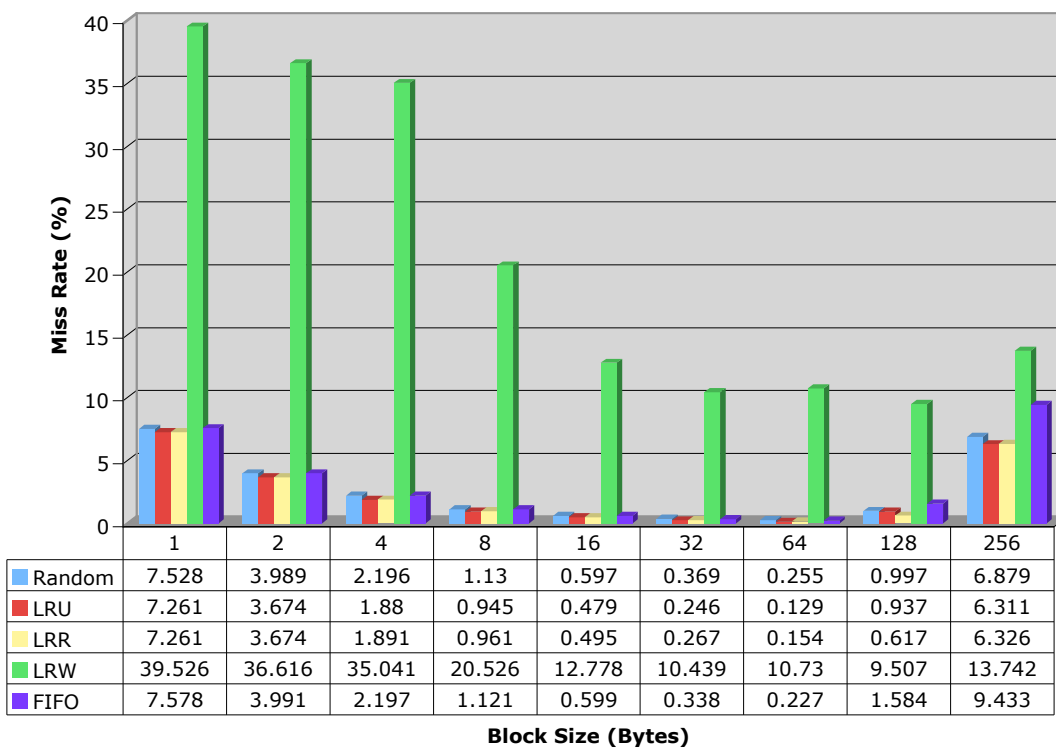
Total Miss Rates

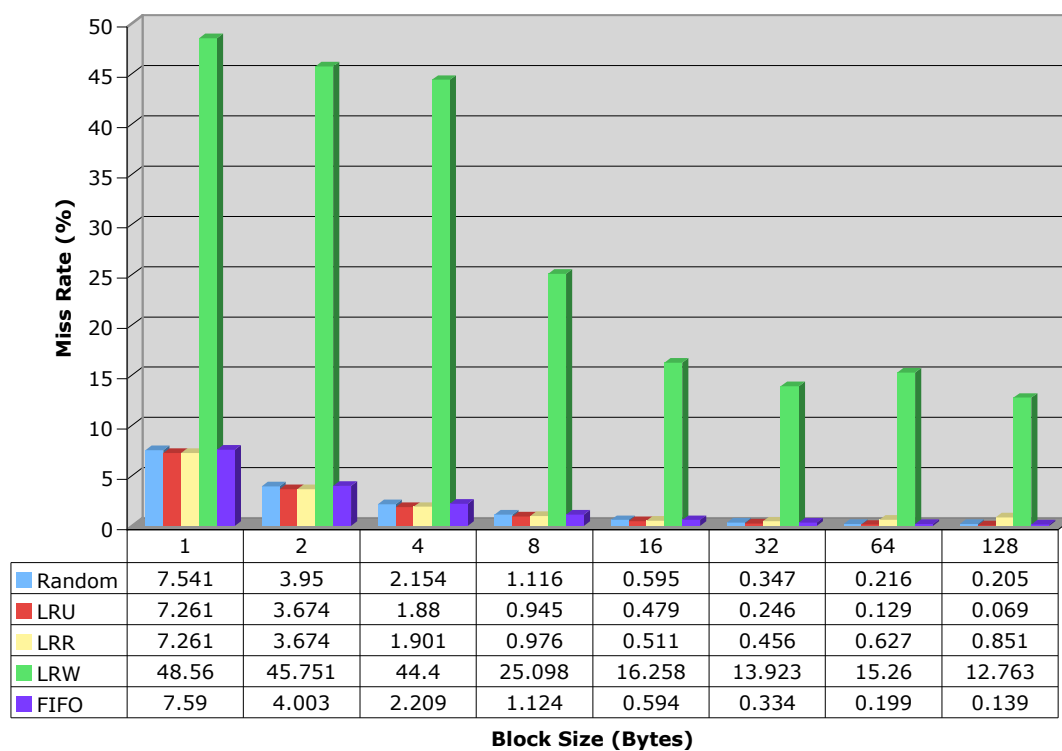
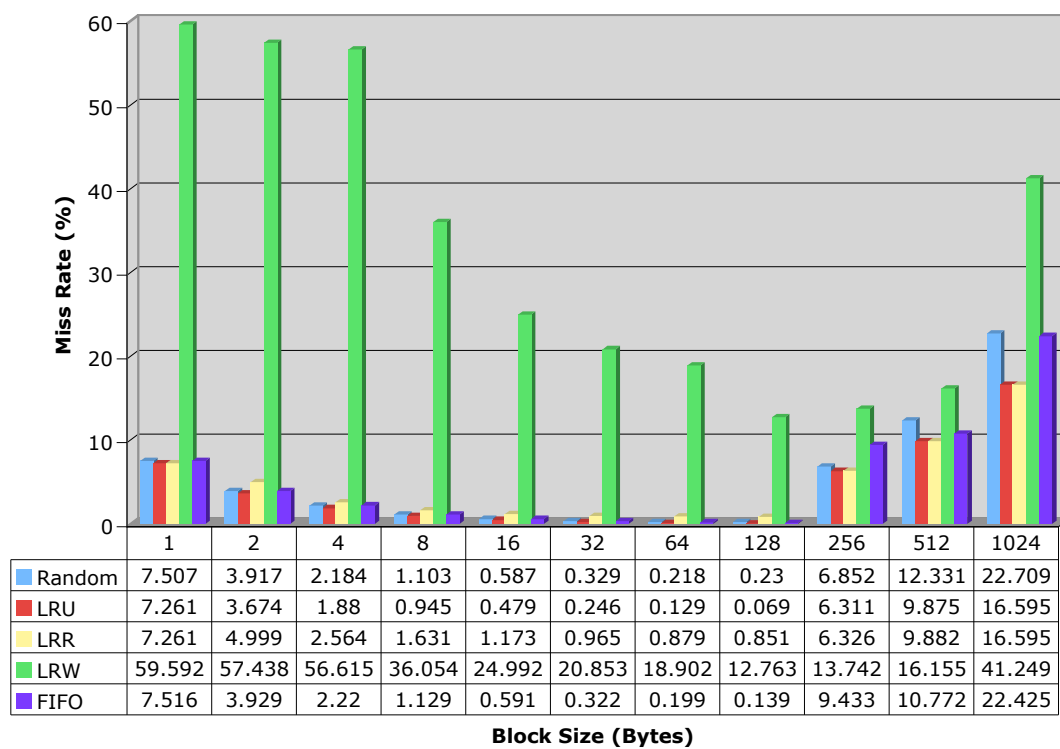
Miss rate versus block size (2048-byte cache, direct mapped, 1-byte words)



Miss rate versus block size (2048-byte cache, 2-way associative, 1-byte words)



Miss rate versus block size (2048-byte cache, 4-way associative, 1-byte words)**Miss rate versus block size (2048-byte cache, 8-way associative, 1-byte words)**

Miss rate versus block size (2048-byte cache, 16-way associative, 1-byte words)**Miss rate versus block size (2048-byte cache, fully associative, 1-byte words)**

Note: Obviously the replacement policy has no bearing for a direct mapped cache, as shown in the first chart above. Also keep in mind that the largest block size for each n-way associativity also corresponds to full associativity, so there is some overlap between the fully associative chart and the others.

Discussion

It must also be kept in mind that although larger blocks may result in fewer compulsory misses (as shall be detailed in the next subsection), larger blocks take longer to fetch. So the above charts may not directly represent performance (as measured by latency or memory bandwidth).

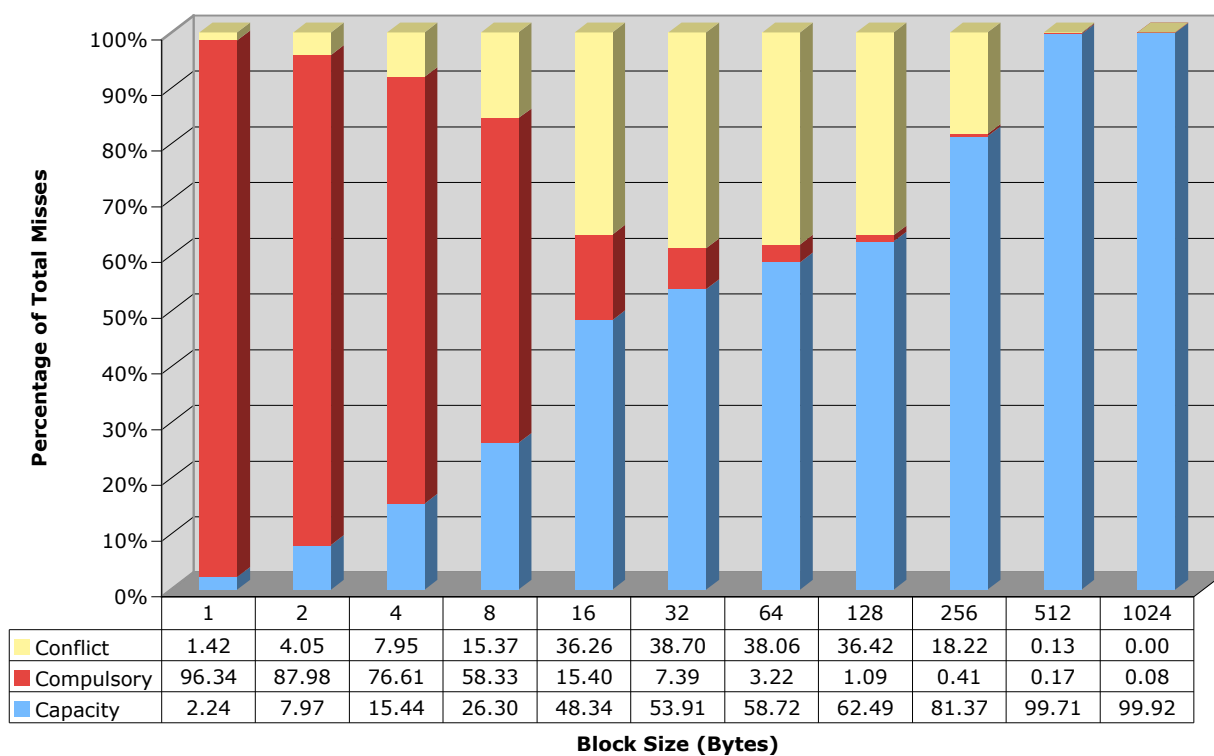
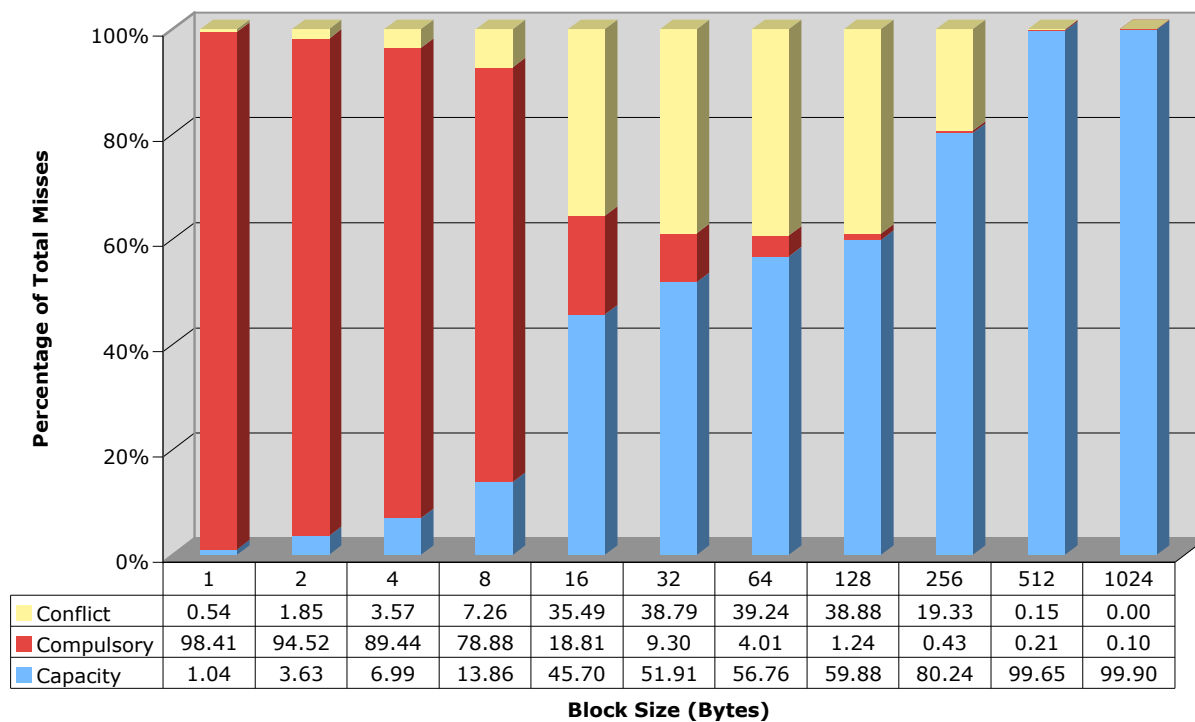
The first thing that is painfully obvious is that LRW again performs terribly. Although interestingly it still follows the same sort of curve as the other policies, but has a “sweet spot” consistently above the other policies - e.g. for the fully associative cache all other policies seem to favour block sizes of around 64 bytes, while LRW favours 256 bytes.

It is interesting to note that performance is much more adversely effected by too-large block sizes, as opposed to too-small block sizes. This is highly dependent on the exact request pattern issued to the cache, but is fairly intuitive - assuming every address is requested more than once, the optimal cache configuration will tend to be one which allows every required address to be in cache while maintaining the largest block size possible (to ensure as many addresses are covered by a single block as possible, thus reducing compulsory misses).

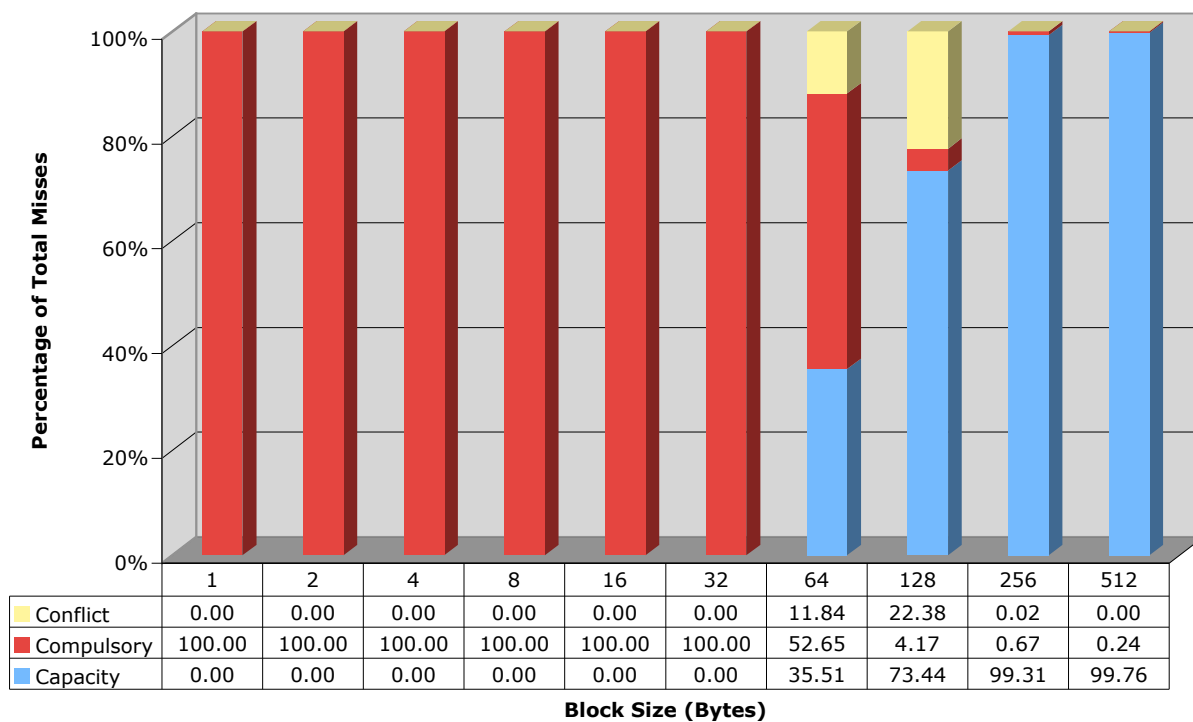
Although it is not entirely supported by the charts above, it seems likely that higher associativities do perform better, and favour larger block sizes, provided the cache is large enough to accommodate a reasonable number of sets. Looking at the fully associative results, as soon as the effective associativity drops below 8-way (128-byte blocks) performance falls off rapidly. If the cache were larger, this would not happen.

Miss Breakdown

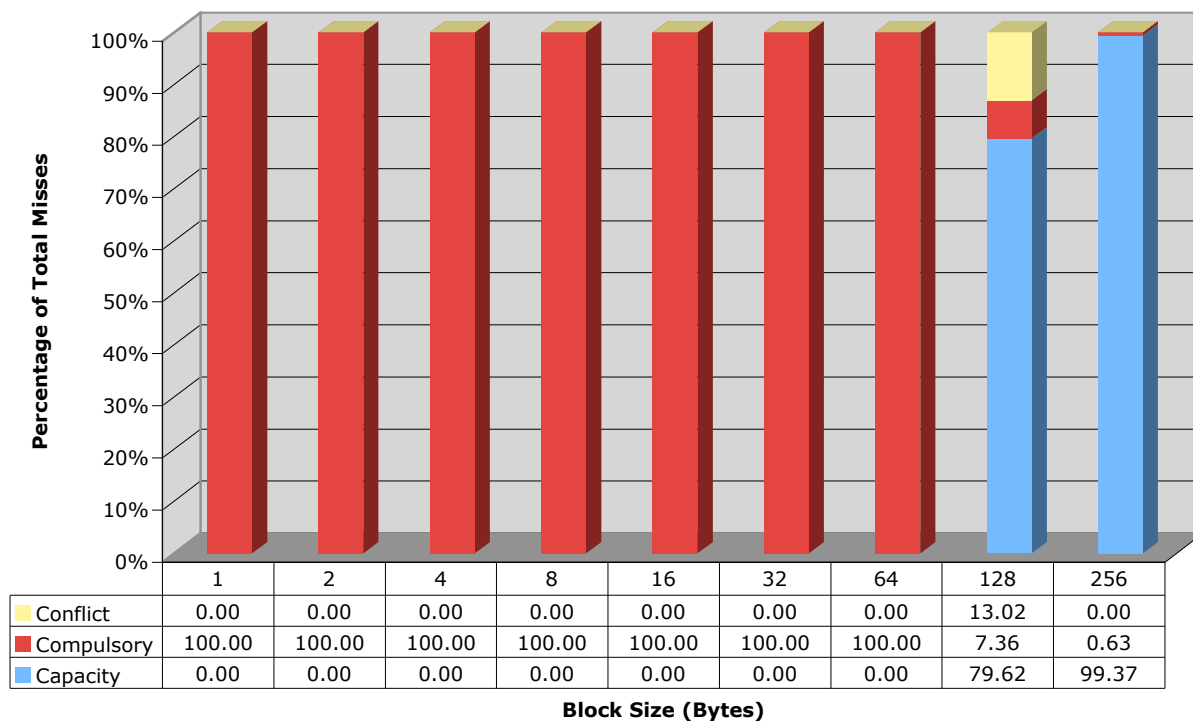
The following charts break down the misses shown for the LRU replacement policy in the previous charts, into capacity, compulsory and conflict misses. All scales are percentages and not directly comparable, as the absolute number of misses does vary significantly between charts and over block size.

Breakdown of Miss Types for LRU (2048-byte cache, direct mapped, 1-byte words)**Breakdown of Miss Types for LRU
(2048-byte cache, 2-way associative, 1-byte words)**

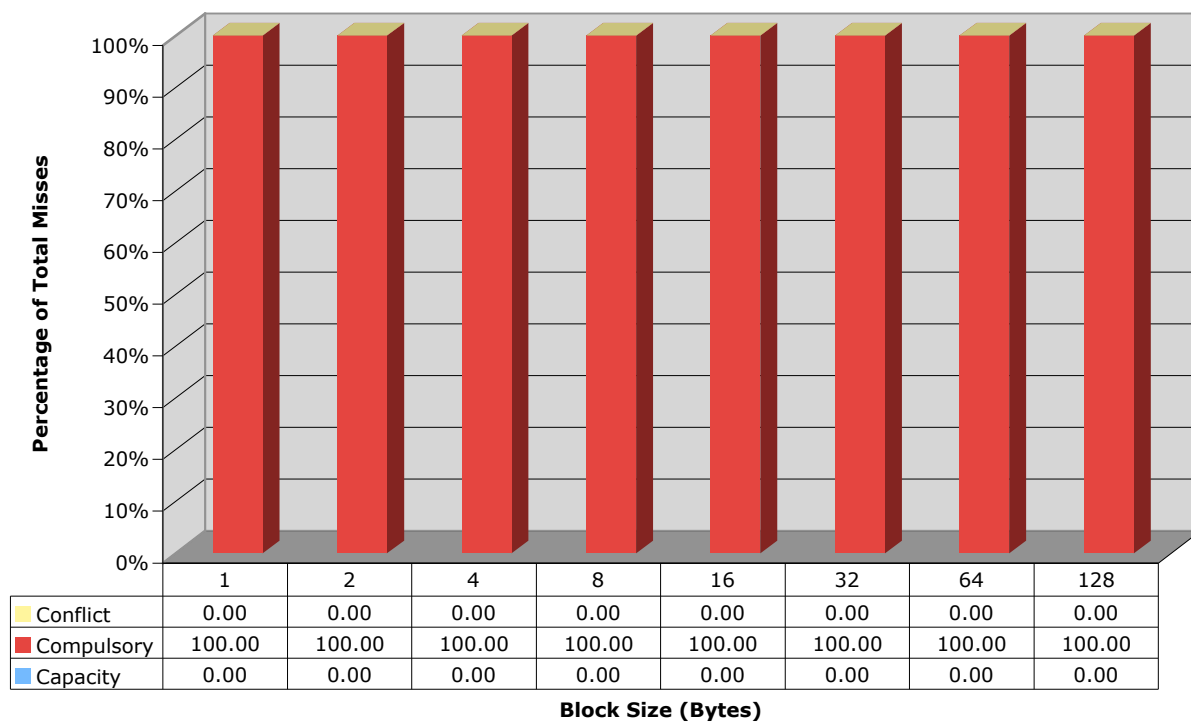
Breakdown of Miss Types for LRU
(2048-byte cache, 4-way associative, 1-byte words)



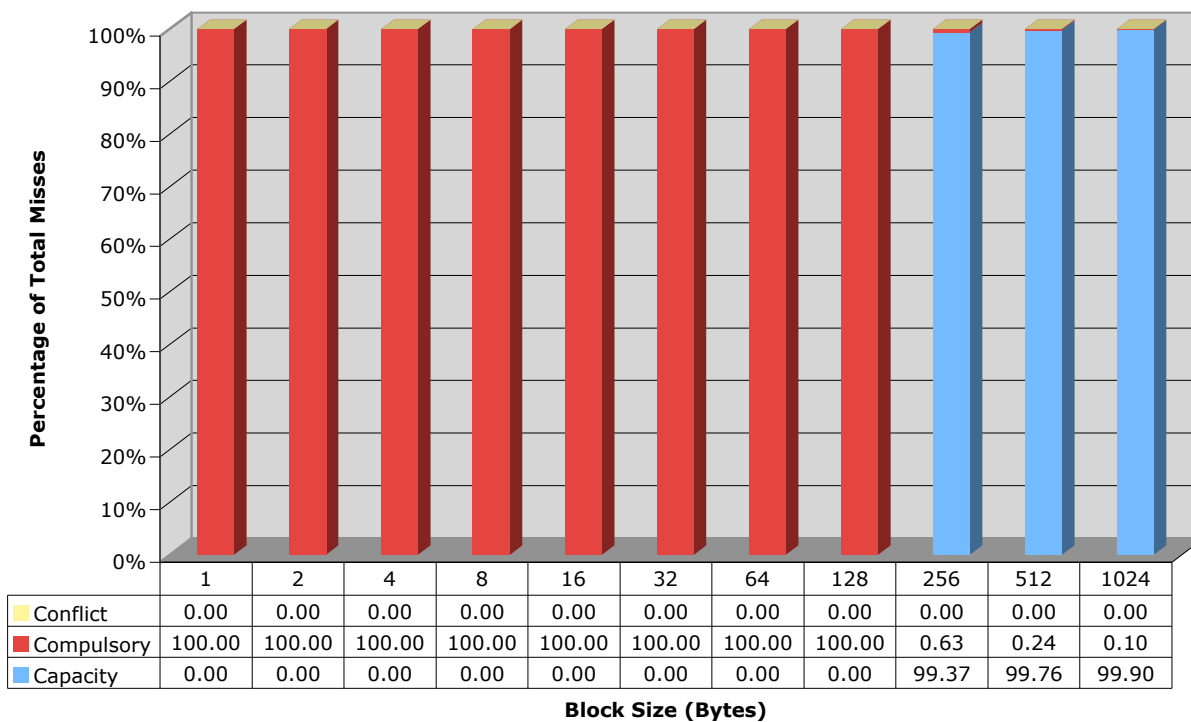
Breakdown of Miss Types for LRU
(2048-byte cache, 8-way associative, 1-byte words)



Breakdown of Miss Types for LRU
(2048-byte cache, 16-way associative, 1-byte words)



Breakdown of Miss Types for LRU
(2048-byte cache, fully associative, 1-byte words)



Discussion

Compulsory misses clearly dominate the lower block sizes, regardless of associativity - indeed, increasing the associativity only increases [slightly] the maximum block size before conflict or capacity misses become significant.

However, the results above may be misleading on their own - remember that the actual total number of misses was generally much higher for very low block sizes (4 or less), as shown in the previous subsection. This is because for very small block sizes, there is very little sharing of blocks by multiple addresses, consequently more unique blocks are required and there are vastly more compulsory misses to load these. For example, remember that a block size of 64-bytes requires 129 compulsory misses, while a block size of 1 byte requires 7,261 compulsory misses.

Also of interest is the fact that, for any real degree of associativity (i.e. not direct mapped), capacity misses dominate over conflict misses. This is relatively good, as capacity misses - as previously discussed - can be easily reduced by increasing cache size, while conflict misses are a product of cache structure (associativity), which is much harder to modify.

Going from direct to 2-way associative, it is interesting to note that the portions are fairly similar. What does change slightly is the total number of misses - as seen in the previous subsection, and indirectly visible on these charts by the larger relative portion of compulsory misses. This is indicating that the extra associativity, even just 2-way, significantly decreases the number of “unnecessary” misses - i.e. conflict and capacity misses. Unfortunately, however, once block size becomes large enough (16 in this case) performance drops off similarly. A large part of this is capacity misses, due to the reduced number of sets.

Ideal Cache

In summary, the ideal approach for designing a cache is:

- Large size
- High associativity
- Reasonably large block size

It is important to recognise at this point the limitations of SimpleCacheSim. Notably, that it does not simulate timing of any sort. All the points above tend to push latencies up and clock speed down, such that there is diminishing returns to increasing them, and at some point performance will in fact begin to decrease. And then even beyond that, there are issues such as pipelining, multiport access (concurrent requests), etc. And once you get to issues such as manufacturing yields and recovery, it's almost anybody's guess.

Nonetheless, given the results presented in this document, it seems a block size of around 32 or 64 bytes is appropriate. A cache size of at least 16kiB is highly desirable - especially 32kiB or above, where there are no conflict or capacity misses. And an associativity of 8 or higher produces best results (although of course with a large enough cache, direct mapping will perform much the same as higher associativities, and would be much simpler and faster in a real implementation).

Looking through the entire 14,100 sets of results produced for this document, the lowest possible number of misses is 16. This occurs only for a block size of 1024 (the largest simulated), for caches at least 16kiB in size (given high associativity - larger caches required for lower associativities), and for any replacement policy (since the 16 misses are compulsory, and there are no evictions). A full listing is included in Appendix II.

In the real world, the typical cache configurations are as shown below. These figures are derived from the specifications for the PowerPC 7400 and 970, the AMD Athlon, the Intel Pentium IV and Pentium M.

	Cache Size	Word Size	Block Size	Associativity
L1	32-128kiB Split	16B-32B/4*4B-4*8B	32-128B	2-8 way
L2	256-1024kiB Unified	8B-32B/4*2B-4*8B	64-128B	8 way
L3	512-2048kiB Unified	8B-16B/4*2B-4*4B	?	?

Note that some of the caches are multiported, allowing multiple requests at different addresses. Thus some L1 caches have an effective word size of 32B, composed of 8 independent requests for standard 4-byte words.

It is worth noting that in the leading edge processors today - notably the Pentium M and future processors from Intel - the block size is decreasing at all levels of cache, down to 32 bytes in the Pentium M (from 128 bytes in the Pentium IV). Smaller block sizes offer better performance when serving multiple threads concurrently, each of which will have a very disparate memory pattern - and parallelism is a high priority for future processors. Combined with modern prefetching (the 970 has eight independent prefetch streams), many compulsory misses can be avoided, reducing or even eliminating that key failing of small block sizes.

One day we may see huge, fully-associative caches with single-word block sizes and all the extra bells and whistles. For the moment, however, there are still some very big trade-offs to be evaluated. And in any case, the performance of the cache is largely dependent on the access pattern thrown at it, making it very difficult to design a cache that pleases everyone. It is strongly evident in modern processors that the caches were designed with very specific aims, aiming to be very good with certain types of programs (e.g. games) while performing miserably for everything else.

Of course, the next breakthrough may be memory which performs as fast as today's L1 or L2 caches, which will change the performance landscape massively.