

Accurate Data Distribution into Blocks may Boost Cache Performance

Dan Truong, François Bodin, André Seznec

► **To cite this version:**

Dan Truong, François Bodin, André Seznec. Accurate Data Distribution into Blocks may Boost Cache Performance. [Research Report] RR-3174, INRIA. 1997. <inria-00073514>

HAL Id: inria-00073514

<https://hal.inria.fr/inria-00073514>

Submitted on 24 May 2006

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

*Accurate data distribution into blocks may boost
cache performance*

Dan N. Truong, François Bodin, André Sezec

N° 3174

Mai 1997

_____ THÈME 1 _____



*apport
de recherche*



Accurate data distribution into blocks may boost cache performance

Dan N. Truong*, François Bodin†, André Seznec‡

Thème 1 — Réseaux et systèmes
Projet CAPS

Rapport de recherche n° 3174 — Mai 1997 — 25 pages

Abstract: Applications often under-utilize cache space and there are no software locality optimization techniques available for non-scientific applications. We propose that data redistribution in memory be used to modify reference patterns to improve locality of references.

To understand the potential of such an approach and to explain where gains come from, we introduce distribution misses, and define a correlation metric to evaluate spatial locality.

Data distribution can help reduce capacity and conflict misses in regular caches, as our experimental results to show. We use as example a profile-based scalar data layout heuristic, which was able to remove up to 76% of the direct-mapped cache miss ratio on some benchmark traces.

Key-words: Cache, Distribution miss, Memory Reference Correlation, Data Layout, Replacement policy, Simulation, Performance, Optimization, Belady, Spatial locality

(Résumé : tsvp)

* dtruong@irisa.fr

† bodin@irisa.fr

‡ seznec@irisa.fr

Un placement adéquat des données en blocs pourrait améliorer les performances des antémémoires

Résumé : Les antémémoires ont été développées afin de réduire le nombre d'accès à la mémoire en s'appuyant sur la localité spatiale et temporelle des références mémoire. Les antémémoires se sont révélées très efficaces, bien qu'étant encore relativement mal maîtrisées: leur mise au point est rendue difficile par la diversité de comportement des applications. Afin de résoudre ce problème, les architectes ont tendance à utiliser des hiérarchies mémoire de plus en plus complexes, ce qui va à l'encontre des contraintes technologiques (temps de cycle, surface disponible, consommation) et économiques (coût des mémoires statiques). La capacité de stockage des antémémoires est souvent sous-utilisée car la localité des références n'est pas bien exploitée. Il n'existe, à ce jour, aucune solution globale qui permette d'optimiser leur utilisation.

Nous analysons en premier lieu les effets de la localité sur le comportement des antémémoires, et rappelons la classification des défauts survenant dans les antémémoires définies par Hill et affinée par Sugumar. Nous rappelons la politique optimale de remplacement des blocs qui génère un taux d'échecs minimal, "*Min*", décrite par Belady. Nous définissons ensuite la notion de défaut lié au **placement des données**. Le placement des données a un impact direct sur la localité des références mémoires: nous montrons qu'il est possible d'améliorer encore le taux d'échecs de *Min* par un placement judicieux des données dans les blocs. Nous mettons aussi en évidence la borne minimum du taux de défauts pour une taille d'antémémoire donnée.

Dans une seconde partie, nous analysons les possibilités d'amélioration du placement des données en blocs, et du placement des blocs en mémoire afin de minimiser les taux d'échecs des antémémoires à correspondance directe. Nous décrivons en premier un algorithme qui permet de trouver, s'il existe, le placement de données qui suit la politique de remplacement de Belady et qui fournit le taux d'échecs le plus faible possible. Cet algorithme a une complexité élevée et ne fournit pas toujours de solution. Nous introduisons les propriétés d'**Activité** d'une donnée et de **Corrélation** entre deux données. La mesure de ces propriétés permet de quantifier la localité des références. Nous décrivons une heuristique basée sur ces mesures pour placer les données en mémoire. Afin d'évaluer le potentiel qu'offre le placement des données, nous avons implémenté et expérimenté cette heuristique sur des traces de programmes C et Fortran. Les résultats montrent qu'en plaçant correctement les données, il est possible de réduire le taux d'échec de façon importante dans un grand nombre de cas, aussi bien pour les antémémoires à correspondance directe (réduction du taux d'échecs entre 25% et 76% sur nos traces) que pour les antémémoires associatives (jusqu'à 64%).

Mots-clé : Antémémoire, Défaut de distribution, Corrélation des références mémoires, Placement mémoire des données, Politique de remplacement, Simulation, Performance, Optimisation, Belady, Localité spatiale

1 Introduction

Modern non-scientific applications tend to under-exploit cache space [BGK95]. Working sets become so large that caches are not able to take advantage of locality with irregular reference patterns. Since cache sizes grow too slowly [BGK96], it is important to regularize the cache reference patterns. A software optimization approach that we call data distribution could be used to adapt applications to the memory hierarchy. It consists in modifying the address of data-elements in memory so that data used simultaneously fit together in the cache.

As a first step towards data-distribution solutions, we propose in this paper a formalism to explain the impact of data distribution. It explains in particular how data distribution can impact on capacity misses. We also present preliminary experimental results, to show that data distribution is a very promising approach for non-scientific applications optimization.

Regular array accesses benefit from efficient loop transformation algorithms. They are source code transformations which modify the reference patterns but do not modify the placement of the data in memory. This approach cannot be used for irregular memory references generated by non-array data: address computations are not index based. Therefore we consider a dual approach consisting in changing the relative position of data elements in memory to improve spatial locality. We call this data-layout optimization approach data distribution, in contrast to loop-transformation techniques.

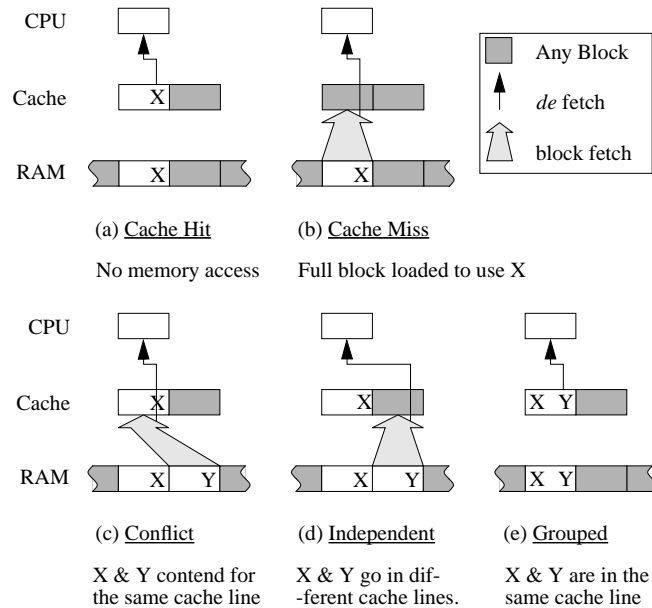


Figure 1: Locality impacts on cache behavior

Let's analyse how data must be laid out in the cache to improve locality¹ of references, by illustrating some simple layout alternatives, as shown in figure 1. Let's look in more detail into reasons of a hit or a miss (figure 1(a,b)): A cache manipulates blocks of memory: A reference to a data-element is a reference to the block holding the data element. Therefore, if we consider two data-elements X and Y with interleaved references, their distribution among the blocks impacts on the cache miss rate. Three possibilities arise.

- Figure 1(c): X and Y are stored in distinct memory blocks which compete for the same cache line. Each reference generates a conflict miss. Temporal locality is counter productive.
- Figure 1(d): X and Y are stored in distinct blocks mapping to different lines. Temporal locality is beneficial. A block-layout optimization must prevent (c) to get (d) to remove conflict misses.
- Figure 1(e): X and Y are stored in the same block. Temporal locality is beneficial, and only a single line is used: The cache benefits from spatial locality. This solution is the best: fewer blocks are used to handle references to these elements. Fewer blocks to store in the cache simultaneously mean fewer cold and capacity misses, and fewer chances of having conflict miss hazards.

With these simple remarks, we see the importance of good spatial locality and proper data-element distribution among blocks.

Data distribution optimization for set-associative caches must be a two steps process. Data elements used simultaneously must be redistributed among blocks to improve spatial locality, and blocks used simultaneously must be laid in memory to improve block temporal locality. Block layout uses the cache address mapping properties of set-associative caches to ensure that blocks used simultaneously load into different lines.

Standard array transformations attempt to improve spatial locality by making stride 1 accesses to arrays. Data-element distribution into blocks was not considered because array elements have a fixed position in the array. If elements are referenced directly, their relative position in memory does not impact on program correctness, so they can be redistributed in memory to improve spatial locality. Redistribution has just to conform to the few layout constraints imposed by compilers (keeping static, stack, heap allocated data respectively together).

In section 2 we present related work. We introduce data element block *distribution misses* in section 3 to account for the impact of data layout on the capacity miss ratio, and define a lower-bound for cache miss-ratio. In section 4 we provide two metrics *Activity* and *Correlation* to quantify locality of reference. We used these metrics to develop a profile based heuristic for data-elements distribution and block layout, and used it to evaluate the potential of such an approach. Our experimental results, shown in section 5, show a great

¹We say that a memory location **benefits from temporal locality** in the cache when it is accessed more than once while it is held in the cache. Furthermore, we say that a memory block **benefits from spatial locality** when many of its data-elements are accessed while it is held in the cache.

potential, with miss rates dropping by as much as 76%. We conclude in section 6 and propose some research directions to exploit this potential.

2 Related work

Many hardware solutions are available to address all aspects of cache performance [CB92, FJ94, DSW95, BS95, Jou90, AP93, JW94]. However, Burger et al. argue that applications often do not use more than 10% of the cache space, and a software managed local memory might be a better solution [BGK95]. However static compiler analysis alone may also be insufficient to take fully advantage of the local memory, and necessitates extra instruction overhead to handle it. We argue that with data distribution, it is possible to combine benefits of software management of local memories with the dynamic adaptation of caches to evolving program behaviour, while saving on the extra overhead incurred by explicit local memory management.

Much work has been done to analyze cache behavior, but most works measure evolution of the miss ratio for different cache parameters [Prz90, Smi82, Smi86]. To our knowledge, besides [BGK95], little work has been done to evaluate the the impact of data layout on locality, and measure the ratio of cache space effectively used. This paper proposes distribution misses to take into account spatial locality, and proposes data-distribution to address the problem. Grimsrud et al. express locality in a reference trace as $Addr(T + \delta_i) = Addr(T) + \Delta_{Addr}$, and uses this to provide a 3D graph, $reference = \mathcal{F}(stride, time)$, to view locality patterns [GAFN94], but do not consider layout optimization. The locality measures we provide in this paper are adapted to data-layout optimization.

Profile based instruction layout in memory to improve I-cache performance has already been explored in many papers [HC89, PH90, Hei94, CG94, TXD95]. The approach is similar in its principles to data-distribution, and efficient layout techniques have already been proposed. Pettis & Hansen [PH90] propose basic-block layout heuristics (they pack frequently used looping sequences) as well as procedure layout (they lay together procedures calling each other to minimize cache interference risks). They also prevent I-cache pollution by removing unused basic blocks from the main instruction stream. Profile-based instruction layout provides about 10% overall speedup by removing I-cache misses, virtual page faults and improving branch prediction. With proper training sets (they use profiling), speedups can be held for any inputs. Instructions are a special case of irregularly referenced data (for the I-cache) and we believe that efficient layout optimization can also be achieved for data caches.

To hide memory latency and to keep the CPU busy, non-blocking caches and software [CKP91] or hardware [Jou90, BC91, FPJ92, PK94] data prefetching techniques have been proposed. Prefetching tends to consume extra memory bandwidth, which in many applications is a critical resource, while non-blocking caches consume register space. Data layout optimization can reduce distribution and conflict misses, reducing at the same time the number of prefetches or pending loads necessary. Furthermore, the locality metrics we propose can be used by software prefetching techniques to select which blocks to prefetch.

Scientific programs have benefited from much more extensive analysis than non-scientific codes. Software array access optimization research has been very active for some years. With loop transformations like tiling [BGS94, CKT94], it is possible to prevent most interference misses [TGJ93] on regular array access patterns. These techniques are sometimes insufficient because interferences can still appear for some array leading sizes [BS95]. Modifying access patterns is not sufficient to eliminate interferences. It can be combined with array data layout optimizations, array leading-size padding and inter-array padding. Padding has been used to align arrays on multiprocessors [GMB95, TLH90, Por89] but can also be used to align arrays in caches [BCJ94]. For non-scientific applications, however, such techniques are irrelevant since complex data structures are not accessed using indexes. The orthogonal approach, data-distribution is better suited.

Although scalars benefit from efficient register allocation techniques [ASU86, CCK90], and new dynamic allocation libraries take advantage of temporal locality by reallocating recently freed blocks [GZH93, BZ93], this is still insufficient to reap full processing power from a modern CPU with non-scientific applications, and new approaches must be developed to improve memory hierarchy performance. This paper attempts to show that data-distribution is a very promising approach for non-scientific code optimization that will have to be investigated thoroughly.

3 Characterizing cache use

We recall the Three Cs miss classification [Hil87] using Sugumar's definition which use a perfect² cache [Sug93] and update it to handle data distribution. Data-distribution impacts not only on block layout in memory, but also on the distribution of data-elements among the memory blocks. This changes block reference patterns, and impacts on the capacity miss ratio. Therefore we provide an extension to the classification, data elements block *distribution misses*. They account for sub-optimal data-element distribution, just like conflict misses account for sub-optimal block replacements. This allows us to evaluate the lower miss ratios boundary attainable with a perfect data distribution and a perfect cache of given cache and line size.

3.1 Sugumar and Hill's classifications

To analyze cache hardware block replacement policies, Hill classified cache misses. We use Sugumar's definitions which are finer, and recall them since they are not as widely used.

Definition 1 *A cold miss or compulsory miss is the first load of a block in the cache, due to the first reference to it since the beginning of the program.*

This condition looks at block being loaded. Subsequent characterizations look at the block replaced.

²We call perfect a fully-associative cache using Belady's optimal *Min* block replacement policy [Bel66]

Definition 2 A memory block is **dead** at a given time if it will never be referenced by the program again, otherwise it is **live**.

Belady [Bel66] has shown that there exists at least one optimal block replacement policy denoted *Min*, which yields the smallest miss ratio possible for a fully-associative cache. On a miss, *Min* replaces any dead cache block, or if there are none, the live block farthest referenced into the future. Therefore, a fully-associative cache using *Min* policy is perfect.

Definition 3 A **capacity miss** is the replacement of the block in the cache that will be referenced the farthest in time³.

Definition 4 A **conflict miss**⁴ or **collision miss** is the replacement of a live block which is not the farthest one referenced in time.

A suboptimal block replacement can happen with suboptimal hardware replacement policies like LRU, or with address mapping. To differentiate these two cases in *set-associative caches*, Sugumar [Sug93] further differentiates conflict misses:

Definition 5 A **Mapping miss** is the replacement of the farthest referenced block in a set, while the farthest referenced block in the cache is not in that set.

Definition 6 A **Replacement miss** is the replacement of a block that is not the farthest one referenced in the set.

A mapping miss is a loss of efficiency because the cache is split into sets, replacement misses are due to a bad replacement policy. All conflict misses in a *direct-mapped cache* are mapping misses, except when cache bypass is allowed: *Min* policy could choose not to load the block in the cache but few processors actually implement a bypass instruction.

The above classification only considers cache blocks, and mapping misses account only for block layout in memory. With data-layout optimization, it is necessary to consider the impact on cache performance of data element distribution into blocks.

3.2 Distribution misses

If we redistribute the data-elements among the different blocks, the reference patterns to the memory blocks will change. Therefore, *Min* generates different replacement choices, and the miss ratio of a perfect cache can change. We define *distribution misses* to account for this and compute a lower bound to the miss ratio of a perfect cache of a given size.

A data-layout is the address mapping of data elements in memory. Shuffling addresses of data elements modifies the data layout. Of all the possible data layouts, there is at least

³Hill's original classification defines as a capacity miss a miss occurring in a fully-associative cache of the same size with the same replacement policy as the one studied – usually LRU. However, Sugumar showed that this definition is too approximative, since it can yield negative conflict miss ratios [Sug93] and using a perfect cache is a better alternative.

⁴Sometimes also called an **interference miss**. This is misleading since cross or self interference misses usually denote conflict and capacity misses within loop nests which do not fully exploit temporal locality of references to an array.

one that generates the lowest miss ratio for a perfect cache (the number of permutations is finite).

Definition 7 *An optimal Data-layout is a data-layout for which a perfect cache yields the lowest miss ratio possible.*

Note: the address offset of a data element within its block does not change block reference patterns, so it does not affect the miss ratio – assuming data elements do not overlap two blocks. The offset within the block will therefore not be considered.

We propose a new miss characterization, which takes into account how data-elements are distributed among the blocks, similar to the definition of conflict misses.

Definition 8 *A data-element block distribution miss is a capacity miss that would not have happened with an optimal data-element distribution into memory blocks.*

Looking back at figure 1(c,d,e), we see that misses generated by (c,d) can be prevented by changing the data-elements distribution to get case (e). The misses generated by (c,d) are distribution misses.

We cannot provide a method for computing exactly distribution misses because we do not know how to easily find an optimal data-element distribution into blocks (we would have to try all address permutations). However, we can easily compute a lower bound for the miss ratio of a perfect cache with an optimal distribution, so we will use this to evaluate the quality of a distribution. Data elements cannot be smaller than a byte, so a perfect cache of size C with 1 byte per cache lines generates a minimal miss ratio $M1$ independent of data distribution.

Property 1 *A perfect cache of size C and line size L generates a miss ratio M greater or equal to $M1/L$, even with a perfect layout.*

If we approximate the miss ratio of a perfect cache with optimal data distribution to $M1/L$, we can split the capacity miss ratio according to our data distribution framework. The capacity miss ratio is the sum of the distribution miss ratio and the fundamental miss ratio:

Definition 9 *We call Fundamental miss ratio $M1/L$. It is not possible with a given cache size and line size L to generate fewer misses.*

Definition 10 *We assimilate $M - (M1/L)$ to the Distribution miss ratio.*

We can make this approximation for two reasons. First, $M1/L$ is a lower miss ratio bound, so it is not possible to get negative distribution miss rates⁵. Secondly, current caches are quite large, and with larger caches $M1/L$ will be closer to the real miss ratio of the perfect cache with optimal layout. An extreme case is when the whole working set fits in the cache: data distribution does not matter any more, since, no matter the reference pattern, the data will be in the cache.

⁵This happened to Sugumar, and led him to reconsider the definition of the 3C's miss rates using Belady's *Min* policy

The miss classification becomes:

$$total\ miss \left\{ \begin{array}{l} Cold \\ Capacity \\ Conflict \end{array} \right\} \left\{ \begin{array}{l} Fundamental \\ Distribution \\ Replacement \\ Mapping \end{array} \right.$$

Data-elements distribution into the blocks is an important factor, it changes the distribution miss ratio, impacting on all cache configurations, from direct-mapped to fully-associative caches. For direct-mapped and set-associative caches, block layout in memory also affects performance by changing mapping miss ratios. Since block layout impact is already covered by Sugumar's mapping-conflict miss definition, the miss classification needs no more refinement to evaluate data distribution optimizations. Data distribution and block mapping effects can combine to provide a great optimization potential. Only fundamental misses cannot be removed.

4 Irregularly referenced data-layout

We have seen that data layout affects the cache miss ratio, and we have defined a miss classification to explain the impact of data layout. In section 4.1 we show how to find for a direct-mapped cache all data distributions which generate only capacity misses, if any one exists. However, measuring miss ratios does not help to evaluate a priori the quality of a data layout, and necessitates a cache simulation to try all possibilities. Instead, we propose in section 4.2, metrics to evaluate the potential for locality of reference generated by program reference patterns. We use these metrics in the scalar data-layout heuristic shown in section 4.3 to generate the preliminary experimental data-distribution results.

4.1 Seeking a Belady conformant layout

A direct-mapped cache uses only the block address to select a cache line. If the data address is chosen without taking into account memory hierarchy organisation, the cache generates mapping misses. However, with data distribution it is possible to choose the addresses of data elements. We propose a heuristic which can find, if they exist, all data distributions which generate only capacity misses for a given direct-mapped cache.

Data distribution can be seen as an attempt to encode a block replacement policy in data element addresses. If there exists data distributions which generate only capacity misses for a given direct-mapped cache, one of them has the lowest capacity miss ratio. That data element distribution is optimal.

The problem of this approach is that there may not be any distribution generating only capacity misses for the direct-mapped cache. Data elements can only be assigned one address in memory. The optimal block replacement policy, on the other hand, may need to map a

```

NextDataReferenced() : returns the de referenced in the execution stream at a given time or  $\emptyset$  at the end.
MisCount() : returns the number of misses generated by a given layout  $\neq \emptyset$ , and  $\infty$  otherwise.
Block() : returns the block of Layout in which a given de is stored, or  $\emptyset$  otherwise.
Line() : returns the set of lines in which the given blocks map to.
Cache() : returns the subset of blocks of Layout held in the cache at a given time.
GetBelady() : returns the subset of blocks of Layout farthest referenced at a given time.
SizeUsed() : returns the memory space used up to store the de(s).
NewBlock() : allocates a new block in memory, holding a given de and mapping into a given cache line.
Layout : List of blocks used. A block at a given address will hold a given set of data elements de

BestLayout = PLACE(0,  $\emptyset$ )

PLACE(t, Layout) {
  de = NextDataReferenced(t, ProgramTrace)
  If de ==  $\emptyset$  // end of trace reached for a layout conforming to Min
  Else If Block(Layout, de)  $\in$  Cache(Layout, t) // de already allocated, cache Hit
    Layout = PLACE(t + 1, Layout)
  Else If Block(Layout, de)  $\neq \emptyset$  // de already allocated, cache Miss
    If Line(Block(Layout, de))  $\subset$  Line(GetBelady(Layout, t))
      Cache(Layout, t) = Cache(Layout, t)  $\cup$  Block(Layout, de) -
        {b | (b  $\in$  Cache(Layout, t))  $\wedge$  (Line(b) == Line(Block(Layout, de))}
      Layout = PLACE(t + 1, Layout)
    Else // Conflict miss: abort to try another layout
      Layout =  $\emptyset$ 
    End If
  Else // de not yet allocated: try all possible de layouts
    BetterLayout =  $\emptyset$ 
    For all B  $\in$  Cache(Layout, t) // Filling all blocks already in the cache, cache hit
      If SizeUsed(B) + SizeUsed(de)  $\leq$  BlockSize
        B = B  $\cup$  {de} // note: the offset address of de within B is unimportant
        TempLayout = PLACE(t + 1, Layout)
        if MisCount(TempLayout) < MisCount(BetterLayout) BetterLayout = TempLayout
        B = B - {de}
      End If
    Loop // Filling all possible blocks not held in the cache, cache miss
    For all B  $\in$  {b | (b  $\in$  Layout)  $\wedge$  (b  $\notin$  Cache(Layout, t))  $\wedge$  (Line(b)  $\subset$  Line(GetBelady(Layout, t))}
      If SizeUsed(B) + SizeUsed(de)  $\leq$  BlockSize
        B = B  $\cup$  {de}
        TempLayout = PLACE(t + 1, Layout)
        if MisCount(TempLayout) < MisCount(BetterLayout) BetterLayout = TempLayout
        B = B - {de}
      End If
    Loop
    For all L  $\in$  Line(GetBelady(Layout, t)) // Allocating a new empty block, cache miss
      B = NewBlock(de, L)
      TempLayout = PLACE(t + 1, Layout  $\cup$  {B})
      if MisCount(TempLayout) < MisCount(BetterLayout) BetterLayout = TempLayout
    Loop
    Layout = BetterLayout
  End If
return Layout } // returns a Min conformant layout or  $\emptyset$ 

```

Figure 2: Optimal data layout search algorithm

block in different cache lines at different times. If we use data addresses to implement the optimal replacement policy on a direct-mapped cache, this is not possible.

If we assume that there exist data distributions which generate only capacity misses for a direct-mapped cache, we can provide a heuristic which can find them by trying a subset of all possible data distributions.

Belady's optimal block replacement policy, *Min*, always replaces the block farthest referenced in time. If all blocks in the cache are live, there is only a single line chosen by *Min*. If there are dead blocks, any one of them is a potential target for replacement. A block distribution generates only capacity misses if blocks are always loaded in a line chosen by *Min*.

Min choices depend on block reference patterns. Data elements distribution into blocks impacts on these access patterns. Therefore to find the optimal layout, an approach is to try all the data element layout solutions which conform to *Min* while using a direct-mapped cache. Although data-element distribution into blocks can be chosen arbitrarily, block layout is dictated by *Min*, reducing the possibilities to analyse.

The algorithm of figure 2 is a recursive tree-like search to find an optimal layout. It fills up blocks by allocating data elements into them as they appear in the trace.

If a data element has not yet been referenced, it can be assigned to any block already in the cache which has enough room to accommodate it. We have a cache hit. It can also be assigned to any block not yet in the cache, but mapping to the cache line chosen by *Min*. We have a capacity miss. The data element cannot be assigned to blocks which are not in the cache and which do not map to the line selected by *Min*, because this generates a conflict miss.

If a data element has been referenced before, it is assigned to a block. If the block is in the cache it's a cache hit. If the block is not in the cache, but maps to the cache line selected by *Min*, then it's a capacity miss. Otherwise, it's a conflict miss, so the layout is inadequate. We have to modify a previous layout choice and assign the data-element to a different block.

If we reach the end of the trace, we have found a layout generating only capacity misses. It may not be the only one generating only capacity misses, so we must continue to try other layout choice in order to find all layouts generating capacity misses. The layout with the fewest capacity misses is the optimal layout.

Even after a full tree search, this algorithm finds the best layout available if this layout generates only capacity misses. This is too restrictive and too expensive. We wish on the contrary to have a heuristic which, even though it may not find the optimal layout, always provides an improved layout. To do so, we must find a layout by evaluating locality of references instead of measuring miss ratios. Therefore, we introduce new metrics to quantify locality of reference patterns over time.

4.2 Quantifying locality of references

To get good cache performance, it is important that memory reference patterns display very good locality properties. Data elements with interleaved references must be redistributed

into the same blocks to reduce distribution misses, and blocks with frequently interleaved references must not contend for the same cache lines, to remove mapping misses hazards. We introduce some notions needed to quantify formally these aspect of locality of references.

Caches use spatial and temporal locality properties to attempt to keep simultaneously variables that are “*used at the same time*” (interleaved, repeated references). To measure this effect, it is necessary to have a notion of time.

Min, measures time as the count of references separating 2 references to an element. However, Belady notes that replacing erroneously a block is less disastrous if the next reference to it is far in time than if it is very close [Bel66]: until its next reference, it won’t interfere with references to other block. Analyzing an unknown number of references into the future is not practical to define global metrics, and we are not constraining ourselves to finding the optimal layout solution, but only to finding a good layout at the lowest expense. Therefore, we introduce a time window to account for this approximation.

Definition 11 We define a **time reference window** W as a time interval of size $\|W\|$. Two references are said to be **close in time** if the time interval between them is smaller than $\|W\|$.

The choice of the window size is a sampling problem: if the window is too small, some interleaved reference patterns, considered too far, will not be seen. If the window is too big we will fall back to the complexity of *Min* measures.

Data elements that are unused while their block is in the cache waste cache space used to store them as well as memory bandwidth used to load them. Elements referenced farther apart than $\|W\|$ have a higher likelihood of necessitating a reloading of the block, so the block should not be considered as benefiting of spatial locality between them.

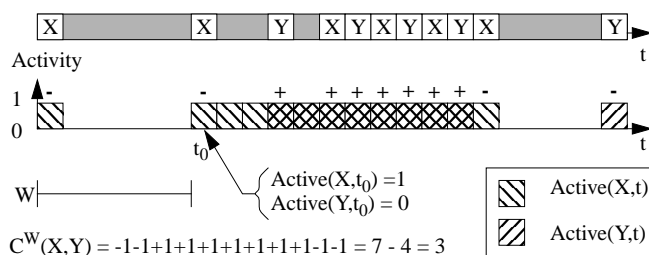


Figure 3: Computing *Correlation* between 2 data-elements

To depict the fact that an element X of memory is referenced or not within $\|W\|$, we define its activity (as shown in figure 3).

Definition 12 We note $Ref(X,t)$ the predicate that is 1 if X is referenced at time t and 0 otherwise.

Definition 13 An element X of memory is **active** at a time t , if t is between two references to X that are not separated by a time interval greater than $\|W\|$.

$A^W(X, t) = 1$ if $\exists(t', t'') \mid (t \in [t', t'']) \wedge (Ref(X, t') = 1, Ref(X, t'') = 1) \wedge (|t'' - t'| \leq ||W||)$, and 0 otherwise.

Now we can provide a correlation measure based on activity, to indicate if 2 elements have often closely interleaved references (see figure 3), and are rarely referenced independantly. Correlation is a global measure over the whole program (time independant) to evaluate the relevance of regrouping two scalars for spatial locality improvements.

Definition 14 *The correlation $C^W()$ of two data-elements is the count of all references done while both are active minus all the references to each while the other element is inactive.*

$$C^W(X, Y) = \sum_{t=0}^{t_{end}^d} ((A^W(X, t) \times Ref(Y, t)) + (Ref(X, t) \times A^W(Y, t)) - ((1 - A^W(X, t)) \times Ref(Y, t)) - (Ref(X, t) \times (1 - A^W(Y, t))))$$

An optimal layout is therefore a layout for which highly correlated data-elements are stored in the same block, and for which highly correlated blocks map to distinct cache lines.

4.3 A direct-mapped cache data-layout heuristic

With the *Activity* and *Correlation* metrics, we developped a profile-based scalar data-elements layout optimization heuristic for direct-mapped caches. We used this heuristic to provide preliminary data-layout results to show the potential gains of data layout.

The algorithm, described in figure 4, is decomposed into four steps:

- *Computing activity* of each data-element at each time t it is referenced during profiling.
- *Redistribution of correlated data-elements into the same blocks.* We repeatedly insert in a block the data element correlating most with it, and recompute block's activity, until it is full. Data elements are not allowed to overlap 2 blocks.
- *Distribution of correlated blocks into distinct cache lines.* We seek the block correlating the most with the blocks already laid-out. We then assign it to the cache line that correlates the least with it. We update cache and line activities and lay another block.
- *Memory layout* is done by assigning blocks any free memory address that maps into the selected cache line (we assume a direct-mapping). We also pick the block offset of data-elements arbitrarily – it does not affect miss ratios.

We obtain a layout that uses a minimum number of memory blocks to satisfy a given number of memory references. The layout also reduces the risks of having blocks used simultaneously in the cache conflicting for the same cache line.

This heuristic always finds an efficient layout solution. However it is still quite costly and cannot be used as-is in an optimizing compiler:

- Profiling needs to store activity information for each data-element over the whole execution trace, this can be too costly for huge programs. Data-elements placement needs $O(n^2)$ correlation computations and blocks layout needs $O(b^2)$, n the number


```

NewBlock(), NewLine(), NewLineList(): Creates a new block, line or list of lines, and resets its Activity.
SizeUsed(): returns the memory space used up by a given element to store the de(s).
Correlate(): Computes the Correlation of 2 elements or  $\int_t A^W dt$  of an element if the other is  $\emptyset$ .

BlockList DATAELEMENTLAYOUTINBLOCKS(deList)
  BlockList =  $\emptyset$  // List of blocks with data-elements layed out
  While deList  $\neq \emptyset$ 
    B = NewBlock()
    BlockList = BlockList  $\cup$  {B}
    While SizeUsed(B) < BlockSize
      Bestde =  $\emptyset$ 
      For all de  $\in$  deList
        If Correlate(de, B)  $\geq$  Correlate(Bestde, B)
          Bestde = de
        End If
      Loop
      B = B  $\cup$  {Bestde} // This modifies the Activity of B
      deList = deList - {Bestde}
    Loop
  Loop
  return BlockList

LineList BLOCKLAYOUTINLINES(BlockList)
  LineList = NewLineList() // List of lists of blocks mapping into the same cache lines
  For i = 1 to CacheLineCount
    L = NewLine() // Creates a list of blocks mapping to the same cache line
    LineList = LineList  $\cup$  {L}
  Loop
  While BlockList  $\neq \emptyset$ 
    BestBlock =  $\emptyset$ 
    For all B  $\in$  BlockList
      If Correlate(B, LineList)  $\geq$  Correlate(BestBlock, LineList)
        BestBlock = B
      End If
    Loop
    BestLine =  $\emptyset$ 
    For all L  $\in$  LineList
      If Correlate(L, BestBlock)  $\leq$  Correlate(BestLine, BestBlock)
        BestLine = L
      End If
    Loop
    BestLine = BestLine  $\cup$  {BestBlock} // This modifies the Activity of BestLine and LineList
    BlockList = BlockList - {BestBlock}
  Loop
  return LineList

Compute Addresses using LineList de and block layout constraints.

```

Figure 4: Data-layout Heuristic

of data-elements, and b the number of blocks. This is too expensive if n and b are big (i.e. a whole program's data).

- This heuristic applies only for scalars, and must be adapted to handle complex data structures.

Taking into account data structuring and compiler data layout constraints necessitates a more complex analysis. However, it can also help the scalar data layout: constraints reduce the number of elements that can be redistributed together, reducing n greatly (for example, redistribution of fields within a C structure). These layout constraints reduce the layout possibilities, and could reduce data layout gains. However, current caches are large, providing a high tolerance on layout. Furthermore, program and compiler layout constraints often reflect locality in data usage. It is locality generated by these data organisation constraints that caches capture. They are however too weak to provide good locality as application data sets grow larger. Data distribution can therefore provide the extra layout constraints needed based on locality properties. More work will be necessary to verify this, however.

To show that distribution heuristics can be developed, we wrote one which lays scalars using profiling data. In the next section we use it to show that indeed, miss rates are improved.

5 Experimentation

We provide experimental results to show the potential of data distribution. Fundamental, distribution, capacity and conflict miss ratios are given for direct-mapped, 2-way and 4-way set associative caches. The capacity miss ratio is the miss ratio of a perfect cache. The fundamental miss ratio is computed using the miss ratio of a perfect cache with 1 byte per line. Cold misses are counted as capacity misses. Results are given for the original and improved layouts.

5.1 Experimental set-up

Experimental results were obtained using traces generated by *Spy*, a tracing tool working under SunOS-4 [Ir192]. The cache simulated is a 16 KB cache with 32 bytes per lines, and true LRU is used for associative caches. We developed a range of tools to gather the miss ratios, figure 5.

- *Belady* provides the miss ratio of a perfect cache (a fully-associative cache using Belady's *Min* replacement policy). It is used to get the capacity miss rates, as well as the fundamental miss rate (we simulate 1 byte lines and divide the miss rate obtained by the line size, 32). *Belady* does not simulate set-associative caches, so we do not present mapping and replacement miss ratios. We do not present cold miss ratio either, since it is irrelevant for our study (divide the number of bytes used by the line size).
- The *CacheSkew* and *NewSim* cache simulators provide direct and set-associative cache miss rates.

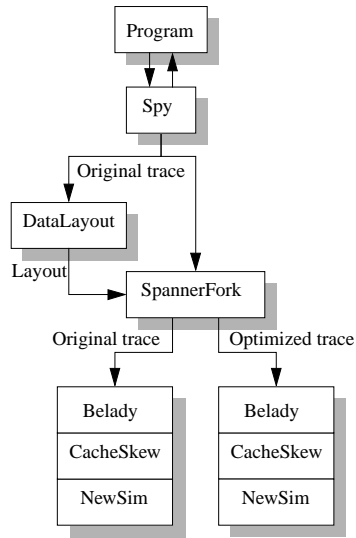


Figure 5: Simulation tools used

- *DataLayout* is the program which optimizes the data layout. it generates a layout information file decoded by *SpannerFork*. *SpannerFork* outputs an optimized layout trace by replacing data addresses with optimized addresses.

A trace simulation approach was chosen, because our heuristic can only layout scalars, and in a trace, all data references are seen as scalar references. To handle data structuring constraints, more work will be necessary. Our goal is not to provide a finished product, but to show the potential behind data distribution. The downside is a slow simulation process since we have to compute correlation between all the scalars (with compiler layout constraints we would have smaller sets of scalars to correlate), so traces are only a few million references wide.

The *DataLayout* program implements a simplified version⁶ of the heuristic presented in section 4.3:

- We used a fixed instead of a sliding time window. This cuts time space into $\|W\|$ sized chunks. References in distinct windows are considered far appart, although they might be close if we used a sliding window (for example $t_i \in W_1$ and $t_j \in W_2$ and $t_i = t_j + 1$).
- One bit is used to mark up if a data-element is referenced in a window. This reduces storage needs of the reference/activity array to 1 bit per *de* per W_i . Of 2 elements

⁶Our goal is not to provide the most efficient layout, but only to show that better layouts can provide low miss rates

active in a window, we cannot tell which one is most referenced. if W is too large, our program will find that all elements correlate. $A^W(de, t) = A^W(de, W_i) = 1$ if $(t \in W_i) \wedge (\exists t' \in W_i | Ref(de, t') = 1)$

- Correlation is simplified to a sum of AND-ed bits.

$$C^W(X, Y) = \sum_{W_i=0}^{W_{End}} (A^W(X, W_i) \wedge A^W(Y, W_i))$$

- to eliminate block activity recomputations during data-elements distribution, the block's activity is approximated to the activity of the first element stored in the block – it is the most referenced element available when the block is created.

5.2 Tracing “ls”

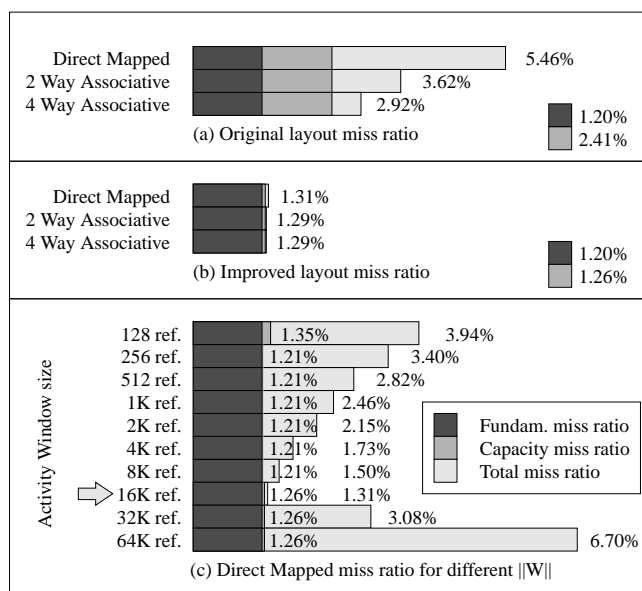


Figure 6: Miss ratios for “ls”

We first attempted our optimization on a very small trace of 102,614 memory references generated by the Unix *ls* command. It references 16,149 distinct data-elements for a total of 39,471 bytes. These preliminary results (figure 6) already give some insights:

- The original capacity miss rate of 2.41%, figure 6(a), is reduced to 1.26% with improved data-element distribution, figure 6(b), a 48% gain. This is due to a 95% reduction of the distribution miss rate (from 1.21% to 0.06%). Note that fewer capacity misses

means fewer block replacements, therefore fewer chances of making a wrong block replacement choice (i.e. conflict miss).

- The direct-mapped cache miss ratio drops from 5.46% (a) to 1.31% (b), a 76% gain. The effects of our data-distribution and block layout optimizations combine to provide a total miss rate which is actually lower than the capacity miss rate of the original layout alone. This shows the dramatic impact of data distribution on cache performance: we even bested the capacity miss ratio.
Contrarily to hardware cache optimizations, with data-layout optimizations, the miss ratio obtained with a perfect cache is not sufficient to evaluate the lowest miss ratio attainable with a real cache.
- Associative caches also benefits from the improved layout. The total miss ratio of the 4 way set-associatif cache drops from 2.92% to 1.29%, a 56% gain. For this example, data distribution is so performant that using an associative cache yields little gain over a direct-mapped cache. In fact, the direct-mapped cache miss ratio (1.31%) is already close to the fundamental miss ratio (1.20%), so even a fully-associative cache could not provide any gain (in fact, it could even degrade performance by generating replacement misses, since LRU is not optimal).
- The heuristic we used to improve data-layout uses a time window to measure data correlation. The choice of the window impacts on the layouts obtained. We tried a range of time windows to see the impact of this choice on miss rates obtained. For the capacity miss, the variation induced is small (c). The window must be large enough (over 256 references) to find data-elements to regroup into blocks, and small enough (less than 8K references) to differentiate between elements that should not be in the cache simultaneously. Since the cache is quite large there is a high tolerance on the distribution of the data elements into blocks. The window impacts much more on block layout. If the window is too small, mapping conflicts arise because blocks are not correlated sufficiently well (repetitive, but far appart reference patterns are not seen). If the window is too large, data-elements distribution becomes too laxist to allow efficient block mapping, and the heuristic also becomes incapable of discriminating between weak and strong correlations (because our implementation uses only a single bit to measure activity). Since our heuristic uses the same window size $||W||$ to distribute data-elements and perform block layout, a window size of 16K references seems the best compromize. We will use this value for all the other experiments.

An optimal data layout can reduce the direct mapped cache miss ratio by up to 78% (by leaving only 1.20% fundamental misses). Our heuristic achieves a 76% reduction, far better than what set-associativity alone can achieve with the original layout.

5.2.1 Tracing “f77” compiler

We traced 2 million memory references of the f77 compiler under Sun-OS 4 while it was compiling 078.swm256. We first discarded 2 million references to remove program startup effects. The fundamental miss rate is 1.35% ($\frac{907,507}{32} = 28,360$ misses).

Table 1: f77 compiler miss ratios for direct mapped (DM), 2-way and 4-way set-associative (SA) caches. Capacity misses are the sum of distribution and fundamental misses. Half the misses are capacity misses, and half the capacity misses are distribution misses. Data distribution reduces distribution and conflict misses by 53% and 38% respectively on a direct-mapped cache, for a total miss reduction of 33%.

| Cache | DM | 2-SA | 4-SA |
|-----------------------------------|--------------|--------------|--------------|
| Original data distribution | | | |
| Fundamental | 1.35% | 1.35% | 1.35% |
| Distribution | 1.38% | 1.38% | 1.38% |
| Capacity | 2.73% | 2.73% | 2.73% |
| Conflict | 2.95% | 1.78% | 1.42% |
| Total | 5.68% | 4.51% | 4.15% |
| Improved data distribution | | | |
| Fundamental | 1.35% | 1.35% | 1.35% |
| Distribution | 0.65% | 0.65% | 0.65% |
| Capacity | 2.00% | 2.00% | 2.00% |
| Conflict | 1.83% | 1.26% | 1.01% |
| Total | 3.83% | 3.26% | 3.01% |

In the original layout, 51% of the capacity misses are distribution misses which could be removed. Our layout heuristic removes 53% of the distribution misses (from 1.38% to 0.65%), reducing the capacity miss rate by 27%, table 1. Furthermore we remove 38% of the conflict misses for the direct mapped cache (2.95% to 1.83%). After layout optimization, there are still conflict misses. An associative cache can remove some more conflicts (1.01% left for a 4-way associative cache).

A perfect data layout can reduce the original direct-mapped cache miss ratio by up to 76% (from 5.68% down to 1.35%). Our layout heuristic manages a 33% reduction (down to 3.83%). By cumulating set-associativity and data distribution, gains reach 47%.

5.2.2 Tracing “078.swm256” Specfp-92 benchmark

After discarding the first million references, we traced 5 million references of 078.swm256 (compiled under Sun-OS 4 using -O4 flag). This program accesses arrays sequentially. The intrinsic miss rate is 0.242% ($404,424/32 = 12,639$ misses).

Table 2: SpecFp 92 078.swm256. When arrays do not interfere, the total miss ratio is close to the fundamental miss ratio.

| Cache | DM | 2-SA | 4-SA |
|-----------------------------------|---------------|---------------|---------------|
| Original data distribution | | | |
| Fundamental | 0.242% | 0.242% | 0.242% |
| Distribution | 0.001% | 0.001% | 0.001% |
| Capacity | 0.243% | 0.243% | 0.243% |
| Conflict | 0.095% | 0.001% | 0.002% |
| Total | 0.338% | 0.244% | 0.245% |
| Improved data distribution | | | |
| Fundamental | 0.242% | 0.242% | 0.242% |
| Distribution | 0.000% | 0.000% | 0.000% |
| Capacity | 0.242% | 0.242% | 0.242% |
| Conflict | 0.010% | 0.010% | 0.010% |
| Total | 0.252% | 0.252% | 0.252% |

We show with this example, that a well behaved scientific application makes almost perfect usage of the cache line space, there are almost no distribution misses, table 2.

Some conflicts appear in the direct mapped cache, easily removed by set-associative caches, or data-distribution. The only data-distribution alternative for scientific codes, padding, may not suffice. Contrarily to non-scientific applications, loop-transformations can be used efficiently to generate tiled or stride-one accesses.

5.2.3 Tracing “gzip”

We traced 2 million references of gzip while it was compressing a trace. We first discarded 2 million references to remove program startup effects. The fundamental miss rate is 0.71% ($476,185/32 = 14,881$ misses). The 141,379 data-elements referenced represent 213KB of data.

Although the 0.71% fundamental miss rate is quite low, GZip 9.57% original capacity miss ratio is extremely high, due to distribution misses, table 3. Our heuristic removes 78% of the distribution misses (from 8.96% to 2%), a 72% reduction of capacity misses.

Improved block mapping reduces conflicts by 43% for the direct-mapped cache (from 11.54% to 6.63%).

The total miss rate is reduced by 56% (from 21.21% to 9.34%). It again beats the original capacity miss rate. Associativity also combines well with data distribution to remove conflicts. The 4 way set-associative cache miss ratio drops to 6.49%.

Table 3: Gzip. The distribution miss ratio is high for the original layout and drops to 2% with a better data distribution.

| Cache | DM | 2-SA | 4-SA |
|-----------------------------------|---------------|---------------|---------------|
| Original data distribution | | | |
| Fundamental | 0.71% | 0.71% | 0.71% |
| Distribution | 8.96% | 8.96% | 8.96% |
| Capacity | 9.67% | 9.67% | 9.67% |
| Conflict | 11.54% | 9.27% | 8.27% |
| Total | 21.21% | 18.94% | 17.94% |
| Improved data distribution | | | |
| Fundamental | 0.71% | 0.71% | 0.71% |
| Distribution | 2.00% | 2.00% | 2.00% |
| Capacity | 2.71% | 2.71% | 2.71% |
| Conflict | 6.63% | 4.54% | 3.78% |
| Total | 9.34% | 7.25% | 6.49% |

5.2.4 Experimental summary

We defined distribution misses and divided capacity misses into distribution and fundamental misses. The fundamental miss ratio is usually quite low (less than 1.35% in all our experiments), leaving much room for cache miss rates improvement. Capacity miss rates are usually quite high due to distribution misses (up to 9.67% in gzip). Data-distribution heuristics can be used effectively to reduce distribution misses.

Layout also impacts on conflict misses. Better data-distribution reduce conflict miss hazards, while improved block layout into memory reduce mapping miss rates. By choosing block addresses in memory, we actually provide compile time hints on the future to the hardware block replacement policy.

We used a simple layout optimization heuristic to show that it is possible to design a tool capable of finding improved layouts. Layout choices made by our heuristic usually provided significant capacity miss reductions of up to 71%, and direct-maped cache miss reductions of up to 76%.

6 Conclusion

In many applications, cache space is under-utilized, yielding excessive miss ratios. Loop transformations can be used to modify array reference patterns, but non-scientific applications tend to use other data types. We show that data-distribution is an interesting approach to provide software locality optimization techniques for the other data types.

Data distribution consists in choosing the relative address locations of data elements in memory. We provide a formalism for data distribution. Since data-distribution optimi-

zations impact on capacity miss ratios, we define distribution misses to account for poor data-element distribution into blocks and compute a lower miss ratio boundary, the fundamental miss ratio. We also define correlation between data elements to evaluate spatial locality. It facilitates the development of data distribution heuristics.

Data distribution for caches is a two faceted problem. Memory blocks layout in memory must be selected according to the cache mapping policy and block reference patterns to reduce mapping-conflict misses. More importantly for non-scientific applications, we have shown that proper data-elements redistribution into blocks impacts significantly on the capacity miss rate (we removed up to 71% of the capacity misses). Data distribution can count on a great potential for miss reduction, since limit is only the fundamental miss ratio. The data distribution heuristic we developed was not able to reach this limit, but already managed significant gains for direct mapped caches, of up to 76%, sometimes even beating the capacity miss ratio of the original layout.

This study is a first step to evaluate the potential of data-layout for cache optimization of non-scientific programs. These first results hint that it is important and possible for compilers to pay attention to how data is placed in memory, and that doing so can remove a significant amount of capacity and conflict misses. Currently compiler data layout constraints are weak, and are never based on data locality considerations. We feel that one of the most promising aspects of data layout optimization for non-scientific codes, is the optimisation of the layout of complex data structures allocated dynamically (C struct and C++ classes). We are currently developing a new allocation library which will handle data-layout optimizations. Other interesting aspects of data-distribution are fast scalar layout, array padding for caches and instruction basic block layout and prefetching.

Table 4: Total miss counts used to compute cache miss ratios

| Cache | Perfect | DM | 2-SA | 4-SA |
|-------------------|---------|--------|--------|--------|
| f77 | | | | |
| Fundamental | 28360 | | | |
| Original | 57211 | 119177 | 94678 | 87213 |
| Improved | 41950 | 80217 | 68262 | 63106 |
| 078.swm256 | | | | |
| Fundamental | 12639 | | | |
| Original | 12655 | 17646 | 12721 | 12787 |
| Improved | 12640 | 13149 | 13147 | 13147 |
| gzip | | | | |
| Fundamental | 14881 | | | |
| Original | 202886 | 444839 | 397174 | 374833 |
| Improved | 56867 | 195872 | 151970 | 136205 |

Acknowledgements

We would like to thank Richard Uhlig and Olivier Michel for proofreading this article and G. Irlam for developing the spy package.

References

- [AP93] A. Agarwal and S. D. Pudar. Column associative caches: A technique for reducing the miss rate of dm caches. *Proceedings of the 20th International Symposium on Computer Architecture*, 1993.
- [ASU86] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers principles, Techniques and Tools*. Addison Wesley, 86.
- [BC91] J. L. Baer and T-F. Chen. An effective on chip preloading scheme to reduce data access penalty. *Proceedings of Supercomputing 91*, 1991.
- [BCJ94] D. F. Bacon, J-H. Chow, and D. R. Ju. A compiler framework for restructuring data declarations to enhance cache and tlb effectiveness. *Proceedings of CASCON'94*, November 1994.
- [Bel66] L. A. Belady. A study of replacement algorithms for a virtual-storage computer. *IBM System Journal*, page 24, 1966.
- [BGK95] D. C. Burger, J. R. Goodman, and A. Kagi. The declining effectiveness of dynamic caching for general purpose multiprocessor. Technical report, University of Winsconsin, 1261, 1995.
- [BGK96] Doug Burger, James R. Goodman, and Alain Kagi. Memory bandwidth limitations of future processors. In *23rd Annual International Symposium on Computer Architecture*, pages 78–89, May 1996.
- [BGS94] D. F. Bacon, S. L. Graham, and O. J. Sharp. Compiler transformations for high-performance computing. *ACM Computing Surveys*, 26(4):345–420, December 1994.
- [BS95] F. Bodin and A. Seznec. Skewed associativity enhances performance predictability. *Proceedings of the 22nd International Symposium on Computer Architecture*, 23(2):265, May 1995.
- [BZ93] D. A. Barrett and B. G. Zorn. Using lifetime predictors to improve memory allocation performance. *Proceedings of the ACM SIGPLAN'93 Conference on Programming Language Design and Implementation*, 28(6):187–196, July 1993.
- [CB92] T. F. Chen and J. L. Baer. Reducing memory latency via non blocking and prefetching caches. *5th symposium on Architectural Support for Programming Languages and Operating Systems*, 1992.

- [CCK90] D. Callahan, S. Carr, and K. Kennedy. Improving register allocation for subscribed variables. *Proceedings of the ACM SIGPLAN'90 Conference on Programming Language Design and Implementation*, 1990.
- [CG94] B. Calder and D. Grunwald. Reducing branch costs via branch alignment. In *Sixth International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 242–251, November 1994.
- [CKP91] D. Callahan, K. Kennedy, and A. Porterfield. Software prefetching. *4th symposium on Architectural Support for Programming Languages and Operating Systems*, 1991.
- [CKT94] S. Carr, K. S. Mc Kinley, and C-W. Tseng. Compiler optimisation for improving data locality. *proceedings of the 4th symposium on Architectural Support for Programming Languages and Operating Systems*, page 11, 94.
- [DSW95] N. Drach, A. Sez nec, and D. Windheiser. Direct-mapped versus set associative pipelined caches. *Proceedings of PACT'95*, 1995.
- [FJ94] K. I. Farkas and N. P. Jouppi. Complexity/performance tradeoffs with non-blocking loads. *Proceedings of the 21st International Symposium on Computer Architecture*, page 12, 1994.
- [FPJ92] J. W. C. Fu, J. H. Patel, and B. L. Janssens. Stride directed prefetching in scalar processors. *Proceedings of Micro-25*, 1992.
- [GAFN94] K. Grimsrud, J. Archibald, R. Frost, and B. Nelson. On the accuracy of memory reference models. In Springer-Verlag, editor, *7th International Conference on Computer Performance Evaluation*, volume 794 of *Lecture notes in computer science*, pages 369–388, May 1994.
- [GMB95] Elana D. Granston, Thierry Montaut, and F. Bodin. Loop transformations to prevent false sharing. *International Journal of Parallel Programming*, 1995.
- [GZH93] D. Grunwald, B. Zorn, and R. Henderson. Improving the cache locality of memory allocation. *Proceedings of the ACM SIGPLAN'93 Conference on Programming Language Design and Implementation*, 28(6):177, July 1993.
- [HC89] W-M Hwu and P. P. Chang. Achieving high instruction cache performance with an optimizing compiler. In *16th Annual International Symposium on Computer Architecture*, pages 242–251, May 1989.
- [Hei94] R. R. Heisch. Trace-directed program restructuring for aix executables. *IBM Journal of Research and Development*, pages 595–603, September 1994.
- [Hil87] M. D. Hill. Aspects of cache memory and instruction buffer performance. Technical Report UCB//CSD-87-381, University of California, Berkeley, 1987.

-
- [Irl92] G. Irlam. “spa” personal communication. the Spa package is available from gordon@cs.adelaide.edu.au, 1992.
- [Jou90] N. P. Jouppi. Improving direct mapped cache performance by the addition of a small fa cache and prefetch buffers. *Proceedings of the 17th International Symposium on Computer Architecture*, 1990.
- [JW94] N. P. Jouppi and S. E. Wilton. Tradeoffs in two level on-chip caching. *Proceedings of the 21st International Symposium on Computer Architecture*, 1994.
- [PH90] Karl Pettis and Robert C. Hansen. Profile guided code positioning. In *ACM SIGPLAN’90 Conference on Programming Languages and Design Implementation*, pages 16–27, June 1990.
- [PK94] S. Palacharla and R. E. Kessler. Evaluating stream buffers as a secondary cache replacement. *Proceedings of the 21st International Symposium on Computer Architecture*, page 10, 1994.
- [Por89] A. K. Porterfield. *Software Methods for Improvement of Cache Performance on Supercomputer Applications*. PhD thesis, CRPC, 1989.
- [Prz90] S. A. Przybylski. *Cache and Memory Hierarchy Design*. Morgan Kaufmann, 1990.
- [Smi82] A. J. Smith. Cache memories. *ACM Computing Surveys*, 14(3):473–530, September 1982.
- [Smi86] A. J. Smith. Bibliography and readings on cpu cache memories and related topics. *Computer Architecture News*, pages 22–42, January 1986.
- [Sug93] R. A. Sugumar. *Multi Configuration Simulation Algorithms for the Evaluation of Computer Architecture Designs*. PhD thesis, University of Michigan, 1993.
- [TGJ93] O. Temam, E. D. Grandston, and W. Jalby. To copy or not to copy: A compile time technique for assessing when data copying should be used to eliminate cache conflicts. *Proceedings of Supercomputing’93*, 1993.
- [TLH90] J. Torrellas, M. S. Lam, and J. L. Hennessy. Shared data placement optimisation to reduce multiprocessor cache miss rates. *Proceedings of the International Conference on Parallel Processing*, 1990.
- [TXD95] J. Torrellas, Chun Xia, and Russel Daigle. Optimizing instruction cache performance for operating system intensive workloads. In *First Symposium on High-Performance Computer Architecture*, page 360369, January 1995.



Unit e de recherche INRIA Lorraine, Technop le de Nancy-Brabois, Campus scientifique,
615 rue du Jardin Botanique, BP 101, 54600 VILLERS L ES NANCY
Unit e de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex
Unit e de recherche INRIA Rh ne-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN
Unit e de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex
Unit e de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

 diteur
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399