

A Fully Empirical Autotuned Dense QR Factorization for Multicore Architectures

Emmanuel Agullo¹, Jack Dongarra², Rajib Nath², and Stanimire Tomov²

¹ LaBRI and INRIA Bordeaux Sud Ouest

² University of Tennessee

Abstract. Tuning numerical libraries has become more difficult over time, as systems get more sophisticated. In particular, modern multicore machines make the behaviour of algorithms hard to forecast and model. In this paper, we tackle the issue of tuning a dense QR factorization on multicore architectures using a fully empirical approach. We exhibit a few strong empirical properties that enable us to efficiently prune the search space. Our method is automatic, fast and reliable. The tuning process is indeed fully performed at install time in less than one hour and ten minutes on five out of seven platforms. We achieve an average performance varying from 97% to 100% of the optimum performance depending on the platform. This work is a basis for autotuning the PLASMA library and enabling easy performance portability across hardware systems.

1 Introduction

The hardware trends have dramatically changed in the last few years. The frequency of the processors has been stabilized or even sometimes slightly decreased whereas the degree of parallelism has increased at an exponential scale. This new hardware paradigm implies that applications must be able to exploit parallelism at that same exponential pace. Applications must also be able to exploit a reduced bandwidth (per core) and a smaller amount of memory (available per core). Numerical libraries, which are a critical component in the stack of high-performance applications, must in particular take advantage of the potential of these new architectures. So long as library developers could depend on ever increasing clock speeds and instruction level parallelism, they could also settle for incremental improvements in the scalability of their algorithms. But to deliver on the promise of tomorrow's petascale systems, library designers must find methods and algorithms that can effectively exploit levels of parallelism that are orders of magnitude greater than most of today's systems offer. Autotuning is therefore a major concern for the whole HPC community and there exist many successful or on-going efforts. The FFTW library [1] uses autotuning techniques to generate optimized libraries for FFT, one of the most important techniques for digital signal processing. Another successful example is the OSKI library [2] for sparse matrix vector products. The PetaBricks [3] library is a general purpose tuning method providing a language to describe the problem to tune. It has several applications ranging from efficient sorting to multigrid optimization. In

the dense linear algebra community, several projects have tackled this challenge on different hardware architectures. The Automatically Tuned Linear Algebra Software (ATLAS) library [4] aims at achieving high performance on a large range of CPU platforms thanks to empirical tuning techniques performed at install time. On graphic processing units (GPUs), among others, [5] and [6] have proposed efficient approaches. FLAME [7] and PLASMA [8] have been designed to achieve high performance on multicore architectures thanks to tile algorithms (see Section 2.1). The common characteristics of all these approaches are that they need intensive tuning to fully benefit from the potential of the hardware.

Tuning a library consists of finding the parameters that maximize a certain metric (most of the time the performance) on a given environment. In general, the term *parameter* has to be considered in its broad meaning, possibly including a variant of an algorithm. The *search space*, corresponding to the possible set of values of the *tunable parameters* can be very large in practice. Depending on the context, on the purpose and on the complexity of the search space, different approaches may be employed. Vendors can afford dedicated machines for delivering highly tuned libraries and have thus limited constraints in terms of time spent in exploring the search space. On the other side of the spectrum, some libraries such as ATLAS aim at being portable and efficient on a wider range of architectures and cannot afford a virtually unlimited time for tuning. Indeed, empirical tuning is performed at install time and there is thus a trade-off between the time the user accepts to afford to install the library and the quality of the tuning. In that case, the main difficulty consists of efficiently pruning the search space. Of course, once a platform has been tuned, the information can be shared with the community so that it is not necessary to tune again the library, but this is an orthogonal problem which we do not address here. Model-driven tuning may allow one to efficiently prune the search space. Such approaches have been successfully designed on GPU architectures, in the case of matrix vector products [2] or dense linear algebra kernels [5,6]. However, in practice, the robustness of the assumptions on the model strongly depends both on the algorithm to be tuned and on the target architecture. There is no clearly identified trend yet but model-driven approaches seem to be less robust on CPU architectures. For instance, even in the single-core CPU case, basic linear algebra algorithms tend to need more empirical search [4]. Indeed, on CPU-based architectures, there are many parameters that are not under user control and difficult to model (different levels of cache, different cache policies at each level, possible memory contention, impact of translation lookaside buffers (TLB) misses, ...) whereas the current generations of GPU provide more control to the user.

In a previous work, we had tackled the issue of maximizing PLASMA performance in order to compare it against other libraries [9]. We first manually pre-selected a combination of parameters based on the performance of the most compute-intensive kernel. We then tried all these combinations for each considered size of matrix to be factorized. This basic tuning approach achieved high performance but required human intervention to pre-select the parameters and days of run to find optimum performance. In the present paper, not only

we now tackle the issue of automatically performing the tuning process but we also present new heuristics that efficiently prune the search space so that the whole tuning process is reduced to one hour or so. We illustrate our discussion with the QR factorization implemented in the PLASMA library, which is representative [9] of all three one-sided factorizations (QR, LU, Cholesky) currently available in PLASMA. Because of the trends expose above, we do *not* rely on a model to tune our library (a detailed motivation based on a cased study can be found in Section 2.3 of our corresponding technical report [10]). Instead, we employ a fully empirical approach and we exhibit few empirical properties that enable us to efficiently prune the search space.

The rest of the paper is organized as follows. Section 2 presents the problem and motivates the outline of our two-step empirical approach (Section 3). Section 4 presents the wide range of hardware platforms used in the experiments to validate our approach. Section 5 describes the first empirical step, consisting of benchmarking the most compute-intensive serial kernels. We propose three new heuristics that automatically pre-select (PS) candidate values for the tunable parameters. Section 6 presents the second empirical step, consisting of benchmarking effective multicore QR factorizations. We propose a new pruning approach, which we call “prune as you go” (PAYG), that enables to further prune the search space and to drastically reduce the whole tuning process. We conclude and present future work directions in Section 7.

2 Problem Description

2.1 Tile QR Factorization

The development of programming models that enforce asynchronous, out of order scheduling of operations is the concept used as the basis for the definition of a scalable yet highly efficient software framework for computational linear algebra applications. In PLASMA, parallelism is no longer hidden inside Basic Linear Algebra Subprograms (BLAS) but is brought to the fore to yield much better performance. We do not present tile algorithms in details (more details can be found [8]) but their principles. The basic idea is to split the initial matrix of order N into $NT \times NT$ smaller square pieces of order NB , called *tiles*. Assuming that NB divides N , the equality $N = NT \times NB$ stands. The algorithms are then represented as a Directed Acyclic Graph (DAG) where nodes represent tasks performed on tiles, either panel factorization or update of a block-column, and edges represent data dependencies among them. More details on tile algorithms can be found [8]. PLASMA currently implements three one-sided (QR, LU, Cholesky) tile factorizations. The DAG of the Cholesky factorization is the least difficult to schedule since there is relatively little work required on the critical path. LU and QR factorizations have exactly the same dependency pattern between the nodes of the DAG, exhibiting much more severe scheduling and numerical (only for LU) constraints than the Cholesky factorization. Therefore, tuning the QR factorization is somehow representative of the work to be done

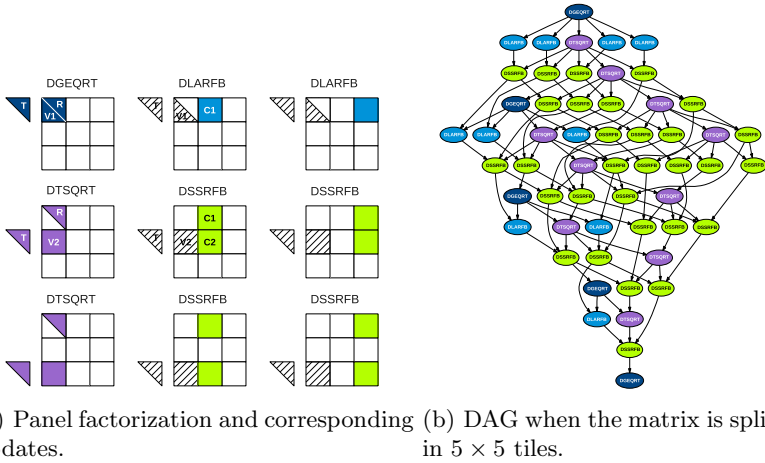


Fig. 1. Tile QR Factorization

for tuning the whole library. In the following, we focus on the QR factorization of square matrices in double precision statically scheduled in PLASMA.

Similarly to LAPACK which was built using a set of basic subroutines (BLAS), PLASMA QR factorization is built on top of four serial kernels. Each kernel indeed aims at being executed sequentially (by a single core) and corresponds to an operation performed on one or a few tiles. For instance, assuming a 3×3 tile matrix, Figure 1(a) represents the first panel factorization (DGEQRT and DTSQRT serial kernels [8]) and its corresponding updates (DLARFB and DSSRFB serial kernels [8]). The corresponding DAG (assuming this time that the matrix is split in 5×5 tiles) is presented in Figure 1(b).

2.2 Tunable Parameters and Objective

The shape of the DAG depends on the number of tiles ($NT \times NT$). For a given matrix of order N , choosing the tile size NB is equivalent to choosing the number of tiles (since $N = NB \times NT$). Therefore, NB is a first tunable parameter. A small value of NB induces a large number of tasks in the DAG and subsequently enables the parallel processing of many tasks. On the other hand, the serial kernel applied to the tiles needs a large enough granularity in order to achieve a decent performance. The choice of NB thus trades off the degree of parallelism with the efficiency of the serial kernels applied to the tiles. There is a second tunable parameter, called inner block size (IB). It trades off memory load with extra-flops due to redundant calculations. With a value $IB = 1$, there are $\frac{4}{3}N^3$ operations as in standard LAPACK algorithm. On the other hand, if no inner blocking occurs ($IB = NB$), the resulting extra-flops overhead may represent 25% of the whole QR factorization (see [8] for more details). The general objective of the paper is to address the following problem.

Problem 1. Given a matrix size N and a number of cores n_{cores} , which tile size and internal blocking size (NB - IB combination) do maximize the performance of the tile QR factorization?

Of course, the performance P we aim at maximizing shall not depend on extra-flops. Therefore, independently of the value of IB , we define $P = \frac{4}{3} \times N^3/t$, where t is the elapsed time of the QR factorization. Note also that we want the decision to be instantaneous when the user requests to factorize a matrix so that the tuning process is to be performed at install time.

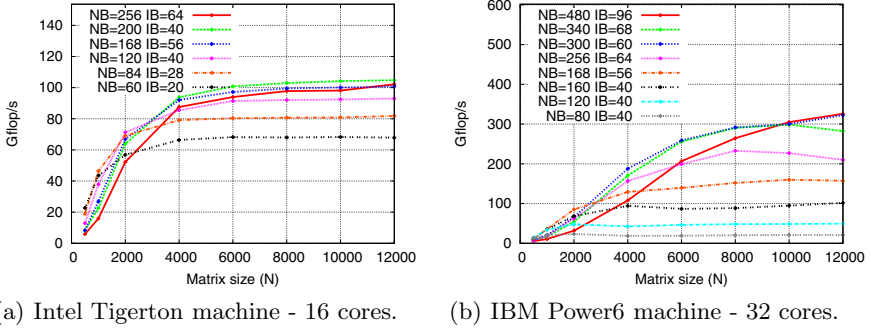


Fig. 2. Performance of the PLASMA QR factorization

In a parallel execution of PLASMA, the optimum tile size depends on the matrix size as shown on a 16 cores execution in Figure 2(a). Indeed, if the matrix is small, it needs to be cut in even smaller pieces to provide work to all the 16 cores even if this induces that the serial kernels individually achieve a lower performance. When the matrix size increases, all the cores may evenly share the work using a larger tile size and thus achieving a higher performance. In a nutshell, the optimum tile size both depends on the number of cores and the matrix size, and its choice is critical for performance. Figure 2(b) shows that the impact is even stronger on a 32 cores IBM Power6 machine. The 80-40 combination is optimum on a matrix of order 500 but only achieves 6.3% of the optimum (20.6 Gflop/s against 325.9 Gflop/s) on a matrix of order 12,000.

3 Two-Step Empirical Method

Given the considerations discussed in introduction and further developed in [10], we do *not* propose a model-driven tuning approach. Instead we use a fully empirical method that effectively executes the factorizations on the target platform. However, not all NB-IB combinations can be explored. Indeed, an *exhaustive search* is cumbersome since the search space is huge. For instance, there are more than 1000 possible NB-IB combinations even if we constrain NB to be an even integer lower than 512 (size where the single core compute-intensive kernel

reaches its asymptotic performance) and if we impose IB to divide NB. Exploring this search space on a matrix of order $N = 10,000$ with 8 cores on the Intel Core Tigerton machine (described in Section 4) would take several days. Therefore, we need to prune the search space. We propose a two-step approach. In Step 1 (Section 5), we benchmark the most compute-intensive serial kernel. This step is fast since the serial kernels operate on tiles, which are of small granularity ($NB < 512$) compared to the matrices to be factorized ($500 \leq N \leq 10000$ in our study). Thanks to this collected data set and a few well chosen empirical properties, we pre-select (PS) a subset of NB-IB combinations. We propose three heuristics for performing that preliminary pruning automatically. In step 2 (Section 6) we benchmark the effective multicore QR factorizations on the pre-selected set of NB-IB combinations. We furthermore show that further pruning (PAYG) can be performed during this step, drastically reducing the whole tuning process.

4 Experimental Environments

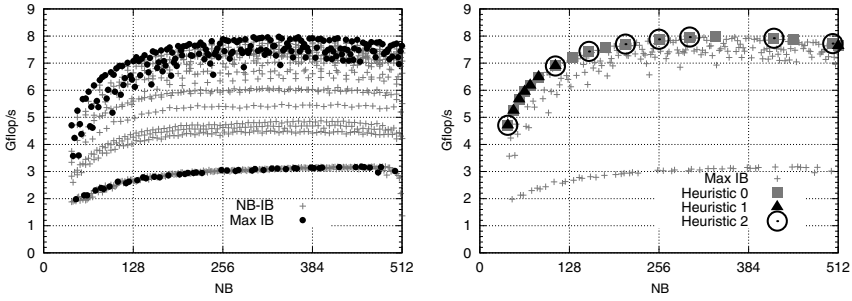
To assess the portability and reliability of our method, we consider seven platforms based Intel EM64T processors, IBM Power and AMD x86_64. **Intel Core Tigerton.** This 16 cores machine is a quad-socket quad-core Xeon E7340 (codename Tigerton) processor, an Intel Core micro-architecture. The processor operates at 2.39 GHz. **Intel Core Clovertown.** This 8 cores server is another machine based on an Intel Core micro-architecture. The machine is composed of two quad-core Xeon X5355 (codename Clovertown) processors, operating at 2.66 GHz. **Intel Core Yorkfield.** This 4 cores desktop is also based on an Intel Core micro-architecture. The machine is composed of one Core 2 Quad Q9300 (codename Yorkfield) processor, operating at 2.5 GHz. **Intel Core Conroe.** This 2 cores desktop is based on an Intel Core micro-architecture too. The machine is composed of one Core 2 Duo E6550 (codename Conroe) processors, operating at 2.33 GHz. **Intel Nehalem.** This 8 cores machine is based on an Intel Nehalem micro-architecture. Instead of having one bank of memory for all processors as in the case of the Intel Core's architecture, each Nehalem processor has its own memory. Nehalem is thus a cache coherent Non Uniform Memory Access (ccNUMA) architecture. Our machine is a dual-socket quad-core Xeon X5570 (codename Gainestown) running at 2.93GHz and up to 3.33 GHz in certain conditions (Intel Turbo Boost technology). The Turbo Boost was activated during our experiments. **AMD Istanbul.** This 48 cores machine is composed of eight hexa-core Opteron 8439 SE (codename Istanbul) processors running at 2.8 GHz. Like the Intel Nehalem, the Istanbul micro-architecture is a ccNUMA architecture. **IBM Power6.** This 32 cores machine is composed of sixteen dual-core IBM Power6 processors running at 4.7 GHz. More details on these platforms can be found in [10]. Note that we do not discuss here the mapping of the tasks onto the cores; this is an orthogonal problem.

5 Step 1: Benchmarking the Most Compute-Intensive Serial Kernel

We explained in Section 2.1 that the tile QR factorization consists of four serial kernels. However, the number of calls to DSSRFB is proportional to NT^3 while the number of calls to the other kernels is only proportional to NT (DGEQRT) or to NT^2 (DTSQRT and DLARFB). Even on small DAGS (see Figure 1(b)), calls to DSSRFB are predominant. Therefore, the performance of this compute-intensive kernel is crucial. DSSRFB’s performance also depends on NB-IB. It is thus natural to pre-select NB-IB pairs that allow a good performance of DSSRFB before benchmarking the QR factorization itself. The practical advantage is that a kernel is applied at the granularity of a tile, which we assume to be bounded by 512 ($NB \leq 512$). Consequently, preliminary benchmarking this serial kernel can be done exhaustively in a reasonable time. *Step 1* thus consists of performing an exhaustive benchmarking of the DSSRFB kernel on all possible NB-IB combinations and then to decide which of these will be kept for further testing in Step 2. Column “Step 1” of Table 1 shows that the total elapsed time for step 1 is acceptable on all the considered architectures (between 16 and 35 minutes). Figure 3(a) shows the resulting set of empirical data collected during step 1 on the Intel Core Tigerton machine. This data set can be pruned a

Table 1. Elapsed time (hh:mm:ss) for Step 1 and Step 2

Machine		Step 1	Step 2	
Architecture	# cores		Heuristic	PS PSPAYG
Conroe	2	00:24:33	0	14:46:37 03:05:41
			1	09:01:08 00:01:58
			2	07:30:53 00:34:47
Yorkfield	4	00:20:57	0	17:40:00 04:48:13
			1	09:30:30 00:05:10
			2	08:01:05 02:58:37
Clovertown	8	00:21:44	0	20:08:43 02:56:25
			1	11:06:18 00:13:09
			2	08:52:24 01:10:53
Nehalem	8	00:16:29	0	06:20:16 01:51:30
			1	06:20:16 01:51:30
			2	06:20:16 01:51:30
Tigerton	16	00:34:18	0	23:29:35 03:15:41
			1	12:22:06 00:08:57
			2	09:54:59 01:01:06
Istanbul	48	00:24:23	0	21:09:27 02:53:38
			1	12:25:30 00:11:01
			2	10:04:46 00:54:51
Power6	32	00:15:23	0	03:06:05 00:25:07
			1	03:06:05 00:25:07
			2	03:06:05 00:25:07



(a) Different NB-IB combinations with (b) Combinations pre-selected (PS) for a common NB value have the same each heuristic. absceisse; “Max IB” represents the one that achieves the maximum performance among them.

Fig. 3. Performance of the DSSRFB serial kernel depending on the NB-IB combination

first time. Indeed, contrary to NB, which trades off parallelism for kernel performance, IB only affects kernel performance but not parallelism. We can thus perform the following orthogonal optimization:

Property 1 (Orthogonal pruning). *For a given NB value, we can safely pre-select the value of IB that maximizes the kernel performance.*

Applying Property 1 to the data set of Figure 3(a) results in discarding all NB-IB pairs except the ones matching “Max IB”, which still represents a large number of combinations. We thus propose and assess three heuristics to further prune the search space. The first considered heuristic is based on the fact that intensive experiments (not reported here) showed the following property.

Property 2 (Convex Hull). *There is consistently an optimum combination on the convex hull of the data set.*

Therefore, **Heuristic 0** consists of pre-selecting the points from the convex hull of the data set (see Figure 3(b)). In general, this approach may still provide too many combinations. Because NB trades off kernel efficiency with parallelism, the gains observed on kernel efficiency shall be considered relatively to the increase of NB itself. Therefore, we implemented **Heuristic 1** that pre-selects the points of the convex hull with a high steepness (or more accurately a point after a segment with a high steepness). The drawback is that all these points tend to be located in the same area as shown in Figure 3(b) corresponding to small values of NB. To correct this deficiency, we consider **Heuristic 2** which first divides the x-axis into iso-segments and pick up the point of maximum steepness on each of these segments (see Figure 3(b) again). Heuristics 1 and 2 are parametrized to select a maximum of 8 combinations. All three heuristics perform a pre-selection (PS) that will be used as test cases for the second step.

6 Step 2: Benchmarking the Whole QR Factorization

6.1 Discretization and Interpolation

We recall that our objective is to immediately retrieve at execution time the optimum NB-IB combination for the matrix size N and number of cores $ncores$ that the user requests. Of course, N and $ncores$ are not known yet at install time. Therefore, the $(N, ncores)$ space to be benchmarked has to be discretized. We decided to benchmark all the powers of two cores (1, 2, 4, 8, ...) plus the maximum number of cores in case it is not a power of two such as on the AMD Istanbul machine. The motivation comes from empirical observation. Indeed, Figures 4(a) and 4(b) show that the optimum NB-IB combination can be finely interpolated with such a distribution. We discretized more regularly the space on N because the choice of the optimum pair is much more sensible to that dimension (see figures 2(a) and 2(b)). We benchmarked $N=500, 1000, 2000, 4000, 6000, 8000, 10000$ ¹. Each run is performed 6 times to attenuate potential perturbations. When the user requests the factorization of parameters that have not been tuned (for instance $N=1800$ and $ncores=5$) we simply interpolate by selecting the parameters of the closest configuration benchmarked at install time ($N=2000$ and $ncores=4$ in that case).

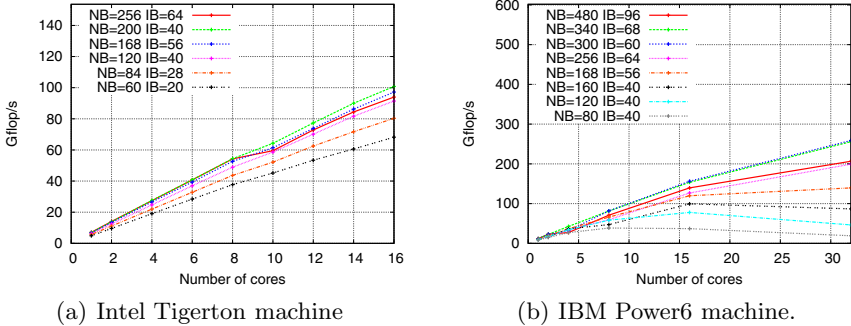


Fig. 4. Strong scalability - $N = 6000$

6.2 Impact of the Pre-selection on the Elapsed Time of Step 2

Column PS (pre-selection) in Table 1 shows the impact of the heuristics on the time required for benchmarking step 2. Clearly Heuristic 0 induces a very long step 2 (up to 1 day). Heuristic 1 and 2 induce a lower time for step 2 (about 10 hours) but that may be still not acceptable for many users.

¹ Except on the IBM Power6 machine where $N=10000$ was not benchmarked.

6.3 Prune as You Go (PSPAYG)

To further shorten step 2, we can perform complementary pruning on the fly. Indeed, Figures 2(a) and 2(b) show the following property.

Property 3 (Monotony). *Let us denote by $P(NB_1, N)$ and $P(NB_2, N)$ the performances obtained on a matrix of order N with tile sizes NB_1 and NB_2 , respectively. If $P(NB_1, N) > P(NB_2, N)$ and $NB_1 > NB_2$, then $P(NB_1, N') > P(NB_2, N')$ for any $N' > N$.*

We perform step 2 in increasing order of N . After having benchmarked the current set of NB-IB combinations on a matrix of order N , we identify all the couples (NB_1, NB_2) that satisfy Property 3 and we remove from the current subset the NB-IB pair in which NB_2 is involved. Indeed, according to Property 3, it would lead to a lower performance than NB_1 on larger values of N which are going to be explored next. We denote this strategy by “PSPAYG” (pre-selection and prune as you go). Column PSPAYG in Table 1 shows that the time for step 2 is dramatically improved with this technique. Indeed, the number of pairs to explore decreases when N increases, that is, when benchmark is costly. For heuristic 2 (values in bold in Table 1), the time required for step 2 is reduced by a factor greater than 10 in two cases (Intel Core Conroe and AMD Istanbul machines).

6.4 Reliability

We employed the following methodology to assess the reliability of the different tuning approaches. We first executed all the discussed approaches on all the platforms with the discretization of the $(N, ncores)$ space proposed in Section 6.1. We then picked up between 8 and 16 $(N, ncores)$ combinations such that half of them were part of the discretized space (for instance $N = 6000$ and $ncores = 32$) and the other half were not part of it (for instance $N = 4200$ and $ncores = 30$) so that the reliability of the interpolation is also taken into account. For each combination we performed an (almost) exhaustive search for reference. Table 2 provides a synthesis of the results. Heuristic 2 coupled with the PSPAYG approach is very efficient since it achieves a high proportion of the performance that would be obtained with an exhaustive search (values in bold). The worst case occurs on the Istanbul machine, with an average relative performance of 97.1% (Column “avg”). However, even on that platform, the optimum NB-IB combination was found in seven cases out of sixteen tests (Column “optimum”).

Column $\frac{PSPAYG}{PS}$ allows to specifically assess the impact of the “prune as you go” method since they compare the average performance obtained with PSPAYG (where pairs can be discarded during step 2 according to Property 3) compared to PS (where no pair is discarded during step 2). The result is clear: pruning during step 2 according to Property 3 does not hurt performance ($\frac{|PS - PSPAYG|}{PS} < 0.3\%$), showing that Property 3 is strongly reliable. Finally, note that on $(N, ncores)$ combinations part of the discretized space, PSPAYG cannot achieve a higher performance than PS since all NB-IB combinations

Table 2. Average performance achieved with a “pre-selection” (PS) method or a “pre-selection and prune as you go” (PSPAYG) method, based on different heuristics (H) applied at step 1. The performance is presented as a proportion of the exhaustive search (ES) or of the prunes search (PS). The column “optimum” indicates the number of times the optimum combination (with respect to the reference method) was found among the number of tests performed.

Machine	H	$\frac{PS}{ES}$ (%)		$\frac{PSPAYG}{ES}$ (%)		$\frac{PSPAYG}{PS}$ (%)	
		avg optimum	optimum	avg optimum	optimum	avg optimum	optimum
Conroe	0	99.67	6/8	99.67	6/8	100	8/8
	1	95.28	0/8	95.28	0/8	100	8/8
	2	99.54	5/8	99.54	5/8	100	8/8
Yorkfield	0	98.63	6/12	98.63	6/12	100	12/12
	1	91.53	0/12	91.59	0/12	100.07	10/12
	2	98.63	6/12	98.63	6/12	100	12/12
Clovertown	0	98.59	8/16	98.35	7/16	99.76	15/16
	1	91.83	0/16	91.83	0/16	100	16/16
	2	98.49	9/16	98.25	8/16	99.76	15/16
Nehalem	0	98.6	8/16	98.9	8/16	100.33	16/16
	1	98.6	8/16	98.9	8/16	100.33	16/16
	2	98.6	8/16	98.9	8/16	100.33	16/16
Tigerton	0	97.36	8/16	97.54	5/16	100.21	12/16
	1	91.61	0/16	91.61	0/16	100	16/16
	2	97.51	8/16	97.79	7/16	100.31	15/16
Istanbul	0	97.17	7/16	97.17	7/16	100	16/16
	1	94.12	2/16	94.12	2/16	100	16/16
	2	97.23	7/16	97.1	7/16	99.87	15/16
Power 6	0	100	16/16	100	16/16	100	16/16
	1	100	16/16	100	16/16	100	16/16
	2	100	16/16	100	16/16	100	16/16

tested with PSPAYG are also tested with PS. However, PSPAYG can achieve a higher performance if ($N, ncores$) was not part of the discretized space because of the interpolation. This is why cases where $\frac{PSPAYG}{PS} > 100\%$ may be observed.

7 Conclusion and Future Work

We have presented a new fully empirical autotuned method for tuning dense linear algebra libraries on multicore architectures. Thanks to three strong empirical properties, we showed that the search space can be efficiently pruned. Our tuning process is automatic, fast (less than one hour and ten minutes on five out of seven platforms) and reliable (average performance varying from 97% to 100% of the optimum). We plan to extend our work to the case of non square matrices and to other factorizations. We will then extend our work to the case of hybrid multicore platforms enhanced with multiple GPU accelerators for which heterogeneity will have to be taken into account.

Acknowledgment. The authors would like to thank Jakub Kurzak, Greg Henry and Clint Whaley for their constructive discussions.

References

1. Frigo, M., Johnson, S.: FFTW: An adaptive software architecture for the FFT. In: Proc. 1998 IEEE Intl. Conf. Acoustics Speech and Signal Processing, vol. 3, pp. 1381–1384. IEEE, Los Alamitos (1998)
2. Choi, J.W., Singh, A., Vuduc, R.W.: Model-driven autotuning of sparse matrix-vector multiply on GPUs. In: Proc. ACM SIGPLAN Symp. Principles and Practice of Parallel Programming (PPoPP), Bangalore, India (January 2010)
3. Ansel, J., Chan, C., Wong, Y.L., Olszewski, M., Zhao, Q., Edelman, A., Amarasinghe, S.: Petabricks: A language and compiler for algorithmic choice. In: ACM SIGPLAN Conference on Programming Language Design and Implementation, Dublin, Ireland (June 2009)
4. Clint Whaley, R., Petitet, A., Dongarra, J.J.: Automated empirical optimizations of software and the atlas project. *Parallel Computing* 27(1-2), 3–35 (2001)
5. Volkov, V., Demmel, J.W.: Benchmarking gpus to tune dense linear algebra. In: SC 2008: Proceedings of the ACM/IEEE Conference on Supercomputing, pp. 1–11. IEEE Press, Piscataway (2008)
6. Tomov, S., Nath, R., Ltaief, H., Dongarra, J.: Dense linear algebra solvers for multicore with gpu accelerators. Accepted for publication at HIPS 2010 (2010)
7. Quintana-Ortí, G., Quintana-Ortí, E., van de Geijn, R., Van Zee, F., Chan, E.: Programming matrix algorithms-by-blocks for thread-level parallelism. *ACM Trans. Math. Softw.* 36(3) (2009)
8. Buttari, A., Langou, J., Kurzak, J., Dongarra, J.: A class of parallel tiled linear algebra algorithms for multicore architectures. *Parallel Computing* 35(1), 38–53 (2009)
9. Agullo, E., Hadri, B., Ltaief, H., Dongarra, J.: Comparative study of one-sided factorizations with multiple software packages on multi-core hardware. In: 2009 International Conference for High Performance Computing, Networking, Storage, and Analysis (SC 2009) (2009)
10. Agullo, E., Dongarra, J., Nath, R., Tomov, S.: A Fully Empirical Autotuned Dense QR Factorization For Multicore Architectures. Research Report 7526, INRIA (February 2011)