



## Exploring New Search Algorithms and Hardware for Phylogenetics: RAxML Meets the IBM Cell

A. STAMATAKIS

*School of Computer and Communication Sciences, École Polytechnique Fédérale de Lausanne,  
Lausanne, Switzerland*

F. BLAGOJEVIC AND D. S. NIKOLOPOULOS

*Department of Computer Science, Center for High-end Computing Systems, Virginia Tech,  
Blacksburg, VA, USA*

C. D. ANTONOPOULOS

*Division of Research and Informatics, Greek Armed Forces, Thessaloniki, Greece*

*Received: 8 December 2006; Revised: 00 Month 0000; Accepted: 17 March 2007*

**Abstract.** Phylogenetic inference is considered to be one of the *grand challenges* in Bioinformatics due to the immense computational requirements. RAxML is currently among the fastest and most accurate programs for phylogenetic tree inference under the Maximum Likelihood (ML) criterion. First, we introduce new tree search heuristics that accelerate RAxML by a factor of 2.43 while returning equally good trees. The performance of the new search algorithm has been assessed on 18 real-world datasets comprising 148 up to 4,843 DNA sequences. We then present the implementation, optimization, and evaluation of RAxML on the IBM Cell Broadband Engine. We address the problems and provide solutions pertaining to the optimization of floating point code, control flow, communication, and scheduling of multi-level parallelism on the Cell.

**Keywords:** phylogenetic inference, maximum likelihood, RAxML, IBM cell

### 1. Introduction

Phylogenetic trees are used to represent the evolutionary history of a set of  $n$  organisms. An alignment with the DNA or protein sequences representing  $n$  organisms can be used as input for phylogenetic inference. In a phylogeny the organisms of the input data are located at the tips (leaves) of the tree and the inner nodes represent extinct common ancestors. The branches of the tree represent the time which was required for the mutation of one species into another, new one. Phylogenetic trees have many important applications in medical and biological research (see [2] for a summary).

Due to the rapid growth of sequence data over the last years it has become feasible to compute large trees which often comprise more than 1,000 organisms and sequence data from several genes (so-called multi-gene alignments). This means that alignments grow in the number of organisms *and* in sequence length.

The basic algorithmic problem computational phylogeny faces is the immense amount of alternative tree topologies which grows exponentially with the number of organisms  $n$ , e.g. for  $n = 50$  there exist  $2.84 * 10^{76}$  alternative trees. In fact, it has only recently been shown that the ML phylogeny problem is NP-hard [6]. In addition, ML-based inference of

52 phylogenies is very memory- and floating point-  
 53 intensive. Since phyloinformatics has definitely  
 54 entered the HPC era by now, the application of high  
 55 performance computing techniques as well as the  
 56 assessment of new CPU architectures can significant-  
 57 ly contribute to the reconstruction of larger and more  
 58 accurate trees. Moreover, typical ML implementa-  
 59 tions exhibit different levels of parallelism and are  
 60 thus well-suited as example applications for exploit-  
 61 ing unconventional and challenging architectures.

62 The Cell Broadband Engine (BE) has been  
 63 developed jointly by Sony, Toshiba and IBM.  
 64 Although the Cell was originally designed for the  
 65 set-top box market, it has evolved into a general-  
 66 purpose processor for high-performance computing,  
 67 server and desktop applications. The Cell BE is a  
 68 heterogeneous multicore processor with nine pro-  
 69 cessing cores: one two-way multithreaded PowerPC  
 70 Processor Element (PPE) and eight Synergistic  
 71 Processor Elements (SPEs). The Cell is well suited  
 72 for data-intensive scientific applications with high  
 73 demand for memory bandwidth. It offers a unique  
 74 assembly of MIMD and SIMD execution capabilities  
 75 and a software-managed memory hierarchy, thus  
 76 providing ample flexibility in selecting programming  
 77 and parallelization models for a given application.

78 According to its specifications [1], the Cell is  
 79 capable of achieving significant performance im-  
 80 provement over conventional multicore CPUs, in-  
 81 cluding an improved Flops per Watt ratio. However,  
 82 due to its unconventional architecture, the develop-  
 83 ment of parallel applications that can exploit all  
 84 advantages of the Cell design, is a challenging task.  
 85 One of the main difficulties is the management of the  
 86 local storage of the SPEs by software. Another  
 87 challenge lies in the distribution of work among  
 88 SPEs, which can be implemented at multiple degrees  
 89 of granularity and with a variety of code and data  
 90 distribution schemes. This paper addresses these  
 91 problems in the realm of algorithms for phylogenetics  
 92 and more specifically RAxML [25]. We present  
 93 optimizations and system software support for vecto-  
 94 rization, control flow parallelization, scheduling, and  
 95 communication on the Cell, along with a stepwise  
 96 evaluation of these optimizations in RAxML. The  
 97 optimizations are generic enough to be reused across  
 98 a wider range of parallel applications.

99 The remainder of this paper is organized as  
 100 follows: First, we review related work on IBM Cell  
 101 portings and tools (Section 2). In Section 3 we

102 provide an overview of RAxML and related work in  
 103 the area of phylogenetics. We also introduce new  
 104 tree search heuristics that accelerate RAxML by  
 105 factor 2.43 while returning equally good trees on real-  
 106 world datasets of 148 up to 4,843 DNA sequences.  
 107 The following Section 4 gives a step-by-step descrip-  
 108 tion of porting and optimizing RAxML on Cell. The  
 109 conventional optimizations we used to speedup the  
 110 execution time include the use of optimized numer-  
 111 ical libraries for the SPEs, double-buffering for  
 112 complete communication/computation overlap, vec-  
 113 torization of floating point operations, and minimi-  
 114 zation of the ratio of PPE-SPE computation via  
 115 function off-loading. We also present new Cell-  
 116 specific optimizations, including the vectorization of  
 117 conditional statements, asynchronous communica-  
 118 tion via direct SPE memory accesses and an event-  
 119 driven scheduling model, which can selectively  
 120 exploit coarse-grain and fine-grain parallelism, in  
 121 response to workload variation. In Section 5 we  
 122 evaluate performance of RAxML on the IBM Cell  
 123 and the IBM Power5 multicore processor. We  
 124 conclude with Section 6.

**2. Related Work on IBM Cell**

125 Recently, several studies were conducted to measure  
 126 performance and develop programming models for  
 127 easier porting as well as optimization of parallel  
 128 applications on the Cell.  
 129

130 Fatahalian et al. [10] developed Sequoia, a program-  
 131 ming language suitable for porting memory-aware  
 132 applications to machines with different memory  
 133 hierarchy configurations. One of the target architec-  
 134 tures in this study is the Cell Broadband Engine. The  
 135 authors used Sequoia to port several programs to Cell  
 136 and obtained memory throughput exceeding 20 GB/s.  
 137 Sequoia currently supports applications which can be  
 138 parallelized via recursive block decomposition,  
 139 whereas our work focuses on less structured applica-  
 140 tions that can be parallelized at multiple layers.

141 Bellens et al. [3] developed a dependence-driven  
 142 programming model for porting sequential applica-  
 143 tions to Cell. Their compiler is capable of generating  
 144 tasks that can potentially be executed in parallel. The  
 145 supporting runtime system creates a dependence  
 146 graph of active tasks during execution and deter-  
 147 mines which tasks can be scheduled for execution on  
 148 the SPEs. This work considers only one layer of  
 149 task-level parallelism and does not explore the

150 implications of task size and available parallelism  
 151 within tasks. Both issues are central in the current  
 152 paper.

153 Kunzman et al. [17] are in the process of adapting  
 154 Charm++ on Cell. Charm++ is a runtime system for  
 155 object-based parallel programming. More specifically,  
 156 Charm++ is a library of machine-independent object  
 157 abstractions for scheduling and communication on  
 158 parallel machines, implemented on both distributed  
 159 and shared memory systems. No results that would  
 160 enable comparisons with our work are reported as of the  
 161 time of writing this paper.

162 Although Cell has been a focal point in numerous  
 163 articles in popular press, published research using  
 164 Cell for real-world scientific applications beyond  
 165 games and streaming computation is scarce. Besides  
 166 our study, Hjelte [15] presents an implementation of  
 167 a smooth particle hydrodynamics simulation on Cell.  
 168 This simulation requires good interactive perfor-  
 169 mance, since it lies on the critical path of real-time  
 170 applications such as interactive simulation of human  
 171 organ tissue, body fluids, and vehicular traffic.  
 172 Benthin et al. [4] present a parallel ray-tracing  
 173 algorithm for Cell.

174 **3. The RAxML Application**

175 Despite the high computational cost, significant  
 176 progress has been achieved over the last years in  
 177 the field of heuristic ML search algorithms with the  
 178 release of programs such as IQPNNI [19], PHYML  
 179 [14], GARLI [31] and RAxML [25, 27] to name only  
 180 a few.

181 RAxML-VI-HPC (v2.2.0; Randomized Axelerated  
 182 Maximum Likelihood version VI for High Performance  
 183 Computing, freely available at [http://icwww.epfl.ch/  
 184 ~stamatak](http://icwww.epfl.ch/~stamatak)) [25] is a program for large-scale ML-based  
 185 [12] inference of evolutionary trees using multiple  
 186 alignments of DNA or AA (Amino Acid) sequences.  
 187 Some of the largest published ML-based phylogenetic  
 188 analyses to date have been conducted with RAxML  
 189 [11, 13, 21]. To the best of our knowledge, RAxML  
 190 has been used to compute trees on the two largest data  
 191 matrices analyzed under ML to date: a 25,057-taxon  
 192 alignment of protobacteria (length: 1,463 nucleotides)  
 193 and a 2,182-taxon alignment of mammals (length:  
 194 51,089 nucleotides).

195 The current version (v2.2.0) implements the new  
 196 search algorithm described in Section 3.2 as well as  
 197 the older algorithm (command line switch *phv - fo*)

outlined in [25]. A recent performance study [25] on 198  
 real world datasets with more than 1,000 sequences 199  
 reveals that it is able to find better trees in less time 200  
 and with lower memory consumption than other 201  
 current ML programs (IQPNNI, PHYML, GARLI). 202

3.1. Parallel Implementations of ML Programs 203

RAxML exploits two levels of parallelism: fine-grain 204  
 loop-level parallelism and coarse-grain embarrassing 205  
 parallelism. 206

RAxML has been parallelized with OpenMP to 207  
 exploit loop-level parallelism. Like every ML-based 208  
 program, RAxML exhibits a source of loop-level 209  
 parallelism in the likelihood functions which typi- 210  
 cally consume over 90% of the overall computation 211  
 time. The OpenMP implementation scales particu- 212  
 larly well on large multi-gene alignments due to 213  
 increased cache efficiency [28]. 214

The MPI version of RAxML exploits the embar- 215  
 rassing parallelism that is inherent to every real- 216  
 world phylogenetic analysis. In order to conduct 217  
 such an analysis (see [13] for an example) a number 218  
 of about 20–200 distinct tree searches (multiple 219  
 inferences) to find a best-scoring tree on the original 220  
 alignment as well as a large amount of 100–1,000 221  
 bootstrap analyses have to be conducted. *Bootstrap* 222  
*Analyses* are required to assign confidence values 223  
 ranging between 0.0 and 1.0 to the inner nodes of the 224  
 best-known ML tree. This allows to determine how 225  
 well-supported certain parts of the tree are and is 226  
 important to draw biological conclusions. Bootstrap- 227  
 ping is essentially very similar to multiple inferen- 228  
 ces. The only difference is that inferences are 229  
 conducted on a randomly re-sampled alignment (a 230  
 certain amount of alignment columns is re-weighted) 231  
 for every bootstrap run. This is performed in order to 232  
 assess the topological stability of the tree under 233  
 slight alterations of the input data. 234

All those individual tree searches, be it bootstrap or 235  
 multiple inferences are completely independent from 236  
 each other and can thus be exploited by a simple 237  
 master-worker scheme. If the dataset is not extremely 238  
 large, this represents the most efficient approach to 239  
 exploit HPC platforms for production runs. 240

Most other parallel implementations of ML pro- 241  
 grams [7, 19, 26, 29, 31] have mainly focused on the 242  
 intermediate level of parallelism (inference parallel- 243  
 ism) which is situated between the loop-level 244  
 parallelism and coarse-grained parallelism currently 245

246 exploited in RAxML. The related work mentioned  
 247 above mainly deals with highly algorithm-specific  
 248 and mostly MPI-based parallelization of various hill-  
 249 climbing, genetic, as well as divide-and-conquer  
 250 search algorithms. Finally, Minh et al. [20] recently  
 251 implemented a hybrid OpenMP/MPI version of  
 252 IQPNNI which exploits loop-level and inference  
 253 parallelism.

254 **3.2. Accelerating the Search Algorithm**

255 This Section describes a novel heuristic optimization  
 256 of the RAxML search algorithm that is available as  
 257 of version 2.2.0. The entire procedure is outlined in  
 258 Fig. 1.

259 The fundamental mechanism that is used to search  
 260 the tree space with RAxML is called Lazy Subtree  
 261 Rearrangement (LSR, for details see [27]). An LSR  
 262 consists in pruning/removing a subtree from the  
 263 currently best tree  $t$  and then re-inserting it into all  
 264 neighboring branches up to a certain distance/radius  
 265 (rearrangement distance) of  $n$  nodes from the pruning  
 266 point ( $n$  typically ranges from 5 to 25). For each  
 267 possible subtree insertion within the rearrangement  
 268 distance, RAxML evaluates the log likelihood score  
 269 of the alternative topology. This is done in a lazy  
 270 way since only the length of the three branches  
 271 adjacent to the insertion point/node will be opti-  
 272 mized. Thus, an LSR only yields an approximate log  
 273 likelihood  $all(t')$  score for each alternative topology

274  $t'$  constructed by an LSR from  $t$ . However, this  $all(t')$   
 275 score can be used to sort the alternative topologies.  
 276 After this fast pre-scoring of a large number of  
 277 alternative topologies, only a very small fraction of  
 278 the best-scoring topologies needs to be optimized  
 279 more exhaustively to improve the overall tree score.  
 280 One iteration of the RAxML hill-climbing algorithm  
 281 consists in performing LSRs on all subtrees of a  
 282 given topology  $t$  for a fixed rearrangement distance  $n$ .  
 283 Thereafter, the branches of the 20 best-scoring trees  
 284 are thoroughly optimized. This procedure of con-  
 285 ducting LSRs on all subtrees and then optimizing the  
 286 20 best-scoring trees, is performed until no improved  
 287 tree is encountered.

288 The main idea of the new heuristics is to reduce  
 289 the number of LSRs performed. This is done by  
 290 using an empirical cutoff-rule that stops the recursive  
 291 descent of an LSR into deeper branches at a higher  
 292 rearrangement distance from the pruning position if  
 293 they do not appear to be promising. Thus, if the  
 294 approximate log likelihood  $all(t')$  for the current  
 295 rearranged tree  $t'$  is worse than the log likelihood  
 296  $ll(t)$  of the currently best tree  $t$  and if the difference  
 297  $\delta(all(t'), ll(t))$  is larger than a certain threshold  
 298  $lh_{cutoff}$  the remaining LSRs below that node are  
 299 omitted. The threshold  $lh_{cutoff}$  is determined as  
 300 follows: During the first iteration of the RAxML  
 301 search algorithm  $lh_{cutoff} = \infty$  which means that no  
 302 cutoffs are made. In the course of this first iteration,  
 303 the differences  $\delta_i(all(t_i), ll(t))$  for all those  $i = 1...m$

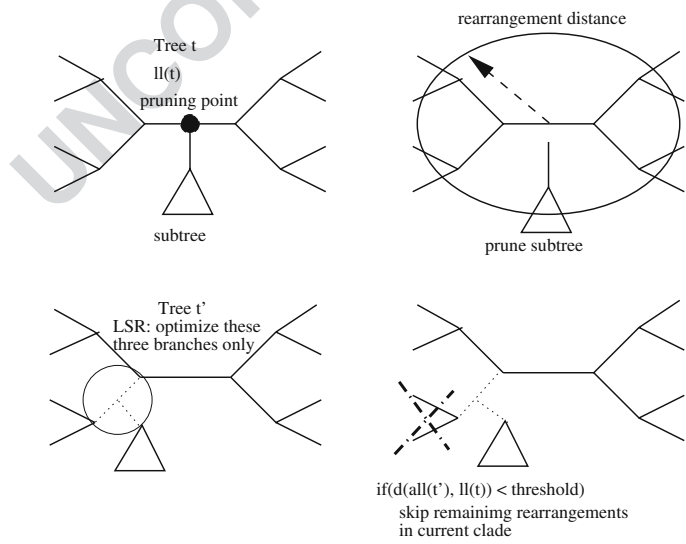


Figure 1. Outline of lazy subtree rearrangements with cutoff procedure.

304 alternative tree topologies  $t_i$  where  $all(t_i) \leq ll(t)$  are  
 305 stored. The threshold  $lh_{cutoff}$  for the next iteration is  
 306 set to the average of  $\delta_i$ , i.e.  $lh_{cutoff} = (\sum_{i=1}^m \delta_i)/m$ . If  
 307 the search computes an LSR for which  $all(t') \leq ll(t)$   
 308 and  $\delta(all(t'), ll(t)) \geq lh_{cutoff}$  it will simply skip the  
 309 remaining LSRs below the current node. Thus, each  
 310 iteration  $k$  of the search algorithm uses a threshold  
 311 value  $lh_{cutoff}$  that has been obtained during the  
 312 previous iteration  $k - 1$ . This allows to dynamically  
 313 adapt  $lh_{cutoff}$  to the specific dataset and to the  
 314 progress of the search. The omission of a large  
 315 amount of unnecessary LSRs that have a high  
 316 probability not to improve the tree yields substantial  
 317 run time improvements and returns equally good  
 318 trees at the same time (see Table 1).

319 **3.2.1. Results.** To assess performance of the new  
 320 heuristics we analyzed 18 real-world datasets com-  
 321 prising 148 to 4,843 sequences from various sources  
 322 (see Acknowledgment). The computations were  
 323 performed under the CAT approximation of rate  
 324 heterogeneity [24], but final tree scores were evalu-  
 325 ated with the standard GAMMA model of rate  
 326 heterogeneity. For each dataset we generated ten  
 327 starting trees and executed the old and new RAxML  
 328 search algorithm on each of those starting trees. We  
 329 performed a total of 380 ML searches which were  
 330 executed on the Infiniband cluster at the TU  
 331 München. The cluster is equipped with 36 quad-  
 332 CPU 2.4 GHz AMD Opteron nodes.

333 Table 1 lists the average final log likelihood values  
 334 (LH-NEW, LH-OLD) for the old and new versions  
 335 of the RAxML search algorithm. In addition, it  
 336 provides the average speedup value per dataset. The  
 337 average speedup over all datasets is 2.43. The slight  
 338 variations in likelihood scores are insignificant.

**4. Porting and Optimizing RAxML on Cell** 339

*4.1. The Cell BE* 340

The Cell BE is a heterogeneous multicore processor. 341  
 The design is inspired by graphics accelerators, 342  
 which are commonly used in set-top boxes for data 343  
 streaming computations. In contrast to conventional 344  
 processors, graphics accelerators typically provide 345  
 vast memory bandwidth and vectorization capabili- 346  
 ties. Therefore, they are able to execute streaming 347  
 computation kernels such as encryption/decryption, 348  
 compression, FIR filters and FFTs far more efficient- 349  
 ly. The Cell BE integrates eight specialized accel- 350  
 erators with a conventional 64-bit host processor on 351  
 the same chip. The host processor is a PowerPC 352  
 SMT core, which runs Linux in a virtualized setting. 353  
 The PowerPC core communicates via a ring network 354  
 with eight accelerator cores, called Synergistic 355  
 Processor Elements, or SPEs. The SPEs and the 356  
 PPE are laid out in equal distances around a ring 357  
 called the Element Interconnect Bus (EIB). The EIB 358  
 can transmit up to 96 bytes per processor cycle, and 359  
 has a maximum bandwidth of 200 GB/s. The PPE is 360  
 also connected to the EIB. An overview of the Cell 361  
 architecture is provided in Fig. 2. 362

The SPEs can issue up to two instructions per 363  
 cycle, one of which can be a 128-bit SIMD 364  
 instruction. Since the SPEs do not include hardware 365  
 for branch prediction they rely on software to predict 366  
 the outcome of branches. In its current implementa- 367  
 tion, the SPE pipeline is optimized for single- 368  
 precision floating-point vector operations, for which 369  
 the processor can sustain a maximum throughput of 370  
 over 230 Gflops. This is an artifact of earlier designs 371  
 of the processor for game consoles, in which double 372

t1.1 *Table 1.* Average final log likelihood values and speedups over ten runs on distinct MP starting trees for 18 real-world datasets.

t1.2	No. of Seq.	LH-NEW	LH-OLD	Speedup	No. of Seq.	LH-NEW	LH-OLD	Speedup
t1.3	148	-69,726.05	-69,725.55	1.83	150	-39,606.06	-39,606.06	2.00
t1.4	218	-134,195.92	-134,199.30	1.98	404	-156,151.42	-156,147.25	2.00
t1.5	498	-219,186.90	-219,186.90	2.21	500	-85,794.87	-85,794.87	2.11
t1.6	628	-50,940.81	-50,938.29	2.68	714	-148,543.21	-148,544.48	2.51
t1.7	994	-348,936.85	-348,936.85	2.34	1,288	-395,999.61	-395,999.61	1.92
t1.8	1,512	-273,435.30	-273,443.51	2.74	1,604	-167,398.16	-167,399.43	3.31
t1.9	1,780	-178,930.53	-178,925.02	2.95	1,908	-149,645.43	-149,645.43	2.89
t1.10	2,000	-364,916.18	-364,914.44	2.65	2,554	-318,488.01	-318,488.01	2.49
t1.11	4,114	-325,621.05	-325,620.54	2.70	4,843	-748,075.04	-748,067.23	2.48

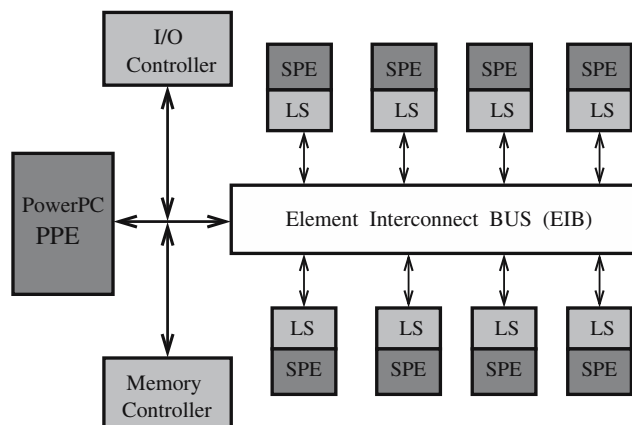


Figure 2. Outline of the CELL BE architecture.

373 precision floating-point calculations are of lesser impor-  
 374 tance than in scientific computing. Nevertheless, double-  
 375 precision throughput is still 21 Gflops, which is higher  
 376 than the throughput of most high-end homogeneous  
 377 multi-core processors. New versions of the processor for  
 378 high-performance computing resolve the bottleneck of  
 379 double-precision floating point operations.

380 Each SPE on the Cell has 128 128-bit general-  
 381 purpose registers and a software-managed fast local  
 382 storage. The use of a software-managed local storage is  
 383 unique in the Cell architecture. In the current version of  
 384 the processor, each SPE has 256 KB of local storage.  
 385 This storage is expandable to up to 4 MB. Local storage  
 386 is locally addressable via direct load and store  
 387 instructions from the owner SPE and globally accessi-  
 388 ble from other SPEs and the PPE via DMA requests.  
 389 Software control of the local storage enables the use of  
 390 customized cache replacement policies on the SPEs.  
 391 The local storage can also be programmed to operate as  
 392 a conventional set-associative cache [8].

393 The SPEs decouple processing from communica-  
 394 tion, via the use of a memory flow controller (MFC).  
 395 Each SPE can issue up to 16 concurrent DMA  
 396 requests, including requests for atomic DMA opera-  
 397 tions, which can be used in lieu of locks, and  
 398 requests for scatter-gather memory operations.

#### 399 4.2. RAXML Porting Strategy

400 We followed a three-step process for porting  
 401 RAXML to Cell:

- 402 1. We identified and off-loaded the compute inten-  
 403 sive parts of RAXML to the SPEs.

2. We optimized the off-loaded code on the SPEs. 404  
 We used optimized numerical libraries for the 405  
 SPEs, double buffering, vectorization of compu- 406  
 tation, vectorization of control flow, and PPE– 407  
 SPE communication optimizations. 408
3. Finally, we developed an event-driven scheduler 409  
 for RAXML’s nested parallel code, along with the 410  
 necessary system software support. 411

The first two steps follow the function off-loading 412  
 programming model, which arguably provides the 413  
 easiest path for porting MPI applications on the Cell 414  
 [8, 9]. The last step is a new contribution to support 415  
 the execution of parallel programs with multiple 416  
 levels of parallelism on the Cell. 417

In the following Sections we describe the porting 418  
 and optimization process, together with speedups 419  
 obtained at each optimization step. 420

#### 421 4.3. Porting the MPI Code

The version of RAXML that we ported on Cell is 422  
 based on the new, improved search heuristics, 423  
 described in Section 3.2. We initially executed the 424  
 MPI version of RAXML on the PPE. Since the PPE 425  
 is a dual-threaded processor, the PPE can execute 426  
 two MPI processes simultaneously. These processes 427  
 share the resources of the PPE, except from the 428  
 register file. Therefore, their parallel execution on 429  
 the PPE is not expected to scale as well as on a two- 430  
 way SMP. Although it is natural to consider using 431  
 the PPE for direct porting of MPI code without 432  
 modifications, this approach is clearly suboptimal. 433  
 Due to the heterogeneity of the Cell, further steps are 434

435 needed to off-load and optimize most, if not all, of  
436 the computation within each MPI process, in order to  
437 exploit the computational capacity of the SPEs.

438 We explored two strategies for off-loading code on  
439 the Cell SPEs. The first strategy is pure task-level  
440 off-loading. Each MPI process running on the PPE  
441 off-loads functions on the SPEs and each function  
442 executes from start to finish on the SPE it is assigned  
443 to. The second strategy is a hybrid task-level and  
444 loop-level parallelization scheme. Functions are off-  
445 loaded from the MPI process to some of the SPEs.  
446 The remaining SPEs are used for work sharing of  
447 parallel loops included in the already off-loaded  
448 functions. These parallelization strategies are de-  
449 scribed further in Section 4.7.

#### 450 4.4. Selecting Functions for Off-loading

451 In order to find the functions in RAxML that are  
452 suitable for SPE execution, we profiled the code  
453 using **gprof**. For all the profiling and benchmarking  
454 runs of RAxML presented in this paper, we use the  
455 input file *phv42\_SC*, which contains 42 organisms,  
456 each represented by a DNA sequence of 1,167  
457 nucleotides. The number of distinct data patterns in  
458 a DNA alignment is in the order of 250. Profiling the  
459 sequential version of RAxML on the IBM Power5  
460 processor shows that 96.6% of the execution time is  
461 spent in four functions:

- 462 • 64.62% of the time is spent in the function  
463 **newview()**. This function computes the partial  
464 likelihood vector [12] at an inner node of the  
465 phylogenetic tree.
- 466 • 26.88% of the time is spent in **makewz()**. This  
467 function optimizes the length of a given branch  
468 with respect to the tree likelihood using the  
469 Newton–Raphson method.
- 470 • 2.87% of the time is spent in **newviewPartial()**,  
471 which optimizes the per-site evolutionary rates for  
472 the GTRCAT approximation (see [24] for details).
- 473 • Finally, 2.29% of the time is spent in  
474 **phvevaluate()**. This function calculates the Log  
475 Likelihood score of the tree at a given branch by  
476 summing over the partial likelihood vector entries.  
477 Note that the log likelihood value is the same at all  
478 branches of the tree if the model of nucleotide  
479 substitution is time-reversible [12].

480 The prerequisite for computing **evaluate()** and  
481 **makewz()** is that the likelihood vectors at the

482 nodes to the right and left of the branch have been  
483 computed. Thus, **makewz()** and **evaluate()** ini-  
484 tially make calls to **newview()** before they can  
485 execute their own computation. The **newview()**  
486 function at an inner node **p** calls itself recursively  
487 when the two children **r** and **q** are not tips (leaves)  
488 and the likelihood array for **r** and **q** has not already  
489 been computed. Consequently, the first candidate  
490 for off-loading is **newview()**. Although **makewz()**  
491 and **evaluate()** are both taking a smaller portion  
492 of execution time than **newview()**, off-loading  
493 these two functions can also lead to significant  
494 speedup (see Section 4.6). Besides the fact that  
495 each function can be optimized and executed faster  
496 on an SPE, having all four functions off-loaded to  
497 an SPE significantly reduces the amount of PPE–  
498 SPE communication.

499 The code executed on the SPE has to be compiled  
500 separately (using the SPE specific compiler) from the  
501 code executed on the PPE. We choose to include all  
502 four off-loaded functions in one code module. This  
503 approach has the advantage of having all off-loaded  
504 functions in the SPE’s local storage during the entire  
505 execution of the program. Consequently, an off-  
506 loaded function can be invoked without introducing  
507 the overhead of moving its code from global memory  
508 to local storage. However, this decision imposes a  
509 trade-off, since the extended code segment in the  
510 SPE local storage reduces the available space for the  
511 heap and stack segments. In the case of RAxML, the  
512 total size of all four off-loaded functions is 110 KB.  
513 The remaining space (146 KB) is large enough to  
514 store the heap and stack segments of each of the four  
515 functions, as well as the buffers needed for commu-  
516 nication–computation overlap.

517 To keep the implementation simple, the call to  
518 each off-loaded function in the original MPI code is  
519 executed with the same signature on the PPE in the  
520 Cell code. We replaced the original body of each  
521 function with communication code needed to transfer  
522 local data used in the function from the PPE to SPEs.  
523 Whenever an off-loaded function is called, the PPE  
524 sends a signal to the SPE thread and waits for the  
525 SPE thread to complete the function and return the  
526 result. While waiting for the SPE code to finish, the  
527 PPE seeks functions for off-loading from other MPI  
528 processes. This process is described in more detail in  
529 Section 4.7.

530 All four off-loaded functions are executed inside a  
531 single SPE thread. The SPE thread is created at the

532 beginning of the program and stays alive during the  
 533 entire program execution. Depending on the content  
 534 of a triggering signal received from the PPE, the  
 535 thread executes one of the four off-loaded functions.  
 536 By having a single thread active during the entire  
 537 program execution, we avoid excessive overhead  
 538 from repeated spawning and joining of threads. The  
 539 SPE threads execute a busy wait for a PPE signal.

540 Data consistency is maintained at the granularity  
 541 of off-loaded functions. Each function is individually  
 542 responsible for collecting local updates and propa-  
 543 gating these updates to global shared memory.

#### 544 4.5. Optimizing Off-loaded Functions by Example 545 of `phvnewview()`

546 Since `newview()` is the most computationally  
 547 expensive function in the code, it becomes the first  
 548 candidate for optimization. We found that naively  
 549 off-loading `newview()` slows down the sequential  
 550 version of the code by a factor of 2.8. Therefore,  
 551 exhaustive optimization of the function on the SPE  
 552 was necessary. Table 2 summarizes the execution  
 553 times of RAXML before and after `newview()` is off-  
 554 loaded. The first column shows the number of  
 555 workers (MPI processes) used in the experiment  
 556 and the amount of work done (number of bootstraps).

557 We profiled the function `newview()` using the  
 558 decremter register of the SPE. We identified four  
 559 code segments that consume almost the entire  
 560 execution time in `newview()`:

- 561 1. Math library functions, such as `exp()` and `log()`.  
 562 These functions are very expensive if their native  
 563 implementations are executed on the SPE. In  
 564 `newview()`, the `exp()` function is used to compute  
 565 the transition probabilities of the nucleotide  
 566 substitution matrix for the branches from the root  
 567 of a subtree to its descendants. The `log()` function  
 568 is used to scale the branch lengths for numerical  
 569 reasons [23].
- 570 2. An `if (...)` statement in the code which  
 571 determines the value of a conjunction of four  
 572 expressions is the second bottleneck. This condi-  
 573 tional statement is used to check if small  
 574 likelihood vector entries need to be scaled to  
 575 avoid numerical underflow (similar operations are  
 576 used in every ML implementation).
- 577 3. Blocking DMA requests also have a significant  
 578 impact on execution time. Whenever data re-

Table 2. Execution time of RAXML (in seconds). t2.1

	Time (s)	
(a)		t2.3
1 worker, 1 bootstrap	28.3	t2.4
2 workers, 8 bootstraps	152.56	t2.5
2 workers, 16 bootstraps	309.53	t2.6
2 workers, 32 bootstraps	622.43	t2.7
(b)		t2.8
1 worker, 1 bootstrap	80.52	t2.9
2 workers, 8 bootstraps	348.36	t2.10
2 workers, 16 bootstraps	696.12	t2.11
2 workers, 32 bootstraps	1,375.09	t2.12

The input file is 42\_SC: (a) The application is executed on the PPE, (b) `newview()`, without optimizations, is off-loaded to one SPE. t2.13

579 required for the computation is not in local storage,  
 580 the program has to wait for the necessary data to  
 581 be fetched from global memory.

- 582 4. All double precision floating point arithmetic  
 583 used for the likelihood vector calculation is not  
 584 natively vectorized or optimized.

585 In the next few subsections we describe perfor-  
 586 mance optimizations for the off-loaded `newview()`  
 587 function. We applied analogous optimizations to the  
 588 remaining off-loaded functions. We do not discuss  
 589 the other functions in more detail due to space  
 590 limitations.

591 **4.5.1. Mathematical Functions.** The first step in  
 592 reducing the execution time of the off-loaded  
 593 function was to replace the expensive math functions  
 594 `exp()` and `log()` with the mathematical functions  
 595 provided by the Cell SDK 1.1. The `exp()` and `log()`  
 596 functions provided by the Cell SDK are implemen-  
 597 tations of numerical methods for exponent and  
 598 logarithm calculation. The `exp()` and `log()` functions  
 599 represent less than 1% of the total number of floating  
 600 point operations executed in the off-loaded function,  
 601 however they account for 56% of execution time in  
 602 `newview()`. Using the implementations of the  
 603 exponent and logarithm functions provided by the  
 604 Cell SDK improves total execution time by 37–41%.  
 605 Table 3 shows the execution times of the four test  
 606 runs of RAXML after `newview()` is off-loaded and  
 607 `exp()` and `log()` functions are replaced with

Exploring New Search Algorithms and Hardware for Phylogenetics

t3.1 *Table 3.* Execution time of RAxML with no off-loading (a), and with *phnewview()* off-loaded and optimized to use numerical functions from the SDK library.

t3.2		Time (s)
t3.3	(b)	
t3.4	1 worker, 1 bootstrap	47.94
t3.5	2 workers, 8 bootstraps	218.17
t3.6	2 workers, 16 bootstraps	433.94
t3.7	2 workers, 32 bootstraps	871.49

t3.8 The input file is 42\_SC.

608 optimized implementations. Notice that the off-  
609 loaded code is still 40–47% slower than the non-  
610 off-loaded code.

611 **4.5.2. Optimizing Conditional Statements.** Function  
612 *newview()* is always invoked at an inner node of the  
613 tree (*p*) which is at the root of a subtree. The main  
614 computational kernel of *newview()* has a **switch**  
615 statement which selects one out of four paths of  
616 execution. If one or both descendants *r* and *q* of *p*  
617 are tips (leaves) the computations of the main for-loop in  
618 *newview()* can be simplified. This optimization leads  
619 to significant performance improvements [25]. Thus,  
620 there are distinct implementations of the main  
621 computational part of *newview()* for the case that *r*  
622 and *q* are tips, *r* is a tip, *q* is a tip, or *r* and *q*  
623 are both inner nodes. Each path in the **switch** statement leads to  
624 a large loop which performs the likelihood vector  
625 calculations. Each iteration of the large loop executes  
626 a large **if()** statement which determines the value of a  
627 conjunction of four arithmetic expressions. This  
628 conditional statement checks if likelihood scaling is  
629 required to prevent numerical underflow. On an SPE,  
630 a mispredicted branch incurs a penalty of 20 cycles  
631 (IBM, Cbe\_tutorial\_v1.1, 2006). Using the decre-  
632 menter register to profile the off-loaded code, and  
633 after optimization of **exp()** and **log()**, we find that  
634 39% of the execution time of *newview()* is spent in  
635 checking the condition of the **if()** statement executed

in the likelihood vector calculation loop. The condi- 636  
tional statement is shown in Fig. 3. *minlikelihood* is 637  
a positive constant and all operands are double 638  
precision floating point values. 639

The **ABS()** function increases the number of 640  
condition checking (each **ABS()** function executes 641  
an additional comparison), therefore the number of 642  
conditions that need to be checked in this statement 643  
is actually eight. 644

On an SPE, two integer numbers can be compared 645  
significantly faster than two double precision floating 646  
point numbers. The advantage of integers is that they 647  
can be compared using the existing SPE intrinsics. 648  
Current SPE intrinsics support comparison of at most 649  
32-bit integer values. Sixty-four bit integers can also 650  
be compared relatively fast, by combining 32-bit 651  
integer intrinsics. The *spu-gcc* compiler automatical- 652  
ly optimizes conditional statements that operate on 653  
integer values, by replacing them with suitable SPE 654  
intrinsics. 655

According to the IEEE standard, all double precision 656  
floating point numbers are “lexicographically ordered.” 657  
In other words, if two double floating point numbers *a* 658  
and *b* are ordered ( $a < b$ ), then their bit patterns will 659  
be ordered in the same way when interpreted as Sign- 660  
Magnitude integers [16]  $(*(unsigned\ long\ long*)a$  661  
 $< *(unsigned\ long\ long*)b)$ . However, this rule can 662  
only be applied if both numbers are greater than 0. In 663  
the conditional statement that we are trying to optimize 664  
all parameters are greater than 0 (we know that 665  
**minlikelihood** is a constant greater than 0). Therefore, 666  
instead of comparing double precision floating point 667  
values, we can optimize the problematic **if()** statement 668  
by casting all our operands to **unsignedlonglong** before 669  
comparing them. To avoid the branch used in the 670  
**ABS()** function, we transform all our operands to 671  
positive numbers, using a bitwise **AND**. 672

This optimization reduces the time spent in the 673  
conditional statement to only 7% of the execution 674  
time of *newview()*. The total execution time of 675  
RAxML is reduced by 20%. Table 4 shows the new 676  
execution times of RAxML after optimizing both 677

```

if (ABS(x3->a) < minlikelihood && ABS(x3->g) < minlikelihood &&
    ABS(x3->c) < minlikelihood && ABS(x3->t) < minlikelihood)
{
    . . .
}

```

Figure 3. Conditional statement which takes 39% of *newview()* execution time, after optimization of numerical functions.

t4.1 *Table 4.* Execution time of RAxML after the floating-point conditional statement is transformed to an integer conditional statement.

t4.2		Time (s)
t4.3	(b)	
t4.4	1 worker, 1 bootstrap	38
t4.5	2 workers, 8 bootstraps	177.68
t4.6	2 workers, 16 bootstraps	354.85
t4.7	2 workers, 32 bootstraps	703.95

t4.8 The input file is 42\_SC.

678 conditional statements and numerical functions in  
 679 **newview()**.

680 **4.5.3. Double Buffering and Memory Manage-**

681 **ment.** The number of iterations in the main loop  
 682 of **newview()** depends on the input alignment length.  
 683 This loop operates on three arrays with equal length  
 684 (likelihood vector at current subtree root, and  
 685 likelihood vectors at left and right child). The loop  
 686 has no loop-carried dependencies and executes as  
 687 many iterations as the length of the three arrays.  
 688 Since the size of the local storage is only 256 KB,  
 689 the arrays cannot be stored permanently in local  
 690 storage. Instead, the arrays are strip-mined and  
 691 processed in blocks. The block size is selected such  
 692 that the computation on a block can overlap  
 693 completely with the memory latency for fetching  
 694 the next block. We use a 2 KB buffer for caching  
 695 each array block, which holds enough data to  
 696 execute 16 loop iterations. We also use double  
 697 buffering to mask memory latency, which amounts  
 698 to 8% of the execution time of the code before  
 699 optimization. Table 5 shows the improved execution  
 700 times of RAxML after optimization of numerical  
 701 operations, optimization of conditionals and memory  
 702 latency overlap, which improved total execution time  
 703 by 3–5%.

704 **4.5.4. Vectorization.** The core of the computation

705 in **newview()** is concentrated in two loops. The first  
 706 loop is executed at the beginning of the function and  
 707 computes the individual transition probability matrices  
 708 for each distinct rate category of the CAT model  
 709 [24]. The number of iterations is relatively small  
 710 (between 1 and 25) and each iteration executes 36  
 711 double precision floating point operations. The  
 712 second loop calculates the likelihood vector. The

length of this vector corresponds to the number of  
 distinct alignment patterns. Each iteration of this  
 loop executes 44 double precision floating point  
 operations for the CAT model.

The kernel of the first loop in **newview()** is shown  
 in Fig. 4a. In Fig. 4b we show the same code  
 vectorized for the SPE. The function **spu\_mul()**  
 multiplies two vectors (in this case the arguments are  
 vectors of doubles). Function **\_exp\_v()** is the vector  
 version of the exponential function mentioned in  
 Section 4.5.1. After vectorization, the number of  
 floating point operations executed in the body of the  
 first loop drops from 36 to 24. An additional  
 instruction is required to create a vector from a  
 scalar element. Due to pointer arithmetic on dynam-  
 ically allocated data structures, automatic vectoriza-  
 tion of this code would be particularly challenging  
 for a compiler.

The second loop is vectorized in a similar way.  
 Figure 5 shows the core of the second loop before  
 and after vectorization. The variables **x1** → **a**,  
**x1** → **c**, **x1** → **g**, **x1** → **t**, belong to the same C  
 structure (**likelihood\_vector**) and are stored in  
 contiguous memory locations. In Fig. 5a we see that  
 only three of these variables are multiplied with the  
 elements of array **left**. Therefore, vectorization  
 cannot be accomplished by simply interpreting the  
 members of the **likelihood\_vector** structure that  
 reside in consecutive memory locations as vectors.  
 We actually need to create vectors using special  
 intrinsics for vector creation, such as **spu\_splats()**.

Vectorization decreases the execution time of the  
 two major loops in **newview()** from 12.8 to 7.3 s.  
 Table 6 summarizes the execution times for RAxML,  
 after vectorization, optimization of numerical func-  
 tions, optimization of conditionals and  
 optimization for memory latency overlap in

t5.1 *Table 5.* Execution time of RAxML with double buffering  
 applied to overlap DMA transfers with computation, after  
 optimization of numerical functions and conditionals.

	Time (s)	t5.2
1 worker, 1 bootstrap	36.29	t5.3
2 workers, 8 bootstraps	169.50	t5.4
2 workers, 16 bootstraps	338.24	t5.5
2 workers, 32 bootstraps	688.04	t5.6
The input file is 42_SC.		t5.7

<p><b>a</b></p> <pre> for( ... ) {     ki = *rptr++;      d1c = exp (ki * lz10);     d1g = exp (ki * lz11);     d1t = exp (ki * lz12);      *left++ = d1c * *EV++;     *left++ = d1g * *EV++;     *left++ = d1t * *EV++;      *left++ = d1c * *EV++;     *left++ = d1g * *EV++;     *left++ = d1t * *EV++;     . . . }                 </pre>	<p><b>b</b></p> <pre> 1: vector double *left_v =     (vector double*)left; 2: vector double lz1011 =     (vector double)(lz10,lz11);     . . . for( ... ) { 3: ki_v = spu_splats(*rptr++);  4: d1cg = _exp_v ( spu_mul(ki_v,lz1011) );     d1tc = _exp_v ( spu_mul(ki_v,lz1210) );     d1gt = _exp_v ( spu_mul(ki_v,lz1112) );      left_v[0] = spu_mul(d1cg,EV_v[0]);     left_v[1] = spu_mul(d1tc,EV_v[1]);     left_v[2] = spu_mul(d1gt,EV_v[2]);     . . . }                 </pre>
---	---

Figure 4. The body of the first loop in `newview()`: **a** Non-vectorized code, **b** Vectorized code.

750 `newview()`. After vectorization, the code which off-  
 751 loads `newview()` on the SPE becomes up to 7%  
 752 faster than the code without off-loading.

753 **4.5.5. PPE–SPE Communication.** Although  
 754 `newview()` dominates the execution time of RAxML,  
 755 each instance of the function executes for only 82μs  
 756 on average. This means that the granularity of the  
 757 function is small and function invocation overhead when  
 758 the function is off-loaded on the SPE is by no means  
 759 negligible, since it requires PPE to SPE communication.  
 760 For our test dataset, the function is called 178,244 times.

761 We originally implemented PPE to SPE commu-  
 762 nication using mailboxes, which offer a programma-  
 763 ble, high-level interface for communication via  
 764 message queues. After further experimentation, we

found that that DMA transfers achieved lower  
 latency than using mailboxes and improved applica-  
 tion execution time by a further 2–11% in our test  
 cases. Table 7 shows the execution times of RAxML  
 when DMA transfers are used instead of mailboxes  
 and after all optimizations described so far. The  
 input dataset is 42\_SC.

It is interesting to note that direct memory-to-  
 memory communication is an optimization which  
 scales with parallelism on Cell, i.e. the impact on  
 performance improves as the code uses more SPEs.  
 As the number of workers and bootstraps executed  
 on the SPEs increases, the code becomes more  
 communication-intensive, due to the fine granularity  
 of the off-loaded functions. Fast communication by  
 DMA therefore becomes critical.

765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780

<p><b>a</b></p> <pre> for( ... ) {     ump_x1_0 = x1-&gt;a;     ump_x1_0 += x1-&gt;c * *left++;     ump_x1_0 += x1-&gt;g * *left++;     ump_x1_0 += x1-&gt;t * *left++;      ump_x1_1 = x1-&gt;a;     ump_x1_1 += x1-&gt;c * *left++;     ump_x1_1 += x1-&gt;g * *left++;     ump_x1_1 += x1-&gt;t * *left++;     . . . }                 </pre>	<p><b>b</b></p> <pre> for( ... ) {     a_v = spu_splats(x1-&gt;a);     c_v = spu_splats(x1-&gt;c);     g_v = spu_splats(x1-&gt;g);     t_v = spu_splats(x1-&gt;t);      l1 = (vector double)(left[0],left[3]);     l2 = (vector double)(left[1],left[4]);     l3 = (vector double)(left[2],left[5]);      ump_v1[0] = spu_madd(c_v,l1,a_v);     ump_v1[0] = spu_madd(g_v,l2,ump_v1[0]);     ump_v1[0] = spu_madd(t_v,l3,ump_v1[0]);     . . . }                 </pre>
--	--

Figure 5. Core of likelihood calculation loop in `newview()`: **a** Non-vectorized code, **b** Vectorized code.

t6.1 *Table 6.* Execution time of RAxML after vectorization, optimization of numerical functions, optimization of conditionals and optimization for memory overlap.

	Time (s)	
t6.3	1 worker, 1 bootstrap	30.13
t6.4	2 workers, 8 bootstraps	143.20
t6.5	2 workers, 16 bootstraps	287.31
t6.6	2 workers, 32 bootstraps	577.1
t6.7	The input file is 42_SC.	

781 *4.6. Off-loading Remaining Functions*

782 After off-loading and optimizing the **newview()**  
 783 function, we proceeded with off-loading the three re-  
 784 maining compute-intensive functions: **makenewz()**,  
 785 **evaluate()**, and **newviewPartial()**.

786 Off-loading **makenewz()** and **evaluate()** was  
 787 straightforward. We repeated the same off-loading  
 788 procedure and analogous optimizations as with  
 789 **newview()**. As described in Section 4.4, all off-  
 790 loaded functions are written to the same SPE code  
 791 module. In this way the off-loaded code remains in  
 792 the local storage during the entire execution of the  
 793 application. Consequently, we avoid the cost of  
 794 repeatedly loading different code modules into local  
 795 storage. Moreover, we reduce PPE–SPE communi-  
 796 cation, since each invocation of **newview()** by either  
 797 **makenewz()** or **evaluate()**, can be performed locally  
 798 on an SPE without involving the PPE.

799 Off-loading **newviewPartial()** was a more chal-  
 800 lenging task. This function is used to re-compute the  
 801 per-site log likelihood values  $ll_i$  at column  $i$  of the  
 802 alignment given a rate  $r_i$  during the per-site  
 803 evolutionary rate optimization process (see [24] for  
 804 details). Since optimizing the  $r_i$  implies changes of  
 805 the likelihood vector values at position  $i$  in the entire  
 806 tree, a complete tree traversal (as opposed to a partial  
 807 tree traversal induced by LSRs as described in  
 808 Section 3.2) must be carried out to obtain the  $ll_i$  for  
 809  $r_i$  at each alignment position  $i$ . Some changes in the  
 810 data structures were required to allow for future  
 811 parallel execution of **newviewPartial()**. Instead of  
 812 using the likelihood vector arrays allocated for  
 813 computations over the whole alignment length with  
 814 **newview()**, **makenewz()**, and **evaluate()** each  
 815 recursive invocation of **newviewPartial()** uses the  
 816 heap to create a private local likelihood vector of  
 817 length one. These changes enable the concurrent

computation of individual per-site log likelihood 818  
 values  $ll_i, ll_j$  at distinct alignment positions  $i \neq j$  819  
 with different rates  $r_i, r_j$ . 820

Due to the fine granularity of **newviewPartial()**, 821  
 which only operates on one single alignment position 822  
 and the high amount of calls per alignment position 823  
 (the  $r_i$  are optimized in a Brent-like procedure [24]) 824  
 we off-loaded the large for-loop that optimizes all  $r_i$  825  
 in **optimizeRateCategories()** to the SPE. This large 826  
 for-loop has already been parallelized in the current 827  
 OpenMP version of RAxML using the OpenMP 828  
**schedule(dynamic)** clause. Dynamic scheduling 829  
 avoids potential load imbalance because the optimi- 830  
 zation of each  $r_i$  requires a different number of 831  
 iterations until convergence. The parallelization of 832  
**newviewPartial()** with OpenMP yielded a perform- 833  
 ance improvement of 16% on the 404 sequence 834  
 dataset with 7,429 distinct patterns from Table 1 on a 835  
 four-way AMD Opteron. This large performance 836  
 improvement is due to the fact that **newviewPartial()** 837  
 consumes a significantly larger part of execution time 838  
 for long multi-gene alignments. 839

With all four functions off-loaded and optimized 840  
 on the SPE, the application was performing 34–36% 841  
 faster. Compared to the initial code which is entirely 842  
 executed on the PPE, the optimized code is 29% 843  
 faster. When more than one MPI processes is used 844  
 and more than one bootstrap is off-loaded to SPEs,  
 the gains from off-loading reach 47%. Table 8 845  
 summarizes the execution times for RAxML when all 846  
 four functions are off-loaded and optimized for the 847  
 42\_SC dataset. 848  
 849

4.7. *Scheduling Multilevel Parallelism* 850

The distribution of resources on the Cell often 851  
 introduces an imbalance between computation sup- 852

t7.1 *Table 7.* Execution time of RAxML after optimizing commu-  
 nication to use DMA transfers and all previously described  
 optimizations.

	Time (s)	
t7.2		
t7.3	1 worker, 1 bootstrap	29.9
t7.4	2 workers, 8 bootstraps	126.87
t7.5	2 workers, 16 bootstraps	253.67
t7.6	2 workers, 32 bootstraps	514.89
t7.7	The input file is 42_SC.	

t8.1 *Table 8.* Execution time of RAxML after off-loading and optimizing four functions: *phvnewview()*, *phvmakewz()*, *phvevaluate()* and *phvevaluatePartial()*.

t8.2		Time (s)
t8.3	1 worker, 1 bootstrap	20.4
t8.4	2 workers, 8 bootstraps	84.41
t8.5	2 workers, 16 bootstraps	165.3
t8.6	2 workers, 32 bootstraps	330.62
t8.7	The input file is 42_SC.	

853 ply and demand. The processor has eight SPEs,  
 854 however, at most two threads can off-load code from  
 855 the PPE to the SPEs at the same time. To overcome  
 856 this limitation, we explored an event-driven  
 857 programming model and the associated system  
 858 software support. In this model, we allow an  
 859 arbitrary number of MPI processes to share the PPE  
 860 and implement a context switching strategy which  
 861 opts for switching the context of a PPE thread upon  
 862 off-loading a function from the current context, in  
 863 anticipation of more opportunities for function off-  
 864 loading in other contexts. The intuition is that off-  
 865 loading a function from a PPE thread is typically  
 866 followed by idle time on the PPE thread, which can  
 867 be overlapped with computation originating in other  
 868 PPE threads.

869 We have extended our event-driven scheduling  
 870 model with an algorithm which tracks SPE utiliza-  
 871 tion and assigns SPEs to off-loaded functions based  
 872 on utilization. More specifically, each function  
 873 receives one or more SPEs, which can be used for  
 874 parallelization of loops enclosed in the function, with  
 875 mechanisms and policies similar to those of  
 876 OpenMP. Parallelization across SPEs assumes that  
 877 the work is distributed to the worker SPE threads  
 878 involved in the computation. All the communication  
 879 among the SPE threads is carried out via their local  
 880 storages. When they are created, all SPE threads  
 881 exchange the base addresses of their local storages  
 882 (this is done through the PPE). As a result, each SPE  
 883 thread owns a structure where it keeps the base  
 884 addresses of the local storages of all other threads  
 885 involved in the computation. A thread running on an  
 886 SPE communicates to other SPE threads by issuing  
 887 DMA requests to their private local storages.

888 The assignment of SPEs to functions depends on  
 889 the number of SPEs which are idling during a  
 890 predefined interval of recent execution. More details  
 891 on our event-driven model are provided in [5].

892 The new scheduling model reduces the execution  
 893 time of one bootstrap by 36%, compared to our  
 894 original static off-loading scheme, with all SPE-  
 895 specific optimizations integrated in the code. The  
 896 reason for the high speedup is the ability to distribute  
 897 loops inside off-loaded functions across SPEs. When  
 898 loop parallelism and task parallelism are exploited  
 899 simultaneously in off-loaded functions, the execution  
 900 time is reduced by up to 63%. Table 9 summarizes  
 901 the execution times of the optimized implementation  
 902 of RAxML with our event-driven programming and  
 903 scheduling model.

5. Comparison with the IBM Power5 904

905 It is useful to compare the performance of the Cell  
 906 against other multicore processors, since such a  
 907 comparison provides valuable insight both for appli-  
 908 cation developers working on adapting their software  
 909 to emerging computer architectures, and to computer  
 910 architects who are looking into improving their  
 911 hardware to address the needs of challenging  
 912 applications.

913 In this Section we compare the performance of  
 914 Cell against an IBM Power5, using our test runs of  
 915 RAxML. The IBM Power5 is a homogeneous dual-  
 916 core processor, where each core is itself a two-way  
 917 simultaneous multithreaded processor. Porting the  
 918 MPI version of RAxML to the Power5 is straight-  
 919 forward, since all that needs to be done is load four  
 920 MPI processes (workers) on the four execution  
 921 contexts of the processor. The Power5 used for this  
 922 experiment runs at 1.65 GHz, and has 32 KB of L1-  
 923 D and L1-I cache, 1.92 MB of L2 cache and 36 MB  
 924 of L3 cache.

925 Figure 6 provides the execution times on the two  
 926 processor types for up to 128 bootstraps, a scale

t9.1 *Table 9.* Execution time of RAxML with the event-driven programming and scheduling model (MGPS) is used.

	Time (s)	t9.2
1 bootstrap	14.1	t9.3
8 bootstraps	26.7	t9.4
16 bootstraps	53.63	t9.5
32 bootstraps	107.2	t9.6

t9.7 The input file is 42\_SC. The number of workers is variable and is selected at runtime by the scheduler.

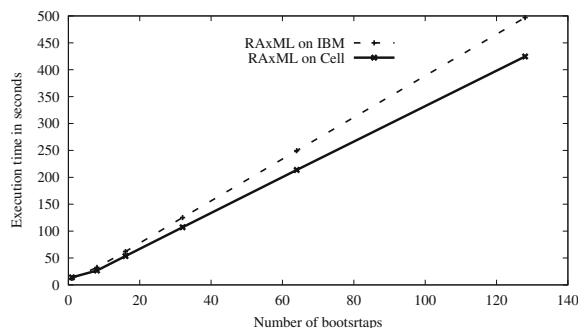


Figure 6. RAXML performance on a multi-core IBM Power5 and the IBM Cell. The number of bootstraps performed covers 1, 8, 16, 32, 64 and 128 runs.

927 which is more representative of real-world data sets  
 928 for RAXML. The Cell outperforms the IBM Power5  
 929 by 15% on average. Although the difference seems  
 930 small, a number of considerations should be taken  
 931 into account. The Power5 has a lower clock  
 932 frequency, but significantly more secondary and  
 933 tertiary cache space available to each core. Further-  
 934 more, double-precision floating point arithmetic is  
 935 unoptimized on the Cell's SPE pipelines, leading to a  
 936 markedly large reduction in processor throughput (up  
 937 to a factor of 10) compared to single-precision  
 938 floating point arithmetic. Furthermore, hardware  
 939 studies of the Cell indicate that the processor is  
 940 significantly more power-efficient than the Power5,  
 941 claiming nominal power consumption in the range of  
 942 27–43 W for the 3.2 GHz model used in this study  
 943 [30], as opposed to a reported 150 W for the Power5  
 944 [18]. Taking these observations into consideration,  
 945 we conclude that the Cell provides a leap forward in  
 946 performance compared to homogeneous, general-  
 947 purpose multicore processors. Our study is the first  
 948 to demonstrate this leap using complex, non-trivial  
 949 parallel code from the field of computational biology.

## 950 6. Conclusion and Future Work

951 We presented an improved heuristic search algorithm  
 952 for RAXML as well as a detailed step-by-step  
 953 description of porting RAXML to IBM Cell. The  
 954 incremental parallelization and optimization method-  
 955 ology presented in this paper can serve as a guideline  
 956 for parallelization of non-trivial applications on Cell.  
 957 The optimizations and system software for multilevel  
 958 parallelization introduced in this paper resolve bottle-  
 959 necks which are common to many applications and  
 960 even multicore architectures other than the Cell.

The porting strategy and methods developed for  
 exploiting fine-grain parallelism in RAXML on the  
 Cell are generally applicable to a broad range of  
 programs for ML-based (GARLI [31], IQPNNI [19],  
 PHYML [14]) and Bayesian (MrBayes [22]) phylo-  
 genetic inference. All these programs spend 90–95%  
 of their total execution time for the evaluation of the  
 likelihood function and face similar problems with  
 respect to memory transfer and function optimiza-  
 tion. The strategies and scheduling techniques for  
 coarse-grain parallelism are—with some modifica-  
 tions—also applicable to the MPI-versions of  
 GARLI and IQPNNI. In fact, there already exist a  
 hybrid MPI/OpenMP version of IQPNNI [20] and an  
 OpenMP parallelization of PHYML (Michael Ott,  
 personal communication).

Future work will focus on improved ways to  
 handle recursions on SPEs which will allow for  
 inference of large real-world datasets.

## Acknowledgments

We would like to thank Olaf Bininda-Emonds,  
 Nicolas Salamin, Josh Wilcox, Daniel Dalevi, Usman  
 Roshan, Chuck Robertson, and Markus Göker for  
 letting us use their often tediously hand-aligned  
 sequence data to assess RAXML performance.

## References

1. IBM, "Cell broadband engine programming tutorial version 1.0," Available at: <http://www-106.ibm.com/developerworks/eserver/library/es-archguide-v2.html>.
2. D. A. Bader, B. M. E. Moret, and L. Vawter, "Industrial Applications of High-performance Computing for Phylogeny Reconstruction," in *Proc. of SPIE ITCOM*, vol. 4528, 2001, pp. 159–168.
3. P. Bellens, J. M. Perez, R. M. Badia, and J. Labarta, "Cells: A Programming Model for the Cell be Architecture," in *Proc. of SC2006*, November 2006.
4. C. Benthin, I. Wald, M. Scherbaum, and H. Friedrich, "Ray Tracing on the CELL Processor," *Technical Report, inTrace Realtime Ray Tracing GmbH, No inTrace-2006-001*, 2006.
5. F. Blagojevic, D. S. Nikolopoulos, A. Stamatakis, and C. D. Antonopoulos, "Dynamic Multigrain Parallelization on the Cell Broadband Engine," in *Proc. of PPOPP 2007*, San Jose, CA, March 2007.
6. B. Chor and T. Tuller, "Maximum Likelihood of Evolutionary Trees: Hardness and Approximation," *Bioinformatics*, vol. 21, no. 1, 2005, pp. 97–106.

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

981

982

983

984

985

986

987 Q3

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

## Exploring New Search Algorithms and Hardware for Phylogenetics

- 1008 7. Z. Du, F. Lin, and U. Roshan, "Reconstruction of Large  
1009 Phylogenetic Trees: A Parallel Approach," *Computational*  
1010 *Biology and Chemistry*, vol. 29, no. 4, 2005, pp. 273–280.  
1011 8. A. E. Eichenberger et al., "Optimizing Compiler for a Cell  
1012 processor," *Parallel Architectures and Compilation Techni-*  
1013 *ques*, September 2005.  
1014 9. D. Pham et al., "The Design and Implementation of a First  
1015 Generation Cell Processor," *Proc. Int'l Solid-State Circuits*  
1016 *Conf. Tech. Digest*, IEEE Press, 2005, pp. 184–185.  
1017 10. K. Fatahalian et al., "Sequoia: Programming the Memory  
1018 Hierarchy," in *Proc. of SC2006*, November 2006.  
1019 11. R. E. Ley et al., "Unexpected Diversity and Complexity of the  
1020 Guerrero Negro Hypersaline Microbial Mat," *Appl. Environ.*  
1021 *Microbiol.*, vol. 72, no. 5, 2006, pp. 3685–3695, May.  
1022 12. J. Felsenstein, "Evolutionary Trees from DNA Sequences: A  
1023 Maximum Likelihood Approach," *J. Mol. Evol.*, vol. 17, 1981,  
1024 pp. 368–376.  
1025 13. G. W. Grimm, S. S. Renner, A. Stamatakis, and V. Hemleben,  
1026 "A Nuclear Ribosomal DNA Phylogeny of Acer Inferred with  
1027 Maximum Likelihood, Splits Graphs, and Motif Analyses of  
1028 606 Sequences," *Evolutionary Bioinformatics Online*, vol. 2,  
1029 2006, pp. 279–294.  
1030 14. S. Guindon and O. Gascuel, "A Simple, Fast, and Accurate  
1031 Algorithm to Estimate Large Phylogenies by Maximum  
1032 Likelihood," *Syst. Biol.*, vol. 52, no. 5, 2003, pp. 696–704.  
1033 15. N. Hjelte, Smoothed Particle Hydrodynamics on the Cell  
1034 Broadband Engine. *Masters Thesis*, June 2006.  
1035 16. W. Kahan, "Lecture Notes on the Status of IEEE Standard 754  
1036 for Binary Floating-point Arithmetic," in *IEEE*, 1997.  
1037 17. D. Kunzman, G. Zheng, E. Bohm, and L. V. Kalé, "Charm++,  
1038 Offload API, and the Cell Processor," in *Proc. of the*  
1039 *Workshop on Programming Models for Ubiquitous Parallel-*  
1040 *ism*, Seattle, WA, USA, September 2006.  
1041 18. Sun Microsystems, Sun UltraSPARC T1 Cool Threads  
1042 Technology, December 2005. [http://www.sun.com/aboutsun/](http://www.sun.com/aboutsun/media/presskits/networkcomputing05q4/T1Infographic.pdf)  
1043 [media/presskits/networkcomputing05q4/T1Infographic.pdf](http://www.sun.com/aboutsun/media/presskits/networkcomputing05q4/T1Infographic.pdf).  
1044 19. B. Q. Minh, L. S. Vinh, A. V. Haeseler, and H. A. Schmidt,  
1045 "pIQPNNI: Parallel Reconstruction of Large Maximum  
1046 Likelihood Phylogenies," *Bioinformatics*, vol. 21, no. 19,  
1047 2005, pp. 3794–3796.
20. B. Q. Minh, L. S. Vinh, H. A. Schmidt, and A. V. Haeseler, 1048  
"Large Maximum Likelihood Trees," in *Proc. of the NIC* 1049  
*Symposium 2006*, 2006, pp. 357–365. 1050  
21. C. E. Robertson, J. K. Harris, J. R. Spear, and N. R. Pace, 1051  
"Phylogenetic Diversity and Ecology of Environmental Arch- 1052  
aea," *Curr. Opin. Microbiol.*, vol. 8, 2005, pp. 638–642. 1053  
22. F. Ronquist and J. P. Huelsenbeck, "MrBayes 3: Bayesian 1054  
Phylogenetic Inference under Mixed Models," *Bioinformatics*, 1055  
vol. 19, no. 12, 2003, pp. 1572–1574. 1056  
23. A. Stamatakis, *Distributed and Parallel Algorithms and* 1057  
*Systems for Inference of Huge Phylogenetic Trees based on* 1058  
*the Maximum Likelihood Method*, PhD thesis, Technische 1059  
Universität München, Germany, October 2004. 1060  
24. A. Stamatakis, "Phylogenetic Models of Rate Heterogeneity: 1061  
A High Performance Computing Perspective," in *Proc. of* 1062  
*IPDPS2006*, HICOMB Workshop, Proceedings on CD, Rhod- 1063  
os, Greece, April 2006. 1064  
25. A. Stamatakis, "RAxML-VI-HPC: Maximum Likelihood-based 1065  
Phylogenetic Analyses with Thousands of Taxa and Mixed 1066  
Models," *Bioinformatics*, vol. 22, no. 21, 2006, pp. 2688–2690. 1067  
26. A. Stamatakis, T. Ludwig, and H. Meier, "Parallel Inference 1068  
of a 10,000-taxon Phylogeny with Maximum Likelihood," in 1069  
*Proc. of Euro-Par 2004*, September 2004, pp. 997–1004. 1070  
27. A. Stamatakis, T. Ludwig, and H. Meier, "RAxML-III: A Fast 1071  
Program for Maximum Likelihood-based Inference of Large 1072  
Phylogenetic Trees," *Bioinformatics*, vol. 21, no. 4, 2005, pp. 1073  
456–463. 1074  
28. A. Stamatakis, M. Ott, and T. Ludwig, "RAxML-OMP: An 1075  
Efficient Program for Phylogenetic Inference on SMPs," 1076  
*PaCT*, 2005, pp. 288–302. 1077  
29. C. Stewart, D. Hart, D. Berry, G. Olsen, E. Wernert, and W. 1078  
Fischer, "Parallel Implementation and Performance of FastD- 1079  
NAml—A Program for Maximum Likelihood Phylogenetic 1080  
Inference," in *Proc. of SC2001*, Denver, CO, November 2001. 1081  
30. D. Wang, "Cell Microprocessor III," *Real World Technolo-* 1082  
*gies*, July 2005. 1083  
31. D. Zwickl, *Genetic Algorithm Approaches for the Phyloge-* 1084  
*netic Analysis of Large Biological Sequence Datasets under* 1085  
*the Maximum Likelihood Criterion*. PhD thesis, University of 1086  
Texas at Austin, April 2006. 1087

## AUTHOR QUERIES

**AUTHOR PLEASE ANSWER ALL QUERIES.**

- Q1. Please provide author bios and photos.
- Q2. Please provide revised date.
- Q3. References were renumbered, please check.

UNCORRECTED PROOF