

Bringing High Performance Computing to Big Data Algorithms

H. Anzt, J. Dongarra, M. Gates, J. Kurzak, P. Luszczek, S. Tomov and I. Yamazaki

1 Abstract ■■■■

AQ1

2 **1 Introduction**

3 *1.1 High Performance Computing Meets Big Data*

4 High Performance Computing (HPC), meaning scientific and engineering comput-
5 ing, with emphasis on simulation, offers decades of experience in crunching numbers
6 at the highest speeds, using machines form the high end of the hardware spectrum.
7 Big Data, meaning data analytics, has been shifted more toward the lower end of
8 that spectrum, where the price/performance ratio is more favorable. Now that Big
9 Data problems enter the mainstream of computing, many solutions from HPC can
10 be applied to Big Data.

H. Anzt · J. Dongarra · M. Gates · J. Kurzak (✉) · P. Luszczek · S. Tomov · I. Yamazaki
Innovative Computing Laboratory, University of Tennessee,
1122 Volunteer Blvd, Knoxville, TN 37996, USA
e-mail: kurzak@icl.utk.edu

H. Anzt
e-mail: hanzt@icl.utk.edu

J. Dongarra
e-mail: dongarra@icl.utk.edu

M. Gates
e-mail: mgates3@icl.utk.edu

P. Luszczek
e-mail: luszczek@icl.utk.edu

S. Tomov
e-mail: tomov@icl.utk.edu

I. Yamazaki
e-mail: iyamazak@icl.utk.edu

© Springer International Publishing AG 2016
A.Y. Zomaya and S. Sakr (eds.), *Handbook of Big Data Technologies*,
DOI 10.1007/978-3-319-49340-4_23

1

This chapter opens with a discussion of the main differences between the hardware/software stacks of Big Data and HPC. Then two prominent HPC workloads are introduced, which happen to be in widespread use in the Big Data domain, the Alternating Least Squares (ALS) algorithm and the Singular Value Decomposition (SVD). Then the main techniques for maximizing the performance of the implementations are discussed. A comprehensive discussion of the implementation details of the ALS algorithm follows. Then a thorough presentation of the implementation details of the SVD algorithm is given. The chapter is concluded with the summary of the most important points.

High Performance Computing: In the 1980s, vector supercomputing dominated high-performance computing, as embodied in the eponymously named systems designed by the late Seymour Cray. The 1990s saw the rise of massively parallel processing (MPPs) and shared memory multiprocessors (SMPs) built by Thinking Machines, Silicon Graphics, and others. In turn, clusters of commodity (Intel/AMD x86) and purpose-built processors (such as IBM's BlueGene), dominated the previous decade.

Today, these clusters are augmented with computational accelerators in the form of coprocessors from Intel and graphical processing units (GPUs) from NVIDIA; they also include high-speed, low-latency interconnects (such as InfiniBand). Storage area networks (SANs) are used for persistent data storage, with local disks on each node used only for temporary files. This hardware ecosystem is optimized for performance first, rather than for minimal cost.

Atop the cluster hardware, Linux provides system services, augmented with parallel file systems (such as Lustre) and batch schedulers (such as PBS and SLURM) for parallel job management. MPI and OpenMP are used for internode and intranode parallelism, augmented with libraries and tools (such as CUDA and OpenCL) for coprocessor use. Numerical libraries (such as LAPACK and PETSc) and domain-specific libraries complete the software stack. Applications are typically developed in Fortran, C, or C++. Figure 1 (right) shows the mainstream HPC system stack.

Big Data: Just a few years ago, the very largest data storage systems contained only a few terabytes of secondary disk storage, backed by automated tape libraries. Today, commercial and research cloud-computing systems each contain many petabytes of secondary storage, and individual research laboratories routinely process terabytes of data produced by their own scientific instruments.

As with high-performance computing, a rich ecosystem of hardware and software has emerged for big data analytics. Unlike scientific computing clusters, data-analytics clusters are typically based on commodity Ethernet networks and local storage, with cost and capacity the primary optimization criteria. However, industry is now turning to FPGAs and improved network designs to optimize performance.

Atop this hardware, the Apache Hadoop system implements a MapReduce model for data analytics. Hadoop includes a distributed file system (HDFS) for managing large numbers of large files, distributed (with block replication) across the local storage of the cluster. HDFS and HBase, an open source implementation of Google's BigTable key-value store, are the big data analogs of Lustre for computational science, albeit optimized for different hardware and access patterns.

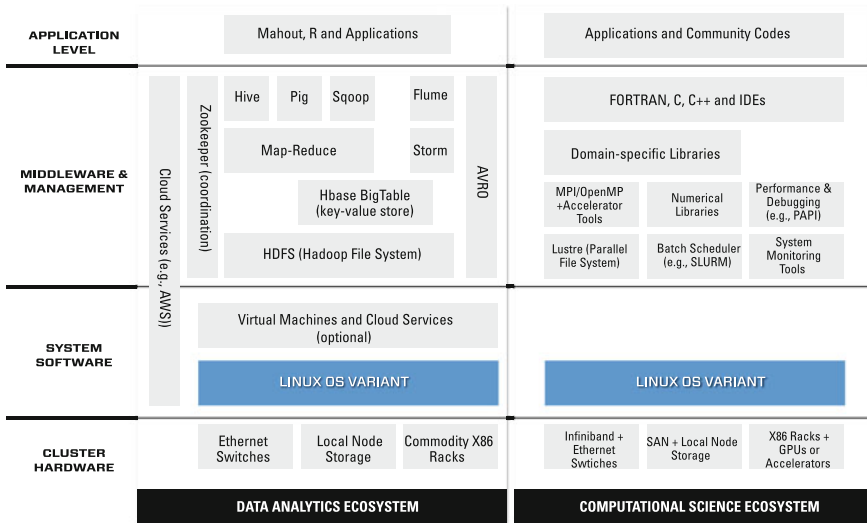


Fig. 1 The mainstream big data stack (left) versus the mainstream HPC stack (right)

Atop the Hadoop storage system, tools (such as Pig) provide a high-level programming model for the two-phase MapReduce model. Coupled with streaming data (Storm and Flume), graph (Giraph), and relational data (Sqoop) support, the Hadoop ecosystem is designed for data analysis. Moreover, tools (such as Mahout) enable classification, recommendation, and prediction via supervised and unsupervised learning. Unlike scientific computing, application development for data analytics often relies on Java and Web services tools (such as Ruby on Rails). Figure 1 (left) shows the mainstream Big Data system stack.

1.2 Application Areas

This chapter discusses HPC implementations of two mainstream Big Data algorithms. While the first one, Alternating Least Squares (ALS), has primarily commercial applications, the second one, Singular Value Decomposition (SVD), is uniformly applicable to a wide range of problems in science, engineering, and commerce.

The **Alternating Least Squares** algorithm provides a classic solution for building a recommender system for e-commerce, and was one of the more successful approaches to the Netflix Prize challenge. The importance of the algorithm is in its ability to deal with systems with implicit feedback, when the user's preference towards some products or content is known, while it is unknown for others. The weighted regularization process employed in the ALS algorithm allows for attaching higher weights to the known values and lower weight to the unknown values,

76 therefore effectively reconstructing the unknown values, as opposed to treating them
 77 as lack of interest. This approach leads to much more accurate recommendations
 78 than the simpler similarity-based algorithms.

79 One of the first open source implementations of the ALS algorithms was pro-
 80 duced in Java, as part of the Mahout machine learning package [48], which relied
 81 on the MapReduce paradigm provided by the Hadoop framework [34, 60]. As the
 82 MapReduce approach is being ousted by the Resilient Distributed Datasets (RDD)
 83 of the Spark framework [67], a faster implementation showed up in the Spark MLlib
 84 library [44]. Also, ALS was one of the first algorithms implemented in the GraphLab
 85 package [37], and also, for some time now, has been available in the Data Analytics
 86 Acceleration Library (DAAL) from Intel [26]. Finally, the first state of the art GPU
 87 implementation was produced by the authors of this chapter [16], and followed by
 88 similar developments from other groups [57].

89 The **Singular Value Decomposition** is ubiquitous in statistics and scientific com-
 90 puting and commonly applied to problems where the matrices are large and substan-
 91 tial computational power is required. Prime examples of application areas include
 92 astrophysics, genomics, climate data analysis, and information retrieval systems. In
 93 astrophysics, the SVD is used on massive datasets from astronomical surveys for
 94 spectral classification, e.g., to predict morphological types using galaxy spectra, and
 95 to select quasar candidates from sky surveys. In genomics, the SVD is routinely used
 96 to analyze genome-wide single-nucleotide polymorphism (SNP) data, for detecting
 97 population structure and potential outliers. In climate data analysis, Empirical orthog-
 98 onal function (EOF) and the SVD are the methods of choice for analyzing spacial
 99 and temporal variability of geophysical data. The SVD is also the primary tools for
 100 latent semantic indexing (LSI) in information retrieval systems, where it is used to
 101 find low-rank approximations to term-document matrices, enabling computation of
 102 query-document similarity scores in low-rank representation, as well as automated
 103 document categorization.

104 Randomized algorithms have been developed for the singular value decompo-
 105 sition [36, 42]. Great surveys of recent developments in randomization algorithms
 106 were published by Halko [21] and Mahoney [41]. In terms of software, singular value
 107 solvers are available in Skylark and Mahout. Skylark is an open-source software
 108 project launched by IBM Research with the objective to develop a set of random-
 109 ized machine learning algorithms that support distributed memory and are accessible
 110 through Python interfaces. Skylark uses a number of sketching transforms to imple-
 111 ment a few randomized linear algebra solvers, including a singular value solver based
 112 on the work by Halko et al. [21]. Mahout is a project of the Apache Software Founda-
 113 tion to produce free implementations of distributed or otherwise scalable machine
 114 learning algorithms focused primarily in the areas of collaborative filtering, cluster-
 115 ing, and classification [48]. In addition to a classic Lanczos SVD algorithm, Mahout
 116 also contains an implementation of a stochastic (randomized) SVD routine [40].

117 **1.3 Tricks of the Trade**

118 Two techniques discussed here and borrowed from the field of High Performance
119 Computing, are automated software tuning and randomization algorithms. The tech-
120 nique of automated software tuning mostly addressed the challenges of programming
121 modern computing devices, such as GPU accelerators, in a way that provides portable
122 performance, i.e., not only allows getting maximum performance from a particular
123 device, but also allows for porting to a new device by retuning rather than rewrit-
124 ing/redesigning the code. The technique of randomization allows dealing with one
125 of the most burning problems of processing Big Data, which is the lagging of IO
126 capabilities behind processing capabilities in modern hardware.

127 **Automated Software Tuning:** Although Moore’s Law has still been in effect in the
128 last few years, the multicore revolution initiated the trend, in processor design, of
129 going away from architectural features that do not directly contribute to processing
130 throughput. This means preference towards shallow pipelines with in-order execution
131 and cutting down on branch prediction and speculative execution. On top of that,
132 virtually all modern architectures require some form of vectorization to achieve top
133 performance, whether it being short-vector SIMD (Single Instruction Multiple Data)
134 extensions of CPU cores, or SIMT (Single Instruction Multiple Thread) pipelines
135 of GPU accelerators. With the landscape of future High Performance Computing
136 populated with complex, hybrid, vector architectures, automated software tuning may
137 provide a path towards portable performance without heroic programming efforts.

138 Automated software tuning was pioneered by projects like ATLAS and Spiral,
139 and is the objective of numerous academic projects, and is also practiced by hard-
140 ware vendors providing libraries like BLAS for their devices. The basic premise
141 is to explore a search space and find the best performers. The search space can be
142 defined by a set of tunable parameters, code transformations, implementation vari-
143 ants, hardware switches, etc. It can then be pruned by applying a set of constraints
144 that eliminate obvious underperformers. Finally, it can be searched to find the win-
145 ners. Exhaustive search, steepest descent methods, genetic algorithms are all valid
146 approaches.

147 **Randomization Algorithms:** The landscape of future High Performance Computing
148 presents an explosive growth in the volume of data, and a relatively dismal growth in
149 the capabilities of communication and IO systems. Under such conditions, it becomes
150 increasingly important to find algorithms that communicate less, and perform IO
151 operations even less. For an important set of problems in numerical computing, a class
152 algorithms emerges that seem to be an answer to these challenges—randomization
153 algorithms.

154 The new classes of random sampling and random projection algorithms offer
155 numerous advantages when dealing with large datasets coming from both scienti-
156 fic applications (astrophysics, genomics, climate modeling), as well as commercial
157 applications (social networks, information retrieval systems, financial transactions).
In many cases, randomized algorithms beat their classical counterparts in terms of

158 accuracy, speed, and robustness. They utilize modern computer architectures bet-
159 ter by exposing higher levels of parallelism than traditional numerical methods. At
160 the same time, they often produce more numerically robust solvers by introducing
161 implicit regularization.

162 2 GPU Acceleration of Alternating Least Squares

163 Web-based services such as movie databases and online retailers increasingly rely
164 on recommendation systems to suggest products to their customers. Collaborative
165 Filtering (CF) is a class of recommendation systems that recommends products based
166 on what other customers with similar interests have enjoyed [17]. It harvests infor-
167 mation collected from a large set of users, which can be either explicit feedback, such
168 as “likes” or product ratings; or implicit feedback, such as purchases, time spent, or
169 search patterns. This yields a large dataset to process, for instance, the Netflix Prize
170 dataset has over 100 million ratings [4].

171 Collaborative Filtering algorithms are based on observation data in a relation
172 matrix R , where each entry denotes how a user rated or interacted with an item. As
173 each user rates only a small subset of the items, most entries are unknown, i.e., the
174 matrix R is sparse. The goal is to determine the unknown values in R for how a user
175 would hypothetically rate every item. Thus it is an instance of the *matrix completion*
176 *problem* [9], to determine the unknown entries of a sparsely sampled matrix. In recent
177 years, latent feature models have assumed a small set of features—such as movie
178 genres—drive users’ interest. However, these latent features are determined by the
179 algorithm, without any explicit, a priori assigned meaning. This small set of features
180 implies the matrix R is (approximately) low-rank.

181 Besides providing an algorithm to complete R , an added benefit of the low-
182 rank model is that it determines R in a compact representation, $R = X^T Y$, taking
183 $O(fm + fn)$ space instead of $O(mn)$ space for m users, n items, and rank $f \ll m, n$.
184 For a site with millions of users and millions of products, this compact representation
185 makes storing and accessing the recommendations database tractable.

186 In addition to recommendation systems, the matrix completion problem occurs
187 in numerous other contexts. Examples include recovery of missing pixels of an
188 image [27], inferring 3D structure from motion of images [10], and determining
189 sensor positions from incomplete distance measurements [8].

190 Various methods exist for computing the matrix completion. Many CF systems
191 used neighborhood models [30]. For low-rank models, Candès and Recht [9] used
192 convex relaxation, and proved that R can be completed if sufficient entries are known.
193 Stochastic gradient descent [8, 50] and alternating least squares (ALS) [27, 70] are
194 popular methods. We will focus on the ALS method, which has adaptations for both
195 explicit [70] and implicit feedback [23].

We propose both multi-core CPU and GPU implementations that are able to exploit the computing power of state-of-the-art processors and accelerators. We compare performance with the open source implementations available in Mahout [1], GraphLab [12], and Spark MLlib [2, 44, 67], and report significant speedups for selected benchmark datasets.

2.1 Explicit Feedback

For explicit feedback, entry r_{ui} of R denotes how user u rated item i . Since users have not rated all items, the goal is to complete the missing entries of R . We assume that R is approximately low-rank, such that $R \approx X^T Y$, where X is $f \times m$ and Y is $f \times n$ for m users, n items, and rank or feature space size f . This latent feature space is small compared to the number of users and items, e.g., from 10 to 100, depending on the application. Column x_u of X represents user u , and column y_i of Y represents item i , such that their inner product yields the rating, $r_{ui} \approx x_u^T y_i$.

Determining X and Y is commonly expressed as an optimization problem, with a summation over known r_{ui} entries,

$$\min_{X,Y} \sum_{\substack{u,i \\ r_{ui} \text{ is known}}} (r_{ui} - x_u^T y_i)^2 + \lambda \left(\sum_u \|x_u\|^2 + \sum_i \|y_i\|^2 \right). \quad (1)$$

Here, λ is a regularization term to avoid overfitting. This can be solved with stochastic gradient descent or alternating least squares.

To solve using ALS, we observe that if X or Y is fixed, the cost function (1) becomes a linear least squares problem. ALS iterates two steps: fixing Y and solving for X , then fixing X and solving for Y . In the first step, fixing Y and finding where the gradient is zero yields

$$(Y D^u Y^T + \lambda I) x_u = Y r_u \quad \text{for } u = 1, \dots, m$$

to solve for each user-factor x_u . Each of the m user-factors can be solved independently, providing a large amount of parallelism. Here, r_u is row u of the R matrix, and D^u is a binary diagonal matrix that selects columns of Y corresponding to known r_{ui} values. Similarly, in the second step, fixing X yields

$$(X D^i X^T + \lambda I) y_i = X r_i \quad \text{for } i = 1, \dots, n$$

to solve for each item-factor y_i , where r_i is column i of the R matrix, and D^i selects columns of X for known r_{ui} values. Experiments have shown that the user- and item-factors typically converge after a few iterations of these two steps [70].

2.2 Implicit Feedback

For implicit feedback, Hu et al. [23] note that a large r_{ui} value does not necessarily indicate a higher preference, but instead gives a higher confidence. For instance, a user may enjoy watching a moderately good TV show every week, yielding a large r_{ui} value, but watch a beloved movie just once or twice, yielding a small r_{ui} value, despite its stronger preference. Therefore, they propose a preference matrix P with binary values,

$$p_{ui} = \begin{cases} 1 & \text{if } r_{ui} > 0, \\ 0 & \text{if } r_{ui} = 0, \end{cases}$$

to indicate whether user u has a preference for item i . Larger r_{ui} values indicate greater confidence in this preference, so a matrix C with entries $c_{ui} = 1 + \alpha r_{ui}$ is introduced that measures the confidence of the preference p_{ui} . Here some minimal confidence is given even to zero entries, while α weights known values more. Hu et al. found $\alpha = 40$ to work well. For implicit feedback, instead of completing the relation matrix R , the goal is to complete the preference matrix as $P \approx X^T Y$. Again, X and Y can be computed by minimizing a cost function,

$$\min_{X,Y} \sum_{u,i} c_{ui} (p_{ui} - x_u^T y_i)^2 + \lambda \left(\sum_u \|x_u\|^2 + \sum_i \|y_i\|^2 \right). \quad (2)$$

The major difference compared to explicit feedback is that the sum is over all u and i , not just those with nonzero r_{ui} values, since some minimal confidence is given even to zero entries. This means there are mn terms, making stochastic gradient descent prohibitively expensive for implicit feedback, whereas for explicit feedback only the nonzero r_{ui} values have terms in (1). Therefore, we apply the alternating least squares algorithm, similar to the explicit feedback case above, yielding

$$\begin{aligned} (YC^u Y^T + \lambda I) x_u &= YC^u p_u && \text{for } u = 1, \dots, m; \\ (XC^i X^T + \lambda I) y_i &= XC^i p_i && \text{for } i = 1, \dots, n; \end{aligned}$$

to solve for each x_u and for each y_i , where C^u is a diagonal matrix of row u of the confidence matrix C , C^i is a diagonal matrix of column i of C , p_u is row u of the preference matrix P , and p_i is column i of P . Pseudocode is given in Algorithm 1.

Algorithm 1 Pseudocode of alternating least square algorithm iterating user-factors and item-factors.

```

function ALS( input:  $\alpha, \lambda, R$ ; output:  $X, Y$  )
    set  $Y$  to random initial guess
    while not converged
        // update user-factors  $X$ 
        for  $u = 1, \dots, m$ 
            solve  $(YC^u Y^T + \lambda I) x_u = YC^u p_u$  for  $x_u$ 
        end
        // update item-factors  $Y$ 
        for  $i = 1, \dots, n$ 
            solve  $(XC^i X^T + \lambda I) y_i = XC^i p_i$  for  $y_i$ 
        end
    end
end function
    
```

259 The two steps, updating the user-factors and the item-factors, are identical except
 260 for swapping the input and output matrices. Therefore, we will subsequently focus
 261 on updating the user-factors, and the item-factors will follow similarly. The explicit
 262 and implicit feedback ALS algorithms are also very similar; we will concentrate on
 263 implicit feedback.
 264

For computational efficiency, the product can be factored as

$$YC^u Y^T = YY^T + \alpha Y R^u Y^T,$$

265
 266 where R^u is a diagonal matrix of row u of R , as shown schematically in Fig. 2.
 267 Since YY^T is the same for all users, it can be computed once per iteration [23],
 268 which is done efficiently using the syr_k (symmetric rank- k update) BLAS routine.
 269 (Explicit feedback lacks the YY^T term.) The remaining term, $\alpha Y R^u Y^T$, involves

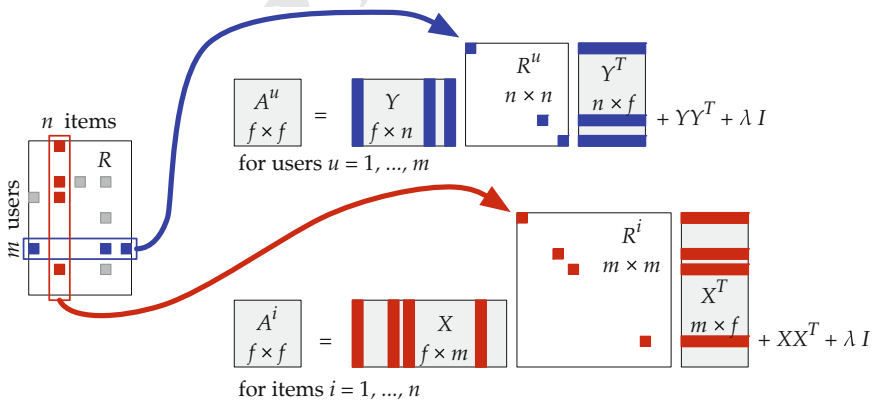
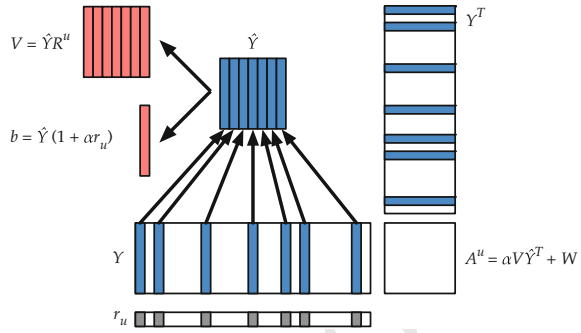


Fig. 2 Diagram of computation of user-factors and item-factors. R is general sparse, R^u and R^i are sparse diagonal, X, Y, A^u, A^i are dense

Fig. 3 Schematic of $A^u = Y R^u Y^T$ and $b = Y C^u p_u$. Shaded boxes in row r_u represent nonzeros; only corresponding shaded columns of Y and rows of Y^T contribute to A^u and b



270 a dense matrix Y and the sparse diagonal matrix R^u , which will require a custom
 271 kernel. Under mild assumptions, $Y C^u Y^T + \lambda I$ is symmetric positive definite (SPD),
 272 allowing us to solve it with the Cholesky factorization.

273 2.3 CPU Implementation

274 In the product $Y R^u Y^T$, the sparse diagonal matrix R^u selects and scales a few columns
 275 of Y , as shown in Fig. 3. Columns of Y corresponding to zeros in R^u are ignored.
 276 As k , the number of nonzeros in R^u , is typically much less than n , the number of
 277 columns of Y , the kernel should take advantage of this sparsity, reducing the cost
 278 from a rank- n update to a rank- k update, with $k \ll n$.

279 For instance, with the Netflix dataset and $f = 64$, the problem is to generate
 280 and solve $m = 480190$ systems, each formed by a 64×64 rank- k update, with
 281 the average $k = 209$ (see Fig. 5). There is not enough parallelism in computing a
 282 single system for an efficient multi-core implementation. Instead, we do a batched
 283 implementation that generates and solves the m systems in parallel. For this, we use
 284 OpenMP to parallelize the loops in Algorithm 2.

285 High efficiency can be attained by relying on optimized Level 3 BLAS routines,
 286 which operate on matrices instead of individual vectors, enabling data reuse and
 287 optimizations for cache efficiency, improving performance to be compute-bound
 288 instead of memory-bound. To use Level 3 BLAS, we copy the relevant columns of Y
 289 to workspaces \hat{Y} and V , with the R^u column scaling included in V , as shown in Fig. 3,
 290 then use a gemm (general matrix-matrix multiply) BLAS call. Since A is symmetric,
 291 work could be reduced by using an extended BLAS routine such as gemmt in Intel
 292 MKL [25] or syrks in NVIDIA cuBLAS [46] instead of gemm.

293 Updating the item-factors is exactly the same, except it uses columns of R instead
 294 of rows of R . For updating the user-factors, we store R in CSR (compressed sparse
 295 row) format, which gives efficient, contiguous access to each row of R , but slow
 296 access to columns of R . For efficiency in updating the item-factors, we also store R

Algorithm 2 Multi-core CPU ALS algorithm.

```

function ALS_CPU( input:  $\alpha, \lambda, R$ ; output:  $X, Y$  )
  set  $Y$  to random initial guess
  while not converged
    // update user-factors  $X$ 
     $W = YY^T + \lambda I$  using syrkl BLAS
    parallel for  $u = 1, \dots, m$ 
      copy columns of  $Y$  corresponding to nonzeros in  $r_u$  to  $\hat{Y}$ 
      copy and scale columns of  $\hat{Y}$  as  $V = \hat{Y}R^u$ 
      accumulate scaled columns of  $\hat{Y}$  as  $b_u = \hat{Y}(1 + \alpha r_u)$ 
       $A^u = \alpha V \hat{Y} + W$  using gemm BLAS (single-threaded)
      solve  $A^u x_u = b_u$  using Cholesky (single-threaded)
    end
    // update item-factors  $Y$ 
     $W = XX^T + \lambda I$  using syrkl BLAS
    parallel for  $i = 1, \dots, n$ 
      copy columns of  $X$  corresponding to nonzeros in  $r_i$  to  $\hat{X}$ 
      copy and scale columns of  $\hat{X}$  as  $V = \hat{X}R^i$ 
      accumulate scaled columns of  $\hat{X}$  as  $b_i = \hat{X}(1 + \alpha r_i)$ 
       $A^i = \alpha V \hat{X} + W$  using gemm BLAS (single-threaded)
      solve  $A^i y_i = b_i$  using Cholesky (single-threaded)
    end
  end
end function

```

297 in CSC (compressed sparse column) format, which gives efficient, contiguous access
 298 to each column of R .

299 Because the number of nonzeros per row can vary significantly (see Fig. 5), there
 300 will be a load imbalance between different processors. This is easily solved by using
 301 the OpenMP dynamic scheduler, adding `schedule (dynamic, NB)`, with a block
 302 size NB. We set $NB = 200$, but performance is not sensitive to the exact value.

303 2.4 GPU Implementation

304 A brief summary of the GPU architecture will help to understand the GPU implemen-
 305 tation. A GPU kernel divides its computation into a grid of thread blocks, and each
 306 thread block into a grid of threads. Within each thread block, threads are not indepen-
 307 dent, but execute the same instructions on different data. Threads can synchronize
 308 and communicate via shared memory, which is a kind of fast, user-controlled cache.
 309 Each thread's local variables are stored in a large register file. Different thread blocks
 310 execute asynchronously, without an easy way to synchronize or communicate. An
 311 NVIDIA Kepler GPU contains up to 15 multiprocessors, each with 192 cores.

312 Due to this GPU architecture, the GPU implementation shown in Algorithm 3
 313 is structured differently than the CPU implementation in Algorithm 2. Each thread
 314 block computes one tile of a matrix A^u and its right-hand side b_u . As with the CPU

315 implementation, a single system has insufficient parallelism to fully occupy all the
 316 GPU's cores. Filling a GPU requires hundreds of thread blocks and tens of thousands
 317 of threads. Therefore, we use a batched implementation, where a single GPU kernel
 318 generates a batch of s matrices using the BATCHED_SPARSE_SYRK routine, then a
 319 batched Cholesky routine factors them, and finally batched triangular solvers solve
 320 the resulting systems. We use the batched Cholesky and triangular solves from the
 321 BEAST project [33]. We used a batch size of $s = 4096$ to balance parallelism with
 322 GPU memory requirements. However, performance is not sensitive to the exact batch
 323 size.

Algorithm 3 GPU implementation of ALS, using batched operations.

```

function ALS_GPU( input:  $\alpha, \lambda, R$ ; output:  $X, Y$  )
  // workspaces:  $A$  is  $f \times f \times s$ ,  $B$  is  $f \times s$ 
  set  $Y$  to random initial guess
  while not converged
    // update user-factors  $X$ 
     $W = YY^T + \lambda I$  using syrk from cuBLAS
    for  $k = 1, \dots, m$  by batch size  $s$ 
      BATCHED_SPARSE_SYRK computes  $A^u = \alpha Y R^u Y^T + W$  and  $b_u = Y C^u p_u$ 
      for  $u = k, \dots, k + s$ 
      BATCHED_CHOLESKY factors  $A^u$  for  $u = k, \dots, k + s$ 
      BATCHED_SOLVE solves  $A^u x_u = b_u$  for  $u = k, \dots, k + s$ 
    end
    // update item-factors  $Y$ 
     $W = X X^T + \lambda I$  using syrk from cuBLAS
    for  $i = 1, \dots, n$  by batch size  $s$ 
      BATCHED_SPARSE_SYRK computes  $A^i = \alpha X R^i X^T + W$  and  $b_i = X C^i p_i$ 
      for  $i = k, \dots, k + s$ 
      BATCHED_CHOLESKY factors  $A^i$  for  $i = k, \dots, k + s$ 
      BATCHED_SOLVE solves  $A^i y_i = b_i$  for  $i = k, \dots, k + s$ 
    end
  end
end function

```

324 The implementation of the BATCHED_SPARSE_SYRK GPU kernel is conceptually
 325 similar to the CPU kernel. Like the CPU kernel, it copies the relevant columns of Y
 326 to a workspace \hat{Y} , in this case stored in GPU shared memory. Instead of copying all
 327 the relevant columns at once, it copies just one block of kb columns at a time and
 328 multiplies these, storing the results in registers, then continues with the next block.
 329 Unlike the CPU version, here the copy and multiply are fused into one kernel. The
 330 multiply is based an optimized gemm GPU kernel [32], which sub-tiles the output
 331 matrix A^u , with each GPU thread computing one entry in each sub-tile (Fig. 4).

332 A few optimizations can be made. Since A^u is symmetric, only the tiles on or
 333 below the diagonal need to be computed; tiles above the diagonal are known by
 334 symmetry. Also, since matrix Y is read-only, it is beneficial to bind its memory
 335 to GPU *texture memory*, which has optimized caching for read-only data. Texture
 336 memory also simplifies the code by dealing with out-of-bounds memory accesses

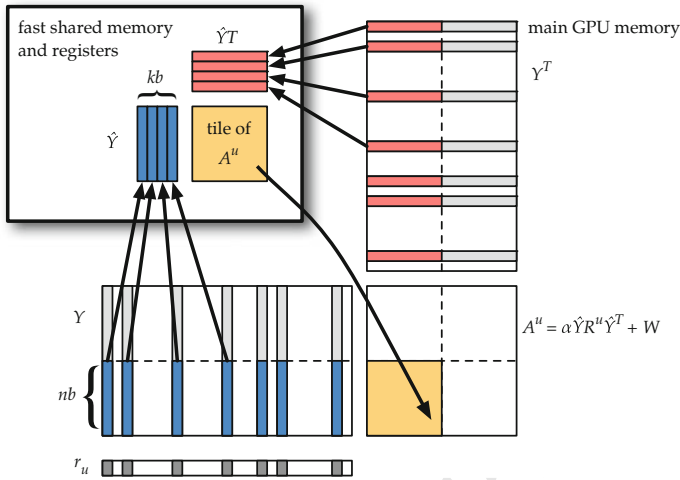


Fig. 4 Schematic of sparse-syrk GPU kernel. A^u is divided into $nb \times nb$ tiles. Block of kb relevant columns are loaded into shared memory and multiplied in registers at a time. At end, tile of A^u in registers is written back to main GPU memory. b_u is also computed (not shown)

337 in hardware—the software can pretend that Y is bigger than it actually is. This
 338 allows for fixed loop bounds and eliminates cleanup code, enabling more compiler
 339 optimizations.

340 2.5 Setup and Datasets

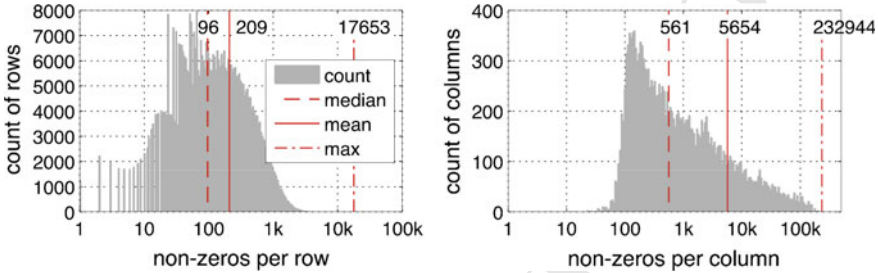
341 For performance comparison, we chose three ALS implementations from popular
 342 data analytics software packages: Mahout version 0.9 [1, 48], GraphLab version
 343 1.3 [12, 37, 38], and Spark MLlib version 1.5 [2]. All results used single precision and
 344 were obtained on a two-socket 2.6 GHz Intel Sandy Bridge E5-2670 with 8 cores per
 345 socket. CPU implementations were linked with Intel’s Math Kernel Library (MKL)
 346 version 11.1.2 [25]. Our GPU implementation ran on an NVIDIA Kepler K40c GPU
 347 with CUDA version 7.0 [47].

348 To compare performance, we target several recommendation datasets that are
 349 available online: Netflix Prize [4], Million Song [6], and Yahoo! Song [53]. For
 350 tuning parameters of the GPU implementation, we employ an autotuning sweep
 351 using the BEAST framework [24], with the EachMovie dataset [43, 53], a smaller
 352 dataset that permits executing a comprehensive set of kernel configurations in a
 353 moderate runtime. Table 1 summarizes properties of the datasets.

354 For the Netflix Prize dataset, we show histograms in Fig. 5 of the number of
 355 nonzeros per row (left) and per column (right). The minimum, median, mean, and
 356 maximum number of nonzeros per row and column are annotated in each graph. As

Table 1 Dataset properties

Dataset	# users	# items	# nonzeros
Netflix prize	480,190	17,771	100,480,508
Million song	1,019,318	384,546	48,373,586
Yahoo! song	130,558	136,736	49,770,695
EachMovie	1,623	61,265	2,811,717

**Fig. 5** Nonzero distribution of rows (*left*) and columns (*right*) of Netflix Prize dataset

357 previously noted, the wide range of nonzeros per row and column means different
 358 users and items incur widely different costs in computing $YC^u Y^T$ and $XC^i X^T$,
 359 potentially leading to load imbalance.

360 2.6 Auto Tuning

361 The sparse-syrk GPU kernel has four tunable parameters: tile size nb , block size
 362 kb , and thread block dimensions dx and dy . The kernel is generalized so that any
 363 value of nb can be used for any feature space size f . The optimal parameters are not
 364 obvious and not easy to derive by an analytical formula. Therefore the factorization
 365 calls for a real autotuning sweep. To achieve high performance, classic heuristic
 366 automatic software tuning methodology is applied, where a large number of kernels
 367 are generated and run, and the fastest ones identified.

368 The BEAST autotuning framework [39] enumerates and tests all possible kernel
 369 configurations. Various constraints are applied to limit the search space. Configu-
 370 rations violating correctness constraints—such as exceeding the maximum shared
 371 memory, or nb not divisible by the thread block dimensions—are eliminated. Sev-
 372 eral heuristic constraints are also applied, for instance, ensuring a compute-intensive
 373 kernel by requiring the ratio of multiply-add instructions to load instructions is at
 374 least 2. While kernels that violate these soft constraints will run correctly, they will
 375 not keep the GPU fully occupied, leading to lower performance.

376 After applying these constraints, BEAST generated 330 kernel configurations to
 377 test. The kernels were tested on the modest sized EachMovie dataset, timing the

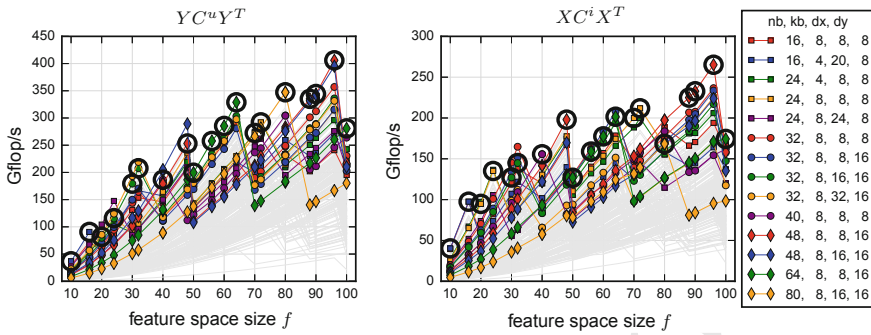


Fig. 6 Performance of all kernels (*gray lines*), highlighting ones that are best for some size. Circled kernel is chosen as best for each size (color figure online)

378 sparse-syrk for both the user-factor and the item-factor matrix generation. Due to
 379 differences in the size of Y and X and the sparsity of R^u and R^i , the performance
 380 was not identical between these two. We ran tests for sizes of f that are multiples of
 381 8 and multiples of 10, from 8 to 100.

382 The performance of all these kernels is plotted in gray in Fig. 6. Kernels that were
 383 best for some size are highlighted with colored markers. For each size f , the circled
 384 kernel was chosen as the best overall kernel.

385 Inspecting the data reveals that no one configuration was optimal across all feature
 386 space sizes. Taking the yellow diamond (80, 8, 16, 16) kernel as an example: for small
 387 f it is a poor performer, but the performance increases as f increases, until it is the
 388 best kernel for $f = 80$, where $f = nb$. For the next size, $f = 88$, its performance
 389 plummets to less than half the optimal performance. This occurs because it goes from
 390 one tile to four tiles covering each matrix A , wasting three large tiles to cover the
 391 extra 8 rows and columns. This saw tooth pattern is evident for all the configurations.

392 While often the best kernel for user-factors (left in Fig. 6) and item-factors (right)
 393 is the same, there are several instances where this is not true due to the difference in
 394 sparsity patterns. In these cases, the kernel with the best geometric mean performance
 395 is chosen as the best compromise between the two.

396 This analysis highlights the need for autotuning. The performance difference
 397 between the best and worst kernels is dramatic—between a factor of 6 and 72 times
 398 for a particular f . Also, the optimal kernel configuration depends heavily on the
 399 size f , and to a lesser extent on the actual dataset. While some kernel configurations
 400 make sense in retrospect, it was infeasible to predict optimal kernels in all cases.

401 2.7 Performance Evaluation

402 Execution time of a single ALS iteration (updating user-factors and item-factors
 403 once) for the three large benchmark databases—Netflix, Million Song, and Yahoo!

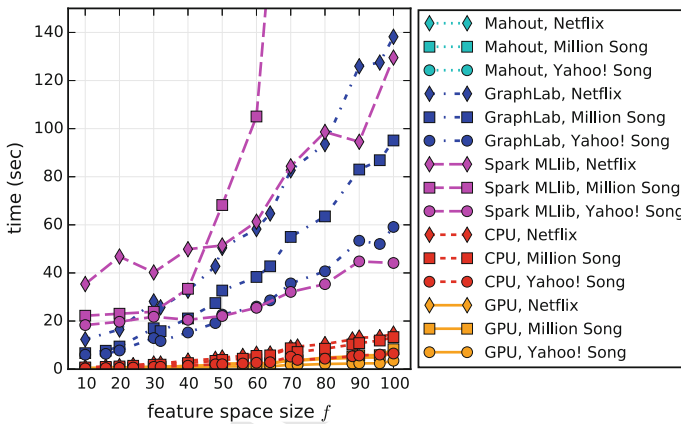
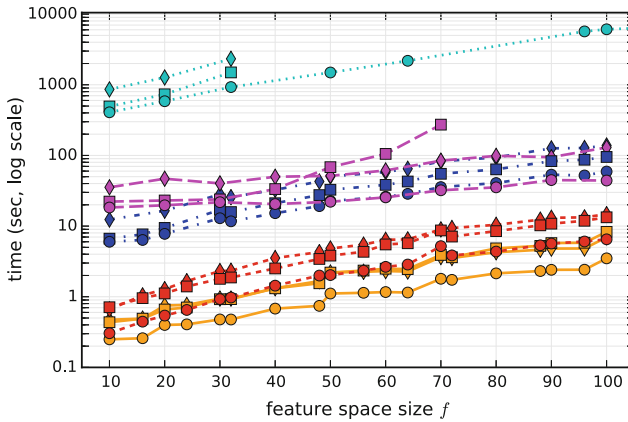


Fig. 7 Time in log scale (*top*) and linear scale (*bottom*) for single ALS iteration, using 16 CPU cores or GPU

404 Song—is presented in Fig. 7, in both log and linear scale. This covers a range of
 405 feature space sizes, all using 16 CPU cores or the GPU. A large performance differ-
 406 ence between implements is evident. Mahout is nearly two orders-of-magnitude
 407 slower than GraphLab and Spark. This is not surprising, as Mahout is written in
 408 Java while GraphLab is a newer implementation written in C++. Spark, while writ-
 409 ten in Scala/Java, links with native optimized BLAS to achieve good performance.
 410 For $f \geq 50$ with the Yahoo and Netflix datasets, Spark had performance compar-
 411 able to GraphLab. However, with the Million Song dataset, the Spark execution time
 412 increased markedly for $f \geq 50$, and it encountered an exception for $f \geq 80$. Our
 413 CPU implementation is 10 times faster than GraphLab and 19 times faster than Spark
 414 MLib, on average.

415 The speedup of our GPU implementation over Mahout, GraphLab, Spark, and our
 416 CPU implementation is given in Fig. 8. The GPU achieves an average speedup of 2.1

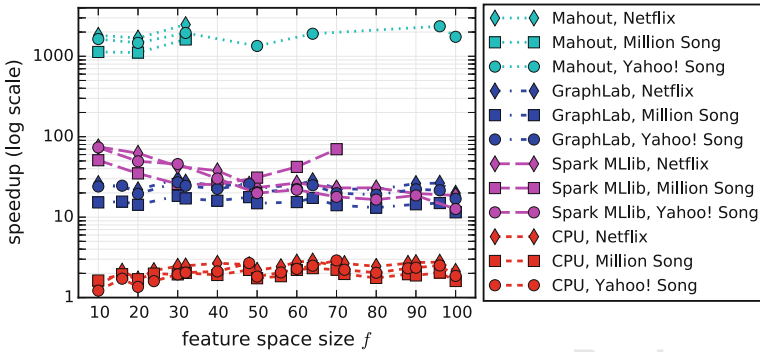


Fig. 8 Speedup in log scale of GPU implementation over Mahout, GraphLab, Spark, and CPU implementations using 16 cores

417 times over our CPU implementation. Compared to GraphLab, the GPU is on average
 418 20.9 times faster, and compared to Spark it is 35.3 times faster. Mahout performs
 419 poorly, taking 1684 times longer, on average, to compute a single ALS iteration.

420 While speedups are similar across datasets, our GPU implementation consistently
 421 gets the best speedups for the Netflix dataset and the least speedups for the Million
 422 Song dataset. This may be because the Million Song dataset has the smallest average
 423 nonzeros-per-row and nonzeros-per-column, with a mean of 47 nonzeros per row and
 424 126 per column, compared to 209 and 5654 for the Netflix dataset (Fig. 5). Having
 425 more nonzeros means a higher floating point operation count in the sparse-syrk
 426 routine to amortize memory reads.

427 We have presented both a multi-core CPU and a GPU implementation for the
 428 alternating least-squares algorithm to compute recommendations based on implicit
 429 feedback datasets. The central kernel involved is sparse_syrk, an algorithm-specific
 430 kernel achieving compute-bound performance for multiplying two dense matrices
 431 scaled by a sparse diagonal matrix. Our results demonstrate the advantage of fully
 432 exploiting the available parallelism by using a batched implementation, along with
 433 using optimized kernels, either from the vendor’s BLAS library or custom auto-tuned
 434 kernels. This yields good performance over several different datasets and a range of
 435 feature space sizes.

436 3 GPU Acceleration of Singular Value Decomposition

437 3.1 Introduction

438 A partial singular value decomposition (SVD) [18] of a sparse matrix is a power-
 439 ful tool for data analysis, where the data is represented as the sparse matrix. The
 440 ability of the SVD to filter out noise and extract the underlying features of the data

441 has been demonstrated in many applications, including Latent Semantic Indexing
 442 (LSI) [5, 13], recommendation systems [13, 55], population clustering [49], and
 443 subspace tracking [28]. The SVD is also used to compute the leverage scores – sta-
 444 tistical measurements for sampling the data in order to reduce the cost of the data
 445 analysis [21].

446 In recent years, the amount of data being generated from the observations,
 447 experiments, and simulations has been growing at unprecedented paces in many
 448 areas of studies, e.g., science, engineering, medicine, finance, social media, and
 449 e-commerce [11, 14]. The algorithmic challenges to analyze such “Big Data” are
 450 exacerbated by its massive volume and wide variety as well as its high veracity and
 451 velocity [35]. Though the SVD has the potential to address the variety and veracity of
 452 the modern data sets, the traditional approaches to computing the partial SVD access
 453 the data repeatedly, e.g., block Lanczos [19]. This is a significant drawback on a
 454 modern computer, where the data access has become significantly more expensive
 455 compared to arithmetic operations, both in terms of time and energy consumptions.
 456 The gap between the communication and computation costs is expected to further
 457 grow on future computers [15, 20], and this high cost of the communication is exacer-
 458 bated by the Big Data. This hardware trend is certainly true for the GPU.

459 3.2 Randomized Algorithms to Compute SVD

460 To address this hardware trend, a randomized algorithm [21] has been gaining atten-
 461 tion since compared to the traditional algorithms, it may require fewer data accesses
 462 to compute the SVD of the matrices arising from the modern applications (see Fig. 9
 463 for an illustration of the algorithm). To compare the performance of different algo-
 464 rithms for computing the truncated SVD, we implemented the framework, which
 465 encapsulates these algorithms on multicore CPUs with multiple GPUs [64]. This
 466 framework not only allows us to develop software whose performance can be tuned
 467 based on domain specific knowledge, but it also allows a user from one discipline
 468 to test an algorithm from another, or to combine the techniques from different algo-
 469 rithms (see Fig. 10 for the list of the algorithms). For example, we studied the per-
 470 formance of a block Lanczos, combining it with communication-avoiding [22, 62]
 471 and thick-restarting [3, 61]; two techniques developed by two different disciplines

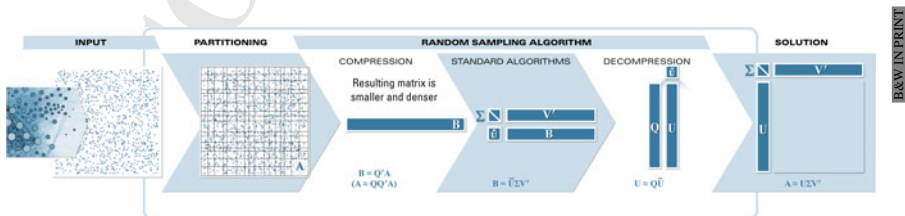
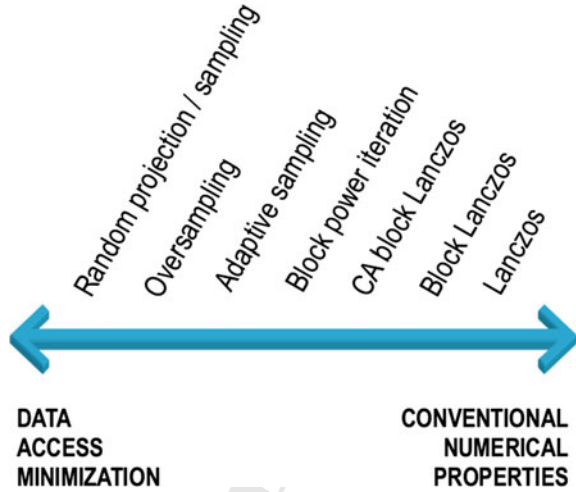


Fig. 9 Randomized algorithm to compute truncated SVD

Fig. 10 Algorithms to compute truncated SVD



472 – computer science and numerical linear algebra. These two techniques allow us to
 473 build the projection subspace with the minimum data access and accelerate the solu-
 474 tion convergence by retaining the useful information when restarting the iteration,
 475 respectively. Hence, compared to the randomized algorithm, Lanczos could build a
 476 projection subspace of the same dimension, which is richer in useful information
 477 with fewer communication phases, and potentially with about the same amount of
 478 data access. Unfortunately, this is possible only when the matrix can be partitioned
 479 well, while many of the matrices from the modern applications cannot be partitioned
 480 in such a way, leading to the significant overheads of the communication-avoiding
 481 technique in term of the computation and storage requirements, as well as the com-
 482 munication volume. Hence, there is a growing interest in a novel algorithm that can
 483 more efficiently compute the SVD of the massive data that are being generated from
 484 many modern applications, and the randomized algorithm is one of such algorithms
 485 with the potential.

486 3.3 Hybrid CPU/GPU Implementation

487 Figure 11 shows the pseudocode of a randomized algorithm to compute the SVD.
 488 Since the computational cost of the randomized algorithm is dominated by the cost
 489 of generating the projection basis vectors, \hat{P} and \hat{Q} , we accelerate this step using
 490 GPUs, while the SVD of the projected matrix B is redundantly computed by each MPI
 491 process on CPU. To generate the basis vectors, the two main computational kernels
 492 of the randomized algorithm are the sparse-matrix dense-matrix multiply (SpMM)
 493 and the orthogonalization. In this subsections, we describe our implementations of
 494 these two kernels on a hybrid CPU/GPU cluster.

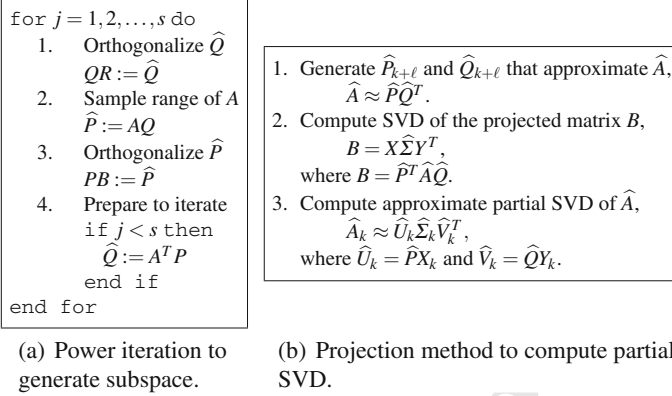
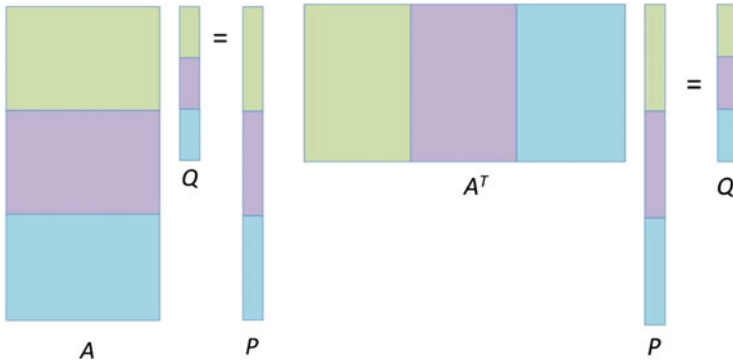


Fig. 11 Randomized algorithm to compute partial SVD based on power iteration

3.3.1 Sparse Matrix Matrix Multiply

To perform SpMM with the matrix A on a hybrid CPU/GPU cluster, we distribute A among the GPUs in a 1D block row format (e.g., using a graph or hypergraph partitioning algorithm). The basis vectors \widehat{P} and \widehat{Q} are then distributed in the same formats. Then, to perform SpMM, each GPU first exchanges the required non-local vector elements with its neighboring GPUs. This is done by first copying the required local elements from the GPU to the CPU, then performing the point-to-point communication among the neighbors using the non-blocking MPI (i.e., `MPI_Isend` and `MPI_Irecv`), and finally copying the non-local vector elements back to the GPU. Then, each GPU computes the local part of the next basis vectors using the CuSPARSE SpMM in the compressed sparse row (CSR) format. This was an efficient communication scheme in our previous studies to develop a linear solver [65], where the coefficient matrix A arising from a scientific or engineering simulation is often sparse and structured, e.g., with three-dimensional embedding. Unfortunately, sparse matrices originating from the modern data sets such as social networks and/or commercial applications have irregular sparsity structures, and have wide ranges of nonzero counts per row. In fact, they often exhibit power-law distributions of nonzeros as they result from scale-free graphs. As a result, this point-to-point communication with all the neighbors at once could be inefficient (in term of time and buffer storage). To alleviate the problem, our current implementation is based on a collective communication scheme. For example, using `MPI_Allgatherv`, each process sends its local vector elements, which are needed by at least one of its neighbors, to all the processes. Though this all-to-all approach requires the buffer to store the receiving messages from all the processes at once, it could obtain a significant speedup over the point-to-point communication, especially when the nonzeros of the matrix follows the power-law distribution.



(a) 1DBR with neighborhood-collective before local SpMM. (b) 1DBC with all-reduce after local SpMM.

Fig. 12 Illustration of matrix and vector distributions for SpMM with A and A^T . The submatrices distributed to the same GPU are colored in the same *color*. In Figure (a) or (b), the sparse matrices A and A^T are distributed either in 1D block row or block column (1DBR or 1DBC in short), respectively

521 Since many matrices of our interests are tall-skinny, to perform SpMM with A^T ,
 522 our current implementation keeps the input and output vectors, \hat{P} and \hat{Q} , in the
 523 1D block row distribution, but distribute A^T in the 1D block column (see Fig. 12b).
 524 Since the columns of A^T are the same as the rows of A on each GPU, we do not
 525 need to separately store A^T and A . In this implementation, each GPU first computes
 526 SpMM with its local parts of A^T and \hat{P} , and then copies the partial result to the
 527 CPU. Then, the MPI process computes the final result \hat{Q} by a global all-reduce, and
 528 copies its local part back to the GPU. Hence, this requires each MPI process to store
 529 the global vectors \hat{Q} . However, when A^T has the power-law distribution, performing
 530 SpMM with A^T in the 1D block row requires each GPU to store the much longer
 531 global vectors \hat{P} . Our performance results have demonstrated the advantage of this
 532 all-reduce communication. Furthermore, partitioning A^T in the 1D block column
 533 often led to a higher performance of SpMM on each GPU as the local submatrix
 534 becomes more square than tall-skinny.

535 3.3.2 Orthogonalization

536 For our experiments in this paper, we used the block classical Gram-Schmidt (CGS)
 537 [18] to orthogonalize a set of vectors against another set of vectors (block orthogonal-
 538 ization, or BOrth in short) and the Cholesky QR (CholQR) [56] to orthogonalize the
 539 set of vectors against each other. In our previous studies, these algorithms obtained
 540 great performance on multiple GPUs on a single compute node [63] or on a hybrid
 541 CPU/GPU cluster [65]. This is because these algorithms can orthogonalize the basis
 542 vectors with a low communication cost. For example, CholQR requires only one

543 global reduction between the GPUs, while most of the local computation is based on
544 BLAS-3 kernels on the GPU.

545 **3.4 Randomized Algorithms to Update SVD**

546 Though the randomized algorithms have the potential to efficiently compute the SVD
547 on the GPUs, there are several obstacles that need to be overcome. In particular, the
548 randomized algorithm may require only a small number of data accesses, but each
549 data access can be expensive due to the irregular sparsity pattern of the matrix and
550 the power-law distribution of its nonzeros. Though several techniques to avoid such
551 communication have been proposed [22], these techniques may not be effective for
552 computing the SVD of the modern data because they often require a significant
553 computational or communication overhead due to the particular sparsity structure of
554 the matrix [64].

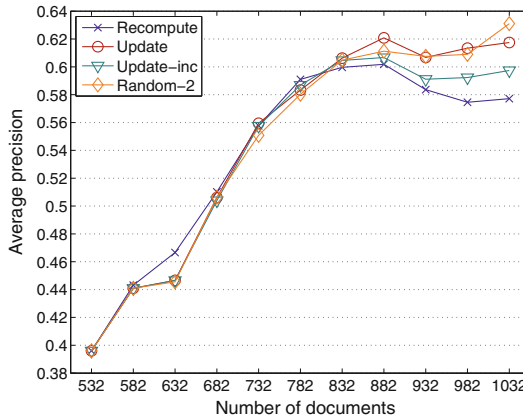
555 To address this challenge, we studied randomized algorithms to update (rather than
556 recompute) the partial SVD as the changes are made to the data set [66]. This is an
557 attractive approach because compared to recomputing it from scratch, the SVD may
558 be updated more efficiently, while in modern applications, the existing data are being
559 constantly updated and new data is being added. Moreover, in some applications,
560 recomputing the SVD may not be possible because the original data, for which the
561 SVD has been already computed, is no longer available. At the same time, in modern
562 applications, the size of the update is significant even though it is much smaller than
563 the massive data that has been already compressed. Therefore, an efficient updating
564 algorithm is needed to address the large volume and high velocity of the modern data
565 sets. Such applications with the rapidly changing data include the communication
566 and electric grids, transportation and financial systems, personalized services on the
567 internet, particle physics, astrophysics, and genome sequencing [11].

568 **3.4.1 Case Studies**

569 To study the potential of the randomized algorithm, we studied its performance for a
570 popular statistical analysis tool, the principal component analysis (PCA) [7]. In PCA,
571 a multidimensional dataset is projected onto a low-dimensional subspace given by the
572 partial SVD such that related items are close to each other in the projected subspace.
573 Here, we show the results from two particular applications of PCA, Latent Semantic
574 Indexing (LSI) and population clustering.

575 For information retrieval by text mining [54], a variant of PCA, Latent Semantic
576 Indexing (LSI) [13], has been shown to effectively address the ambiguity caused
577 by the synonymy or polysemy, which are difficult to address using a traditional
578 lexical-matching [31]. Figure 13a compares the average 11-point interpolated preci-
579 sions [29] after adding different numbers of documents from the MEDLINE matrix.
580 Our test matrices are the term-document matrices generated using the Text to Matrix

Fig. 13 Case studies with randomized algorithms for LSI ($k = 50$)



(a) Latent semantic indexing.

	Total number of documents, $n + d$							
Method	700	800	900	1000	1100	1200	1300	1400
Recompute	26.7	30.9	32.0	32.5	32.7	31.3	30.8	29.8
Update	26.7	29.8	30.1	30.7	31.5	30.7	30.4	29.7
Update-inc	26.7	29.8	30.1	30.6	30.9	30.1	29.8	29.5
Random-1	26.7	29.0	29.9	31.9	31.9	30.9	29.5	28.6
Random-2	26.7	29.6	29.6	30.0	31.0	30.1	30.0	29.7
Random-3	26.7	29.6	28.2	28.2	27.9	27.4	26.8	25.8

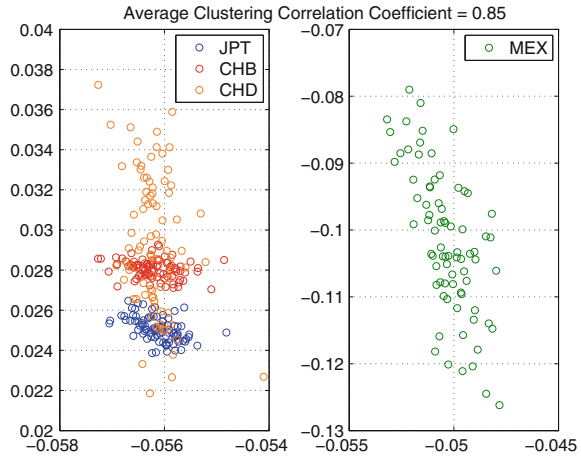
(b) Average 11-point interpolated precision for 6916-by-1400 CRANFIELD matrix with 225 queries, $n = 700$.

581 Generator (TMG)¹ and the TREC dataset,² and are preprocessed using the $1 \times n$. bpx
 582 weighing scheme [29]. These are the standard test matrices and were used in the pre-
 583 vious studies [59, 68]. For our studies, we first performed 20 power iterations of
 584 the randomized algorithm to compute the rank- k approximation of the matrix \hat{A}
 585 representing the first 700 documents. Then, the figure shows the average precision
 586 after new columns are added (e.g., under the column labeled “1000,” 300 docu-
 587 ments were added). To recompute the partial SVD of the matrix, we performed 20
 588 power iterations, while the randomized algorithm used the oversampling parameter
 589 set to be $\ell = k$ (i.e., $r = 2k$), and performed two iterations that access the matrix
 590 three times. Since the basis vectors \hat{P} and \hat{Q} approximate the ranges of \hat{A} and \hat{A}^T ,
 591 respectively, the randomized algorithm accesses the matrix at least twice. Then, they
 592 access the matrix one more time to compute the projected matrix B . We let the
 593 incremental update algorithm (Update-inc) add $k + \ell$ columns at a time such that it
 594 requires about the same amount of memory as the randomized algorithm. We see that
 595 with only three data passes, the randomized algorithm obtained similar precisions as
 596 those of the updating algorithm. In some cases, the updating and randomized algo-
 597 rithms obtained higher precisions than recomputing the SVD, while the precisions of

¹<http://scgroup20.ceid.upatras.gr:8000/tmg>.

²<http://ir.dcs.gla.ac.uk/resources>.

Fig. 14 Case studies with randomized algorithms for population clustering



(a) Population clustering.

	JPT+MEX + ASW + GIH + CEU			
Recompute	1.00	1.00	1.00	0.97
No update	1.00	0.81	0.84	0.67
Update-inc	1.00	1.00	0.89	0.70
Random-1	1.00	0.95	0.92	0.86

(b) Average correlation coefficients of population clustering based on the five dominant singular vectors, where 83 African ancestry in south west USA (ASW), 88 Gujarati Indian in Houston (GIH), and 165 European ancestry in Utah (CEU) were incrementally added to the 116,565 SNP matrix of 86 Japanese in Tokyo and 77 Mexican ancestry in Los Angeles, USA (JPT and MEX). Random-1 iterated twice with $\ell = k$.

598 the incremental update slightly deteriorated at the end. Such phenomena were also
599 reported in the previous studies [58, 68].

600 PCA has been also successfully used to extract the underlying genetic structure of
601 human populations [45, 51, 52]. To study the potential of the randomized algorithm,
602 we used it to update the SVD, when a new population is added to the population
603 dataset from the HapMap project.³ Figure 14 shows the correlation coefficient of
604 the resulting population cluster, which is computed using the k -mean algorithm
605 of MATLAB in the low-dimensional subspace given by the dominant left singular
606 vectors. We randomly filled in the missing data with either -1 , 0 , or 1 with the
607 probabilities based on the available information for the SNP. We let the randomized
608 algorithm iterate twice, and with only the three data passes, the randomized algorithm

³<http://hapmap.ncbi.nlm.nih.gov>.

609 improved the clustering results, potentially reducing the number of times the SVD
610 must be recomputed.

AQ3

611 3.4.2 Performance Studies

612 We now study the performance of the randomized algorithm on the Tsubame Com-
613 puter at the Tokyo Institute of Technology.⁴ Each of its compute nodes consists of
614 two six-core Intel Xeon CPUs and three NVIDIA Tesla K20Xm GPUs. We com-
615 piled our code using the GNU `gcc` version 4.3.4 compiler and the CUDA `nvcc`
616 version 6.0 compiler with the optimization flag `-O3`, and linked it with Intel's Math
617 Kernel Library (MKL) version `xe2013.1.046`.

618 Figure 15a compares the strong parallel scaling of the randomized algorithm with
619 that of the current state-of-the-art updating algorithm [68]. Clearly, the state-of-the-
620 art algorithm can spend significantly longer time in the orthogonalization, leading
621 to a great speedup obtained by the randomized algorithm (i.e., the speedups of up to
622 14.1). At the same time, the speedup decreased on a larger number of GPUs. This
623 is because the execution time of the randomized algorithm is dominated by SpMM,
624 whose strong parallel scaling suffered from the increasing inter-GPU communication
625 cost for this relatively small-scale matrix that was used for this study. On the other
626 hand, the updating algorithm was still spending a significant amount of its execution
627 time for the orthogonalization which was still compute intensive and scaled over
628 the small number of the GPUs. On a larger number of GPUs, compared to the
629 randomized algorithm, the updating algorithm is expected to suffer from the greater
630 communication latency.

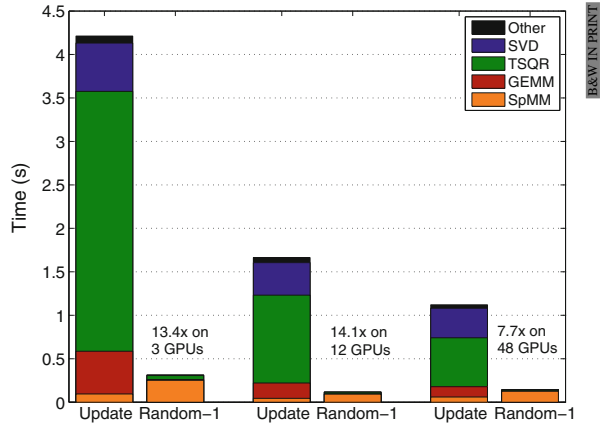
631 Figure 15b shows the weak parallel scaling results for the document-document
632 matrix used in a previous LSI study [69]. The matrix row contains 2,559,430 doc-
633 uments, and each column contains about 4, 176 nonzero entries. The weak parallel
634 scaling results, in particular, show the advantages of the randomized algorithm due
635 to its ability to compress the desired information into a small projection subspace
636 using a small number of data passes. For the updating algorithm, the accumulated
637 cost of the SVDs of the projected matrices also became significant.

638 4 Conclusions

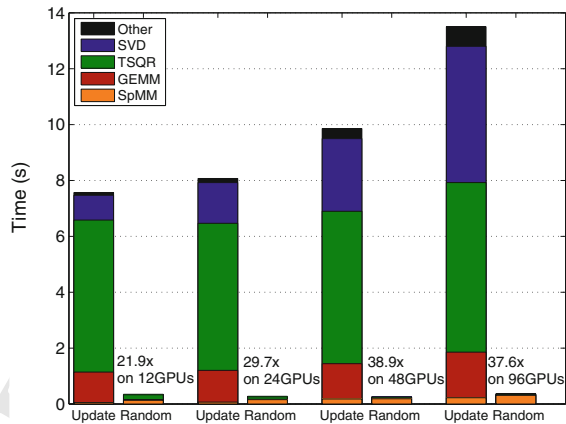
639 In this chapter, two mainstream Big Data algorithms were discussed: the Alternating
640 Least Squares algorithm for solving the matrix completion problem and the Singular
641 Value Decomposition algorithm for computing a low-rank approximation of a matrix,
642 both of which pose significant challenges when offloading to a GPU or a computing
643 cluster with multiple GPUs.

⁴<http://tsubame.gsfc.titech.ac.jp>.

Fig. 15 Performance studies with randomized algorithms



(a) Strong parallel scaling.



(b) Weak parallel scaling.

644 In the case of the ALS algorithm, the technique of automatic software tuning
 645 was used to achieve top performance, leading to an order of magnitude performance
 646 advantage over mainstream open source packages, GraphLab and Spark MLlib, and
 647 three orders of magnitude advantage over Mahout (Hadoop), when using a single
 648 GPU as opposed to a multicore CPU (16 cores).

649 In the case of the SVD algorithm, the technique of random projection was applied
 650 to implement the algorithm efficiently on a computing cluster with up to 48 GPUs,
 651 and also to implement an algorithm for updating a previously computed factorization
 652 upon arrival of new data. In this case, the algorithmic innovations also lead to an order
 653 of magnitude performance advantage.

654 Both case studies show the kind of impact that cutting-edge HPC techniques can
 655 have on the world of Big Data by enabling efficient use of accelerators, which leads
 656 to massive performance improvements.

657 References

- 658 1. Apache, Mahout version 0.9 (2015a). <https://mahout.apache.org/>
- 659 2. Apache, Spark version 1.5 (2015b). <http://spark.apache.org/>
- 660 3. J. Baglama, L. Reichel, Augmented implicitly restarted Lanczos bidiagonalization methods.
 661 SIAM J. Sci. Comput. **27**, 19–42 (2005)
- 662 4. J. Bennett, S. Lanning, The netflix prize, in *Proceedings of the KDD Cup Workshop 2007*
 663 (ACM, New York, 2007), pp 3–6. [http://www.cs.uic.edu/~liub/KDD-cup-2007/NetflixPrize-](http://www.cs.uic.edu/~liub/KDD-cup-2007/NetflixPrize-description.pdf)
 664 [description.pdf](http://www.cs.uic.edu/~liub/KDD-cup-2007/NetflixPrize-description.pdf)
- 665 5. M.W. Berry, Large scale sparse singular value computations. Int. J. Supercomput. Appl. **6**,
 666 13–49 (1992)
- 667 6. T. Bertin-Mahieux, D.P. Ellis, B. Whitman, P. Lamere, The million song dataset, in *Proceedings*
 668 *of the 12th International Conference on Music Information Retrieval (ISMIR)* (2011)
- 669 7. C. Bishop, *Pattern Recognition and Machine Learning* (Springer, New York, 2006)
- 670 8. P. Biswas, T.C. Lian, T.C. Wang, Y. Ye, Semidefinite programming based algorithms for sensor
 671 network localization. ACM Trans. Sensor Networks (TOSN) **2**(2), 188–220 (2006)
- 672 9. E.J. Candès, B. Recht, Exact matrix completion via convex optimization. Found. Comput.
 673 Math. **9**(6), 717–772 (2009)
- 674 10. P. Chen, D. Suter, Recovering the missing components in a large noisy low-rank matrix: appli-
 675 cation to SFM. IEEE Trans. Pattern Anal. Mach. Intell. **26**(8), 1051–1063 (2004)
- 676 11. Committee on the Analysis of Massive Data, Committee on Applied and Theoretical Statistics,
 677 Board on Mathematical Sciences and Their Applications, Division on Engineering and Physical
 678 Sciences, National Research Council (2013). *Frontiers in Massive Data Analysis*. The National
 679 Academies Press
- 680 12. Dato, GraphLab version 1.3 (2015). https://dato.com/products/create/open_source.html
- 681 13. S. Deerwester, S. Dumais, G. Furnas, T. Landauer, R. Harshman, Indexing by latent semantic
 682 analysis. J. Am. Soc. Inf. Sci. **41**, 391–407 (1990)
- 683 14. DOE Office of Science, Synergistic challenges in data-intensive science and exascale com-
 684 puting. DOE Advanced Scientific Computing Advisory Committee (ASCAC) (2013). Data
 685 Subcommittee Report
- 686 15. S.H. Fuller, L.I. Millett, *The Future of Computing Performance: Game Over Or Next Level?*
 687 (National Academy Press, Washington, DC, 2011)
- 688 16. M. Gates, H. Anzt, J. Kurzak, J. Dongarra, Accelerating collaborative filtering using concepts
 689 from high performance computing, in *2015 IEEE International Conference on Big Data (Big*
 690 *Data)* (IEEE, 2015), pp. 667–676
- 691 17. D. Goldberg, D. Nichols, B.M. Oki, D. Terry, Using collaborative filtering to weave an infor-
 692 mation tapestry. Commun. ACM **35**(12), 61–70 (1992)
- 693 18. G. Golub, C. van Loan, *Matrix Computations*, 4th edn. (The Johns Hopkins University Press,
 694 Baltimore, 2012)
- 695 19. G. Golub, F. Luk, M. Overton, A block Lanczos method for computing the singular values and
 696 corresponding singular vectors of a matrix. ACM Trans. Math. Softw. **7**, 149–169 (1981)
- 697 20. S. Graham, M. Snir, C. Patterson, *Getting Up to Speed: The Future of Supercomputing* (The
 698 National Academies Press, Washington, DC, 2004)
- 699 21. N. Halko, P. Martinsson, J. Tropp, Finding structure with randomness: probabilistic algorithms
 700 for constructing approximate matrix decompositions. SIAM Rev. **53**(2), 217–288 (2011)
- 701 22. M. Hoemmen, Communication-avoiding Krylov subspace methods. Ph.D. thesis, University
 702 of California, Berkeley (2010)

- 703 23. Y. Hu, Y. Koren, C. Volinsky, Collaborative filtering for implicit feedback datasets, in *IEEE*
704 *International Conference on Data Mining (ICDM)* (2008), pp. 263–272
- 705 24. Innovative Computing Lab, BEAST (2015). <http://icl.utk.edu/beast/>
- 706 25. Intel Corp, Developer Reference for Intel Math Kernel Library (2015). <https://software.intel.com/en-us/articles/mkl-reference-manual>
- 707
- 708 26. Intel Corp, Intel Data Analytics Acceleration Library 2016, Developer Guide (2016)
- 709 27. P. Jain, P. Netrapalli, S. Sanghavi, Low-rank matrix completion using alternating minimization,
710 in *Proceedings of the Forty-Fifth annual ACM Symposium on Theory of Computing* (ACM,
711 2013), pp 665–674
- 712 28. I. Karasalo, Estimating the covariance matrix by signal subspace averaging. *IEEE Trans.*
713 *Acoust. Speech Signal Process.* **34**(1), 8–12 (1986)
- 714 29. T. Kolda, D. O’Leary, A semidiscrete matrix decomposition for latent semantic indexing infor-
715 mation retrieval. *ACM Trans. Inf. Syst.* **16**(4), 322–346 (1998)
- 716 30. Y. Koren, Factorization meets the neighborhood: a multifaceted collaborative filtering model,
717 in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery*
718 *and Data Mining, KDD’08* (ACM, New York, 2008), pp. 426–434
- 719 31. R. Krovetz, W.B. Croft, Lexical ambiguity and information retrieval. *ACM Trans. Inf. Syst.*
720 **10**(2), 115–141 (1992)
- 721 32. J. Kurzak, S. Tomov, J. Dongarra, Autotuning gemm kernels for the Fermi GPU. *IEEE Trans.*
722 *Parallel Distrib. Syst.* **23**(11), 2045–2057 (2012)
- 723 33. J. Kurzak, H. Anzt, M. Gates, J. Dongarra, Implementation and tuning of batched Cholesky
724 factorization and solve for NVIDIA GPUs. *Trans. Parallel Distrib. Syst.* (2015). doi:[10.1109/](https://doi.org/10.1109/TPDS.2015.2481890)
725 [TPDS.2015.2481890](https://doi.org/10.1109/TPDS.2015.2481890)
- 726 34. C. Lam, *Hadoop in Action* (Manning Publications Co., Stamford, 2010)
- 727 35. D. Laney, 3D data management: controlling data volume, velocity, and variety. Application
728 Delivery Strategies by META Group Inc., File: 949 (2001)
- 729 36. E. Liberty, F. Woolfe, P.G. Martinsson, V. Rokhlin, M. Tygert, Randomized algorithms for the
730 low-rank approximation of matrices. *Proc. National Acad. Sci.* **104**(51), 20167–20172 (2007)
- 731 37. Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, J.M. Hellerstein, GraphLab: a new
732 framework for parallel machine learning. CoRR abs/1006.4990 (2010). [http://arxiv.org/abs/](http://arxiv.org/abs/1006.4990)
733 [1006.4990](http://arxiv.org/abs/1006.4990)
- 734 38. Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, J.M. Hellerstein, Distributed
735 GraphLab: a framework for machine learning and data mining in the cloud. *Proc. VLDB*
736 *Endow.* **5**(8), 716–727 (2012)
- 737 39. P. Luszczek, M. Gates, J. Kurzak, A. Danalis, J. Dongarra, Search space generation and pruning
738 system for autotuners, in *International Workshop on Automatic Performance Tuning (iWAPT*
739 *2016)* (2016, submitted)
- 740 40. D. Lyubimov, Command line interface, stochastic SVD. Technical report, The Apache Soft-
741 ware Foundation (2014). [https://mahout.apache.org/users/dim-reduction/ssvd.page/SSVD-](https://mahout.apache.org/users/dim-reduction/ssvd.page/SSVD-CLI.pdf)
742 [CLI.pdf](https://mahout.apache.org/users/dim-reduction/ssvd.page/SSVD-CLI.pdf)
- 743 41. M.W. Mahoney, Randomized algorithms for matrices and data. *Found. Trends® Mach. Learn.*
744 **3**(2), 123–224 (2011)
- 745 42. P.G. Martinsson, V. Rokhlin, M. Tygert, A randomized algorithm for the approximation of
746 matrices. Technical report, DTIC Document (2006)
- 747 43. P. McJones, Eachmovie collaborative filtering data set. DEC Systems Research Center 249
748 (1997)
- 749 44. X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai,
750 M. Amde, S. Owen et al., MLlib: Machine learning in Apache Spark (2015). arXiv preprint
751 [arXiv:150506807](https://arxiv.org/abs/150506807)
- 752 45. P. Menozzi, A. Piazza, L. C-Sforza, Synthetic maps of human gene frequencies in Europeans.
753 *Science* **201**, 786–792 (1978)
- 754 46. NVIDIA Corp, cuBLAS Library User Guide, v7.0 (2015a)
- 755 47. NVIDIA Corp, CUDA C Programming Guide, v7.0 (2015b)

- 756 48. S. Owen, R. Anil, T. Dunning, E. Friedman, *Mahout in Action* (Manning Publications Co.,
757 Greenwich, 2011)
- 758 49. P. Paschou, E. Ziv, E. Burchard, S. Choudhry, W. R-Cintron, M. Mahoney, P. Drineas, PCA-
759 correlated SNPs for structure identification in worldwide human populations. *PLoS Genet.* **3**,
760 1672–1686 (2007)
- 761 50. A. Paterek, Improving regularized singular value decomposition for collaborative filtering, in
762 *Proceedings of KDD Cup and Workshop* (2007), pp. 39–42
- 763 51. N. Patterson, A. Price, D. Reich, Population structure and eigenanalysis. *PLoS Genet.* **2**(12),
764 2074–2093 (2006)
- 765 52. A. Price, N. Patterson, R. Plenge, M. Weinblatt, N. Shadick, D. Reich, Principal components
766 analysis corrects for stratification in genome-wide association studies. *Nature Genet.* **38**(8),
767 904–909 (2006)
- 768 53. R.A. Rossi, N.K. Ahmed, The network data repository with interactive graph analytics and
769 visualization, in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*
770 (2015). <http://networkrepository.com>
- 771 54. G. Salton, M. McGill, *Introduction to Modern Information Retrieval* (McGraw-Hill, New York,
772 1983)
- 773 55. B. Sarwar, G. Karypis, J. Konstan, J. Riedl, Analysis of recommendation algorithms for e-
774 commerce, in *Proceedings of the 2nd ACM Conference on Electronic Commerce* (2000), pp
775 158–167
- 776 56. A. Stathopoulos, K. Wu, A block orthogonalization procedure with constant synchronization
777 requirements. *SIAM J. Sci. Comput.* **23**(6), 2165–2182 (2002)
- 778 57. W. Tan, L. Cao, L.L. Fong, Faster and cheaper: Parallelizing large-scale matrix factorization
779 on gpus. CoRR abs/1603.03820 (2016). <http://arxiv.org/abs/1603.03820>
- 780 58. J. Tougas, R. Spiteri, Updating the partial singular value decomposition in latent semantic
781 indexing. *Comput. Statist. Data Anal.* **52**, 174–183 (2007)
- 782 59. E. Vecharynski, Y. Saad, Fast updating algorithms for latent semantic indexing. *SIAM J. Matrix*
783 *Anal. Appl.* **35**(3), 1105–1131 (2014)
- 784 60. T. White, *Hadoop: The Definitive Guide* (O’Reilly Media, Inc., Sebastopol, 2012)
- 785 61. K. Wu, H. Simon, Thick-restart Lanczos method for large symmetric eigenvalue problems.
786 *SIAM J. Matrix Anal. Appl.* **22**(2), 602–616 (2000)
- 787 62. I. Yamazaki, K. Wu, A communication-avoiding thick-restart lanczos method on a distributed-
788 memory system, in *Proceedings of the 2011 International Conference on Parallel Processing,*
789 *Euro-Par’11* (Springer, Berlin, 2012), pp. 345–354
- 790 63. I. Yamazaki, H. Anzt, S. Tomov, M. Hoemmen, J. Dongarra Improving the performance of
791 CA-GMRES on multicores with multiple GPUs, in *Proceedings of the IEEE International*
792 *Parallel and Distributed Symposium (IPDPS)* (2014a), pp. 382–391
- 793 64. I. Yamazaki, T. Mary, J. Kurzak, S. Tomov, Access-averse framework for computing low-rank
794 matrix approximations, in *Proceedings of the International Workshop on High Performance*
795 *Big Graph Data Management, Analysis, and Minig* (2014b), pp. 70–77
- 796 65. I. Yamazaki, S. Rajamanickam, E. Boman, M. Hoemmen, M. Heroux, S. Tomov, Domain
797 decomposition preconditioners for communication-avoiding Krylov methods on a hybrid
798 CPU/GPU cluster, in *Proceedings of the International Conference for High Performance Computing,*
799 *Networking, Storage and Analysis (SC)* (2014c), pp. 933–944
- 800 66. I. Yamazaki, J. Kurzak, P. Luszczek, J. Dongarra, Randomized algorithms to update partial
801 singular value decomposition on a hybrid CPU/GPU cluster, in *Proceedings of the International*
802 *Conference for High Performance Computing, Networking, Storage and Analysis (SC)* (2015),
803 pp. 345–354
- 804 67. M. Zaharia, M. Chowdhury, M.J. Franklin, S. Shenker, I. Stoica, Spark: cluster computing
805 with working sets, in *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud*
806 *Computing*, vol. 10 (2010), p.10
- 807 68. H. Zha, H. Simon, On updating problems in latent semantic indexing. *SIAM J. Sci. Comput.*
808 **21**(2), 782–791 (1999)



- 809 69. H. Zha, O. Marques, H. Simon, Large-scale SVD and subspace-based methods for information
810 retrieval, in *Solving Irregularly Structured Problems in Parallel*, vol. 1457, Lecture Notes in
811 Computer Science, ed. by A. Ferreira, J. Rolim, H. Simon, S.-H. Teng (Springer, Heidelberg,
812 1998), pp. 29–42
- 813 70. Y. Zhou, D. Wilkinson, R. Schreiber, R. Pan, Large-scale parallel collaborative filtering for
814 the netflix prize in *Proceedings of the 4th International Conference on Algorithmic Aspects in*
815 *Information and Management, AAIM'08* (Springer, Berlin, 2008), pp. 337–348

Author Queries

Chapter 23

Query Refs.	Details Required	Author's response
AQ1	As abstract is mandatory for this chapter, please provide.	
AQ2	Please check and confirm if the inserted citation of Fig. 4 is correct. If not, please suggest an alternate citation. Please note that figures should be cited sequentially in the text.	
AQ3	Please note that mismatch has been found between author tex and pdf regarding the figure citation "Figure 14a". Hence, we have followed author tex. Kindly check and confirm.	

MARKED PROOF

Please correct and return this set

Please use the proof correction marks shown below for all alterations and corrections. If you wish to return your proof by fax you should ensure that all amendments are written clearly in dark ink and are made well within the page margins.

<i>Instruction to printer</i>	<i>Textual mark</i>	<i>Marginal mark</i>
Leave unchanged	... under matter to remain	Ⓟ
Insert in text the matter indicated in the margin	∧	New matter followed by ∧ or ∧ [Ⓢ]
Delete	/ through single character, rule or underline or ┌───┐ through all characters to be deleted	Ⓞ or Ⓞ [Ⓢ]
Substitute character or substitute part of one or more word(s)	/ through letter or ┌───┐ through characters	new character / or new characters /
Change to italics	— under matter to be changed	↙
Change to capitals	≡ under matter to be changed	≡
Change to small capitals	≡ under matter to be changed	≡
Change to bold type	~ under matter to be changed	~
Change to bold italic	≈ under matter to be changed	≈
Change to lower case	Encircle matter to be changed	≡
Change italic to upright type	(As above)	⊕
Change bold to non-bold type	(As above)	⊖
Insert 'superior' character	/ through character or ∧ where required	Υ or Υ under character e.g. Υ or Υ
Insert 'inferior' character	(As above)	∧ over character e.g. ∧
Insert full stop	(As above)	⊙
Insert comma	(As above)	,
Insert single quotation marks	(As above)	Ƴ or ƴ and/or ƶ or Ʒ
Insert double quotation marks	(As above)	ƶ or Ʒ and/or Ƶ or ƴ
Insert hyphen	(As above)	⊥
Start new paragraph	┌	┌
No new paragraph	┐	┐
Transpose	└┘	└┘
Close up	linking ○ characters	○
Insert or substitute space between characters or words	/ through character or ∧ where required	Υ
Reduce space between characters or words		↑