

Scalable Language Processing Algorithms for the Masses: A Case Study in Computing Word Co-occurrence Matrices with MapReduce

Jimmy Lin

The iSchool, University of Maryland
National Center for Biotechnology Information, National Library of Medicine
jimmylin@umd.edu

Abstract

This paper explores the challenge of scaling up language processing algorithms to increasingly large datasets. While cluster computing has been available in commercial environments for several years, academic researchers have fallen behind in their ability to work on large datasets. I discuss two barriers contributing to this problem: lack of a suitable programming model for managing concurrency and difficulty in obtaining access to hardware. Hadoop, an open-source implementation of Google's MapReduce framework, provides a compelling solution to both issues. Its simple programming model hides system-level details from the developer, and its ability to run on commodity hardware puts cluster computing within the reach of many academic research groups. This paper illustrates these points with a case study in building word co-occurrence matrices from large corpora. I conclude with an analysis of an alternative computing model based on renting instead of buying computer clusters.

1 Introduction

Over the past couple of decades, the field of computational linguistics (and more broadly, human language technologies) has seen the emergence and later dominance of empirical techniques and data-driven research. Concomitant with this trend is a coherent research thread that focuses on exploiting increasingly-large datasets. Banko and Brill (2001) were among the first to demonstrate the importance of dataset size as a significant factor governing prediction accuracy in a supervised machine learning

task. In fact, they argued that size of training set was perhaps more important than the choice of machine learning algorithm itself. Similarly, experiments in question answering have shown the effectiveness of simple pattern-matching techniques when applied to large quantities of data (Brill et al., 2001; Dumais et al., 2002). More recently, this line of argumentation has been echoed in experiments with Web-scale language models. Brants et al. (2007) showed that for statistical machine translation, a simple smoothing technique (dubbed *Stupid Backoff*) approaches the quality of the Kneser-Ney algorithm as the amount of training data increases, and with the simple method one can process significantly more data.

Challenges in scaling algorithms to increasingly-large datasets have become a serious issue for researchers. It is clear that datasets readily available today and the types of analyses that researchers wish to conduct have outgrown the capabilities of individual computers. The only practical recourse is to distribute the computation across multiple cores, processors, or machines. The consequences of failing to scale include misleading generalizations on artificially small datasets and limited practical applicability in real-world contexts, both undesirable.

This paper focuses on two barriers to developing scalable language processing algorithms: challenges associated with parallel programming and access to hardware. Google's MapReduce framework (Dean and Ghemawat, 2004) provides an attractive programming model for developing scalable algorithms, and with the release of Hadoop, an open-source implementation of MapReduce lead

by Yahoo, cost-effective cluster computing is within the reach of most academic research groups. It is emphasized that this work focuses on large-data algorithms from the perspective of academia—colleagues in commercial environments have long enjoyed the advantages of cluster computing. However, it is only recently that such capabilities have become *practical* for academic research groups. These points are illustrated by a case study in building large word co-occurrence matrices, a simple task that underlies many NLP algorithms.

The remainder of the paper is organized as follows: the next section overviews the MapReduce framework and why it provides a compelling solution to the issues sketched above. Section 3 introduces the task of building word co-occurrence matrices, which provides an illustrative case study. Two separate algorithms are presented in Section 4. The experimental setup is described in Section 5, followed by presentation of results in Section 6. Implications and generalizations are discussed following that. Before concluding, I explore an alternative model of computing based on renting instead of buying hardware, which makes cluster computing practical for *everyone*.

2 MapReduce

The only practical solution to large-data challenges today is to distribute the computation across multiple cores, processors, or machines. The development of parallel algorithms involves a number of tradeoffs. First is that of cost: a decision must be made between “exotic” hardware (e.g., large shared memory machines, InfiniBand interconnect) and commodity hardware. There is significant evidence (Barroso et al., 2003) that solutions based on the latter are more cost effective—and for resource-constrained academic NLP groups, commodity hardware is often the only practical route.

Given appropriate hardware, researchers must still contend with the challenge of developing software. Quite simply, parallel programming is difficult. Due to communication and synchronization issues, concurrent operations are notoriously challenging to reason about. Reliability and fault tolerance become important design considerations on clusters containing large numbers of unreliable com-

modity parts. With traditional parallel programming models (e.g., MPI), the developer shoulders the burden of explicitly managing concurrency. As a result, a significant amount of the programmer’s attention is devoted to system-level details, leaving less time for focusing on the actual problem.

Recently, MapReduce (Dean and Ghemawat, 2004) has emerged as an attractive alternative to existing parallel programming models. The MapReduce abstraction shields the programmer from having to explicitly worry about system-level issues such as synchronization, inter-process communication, and fault tolerance. The runtime is able to transparently distribute computations across large clusters of commodity hardware with good scaling characteristics. This frees the programmer to focus on solving the problem at hand.

MapReduce builds on the observation that many information processing tasks have the same basic structure: a computation is applied over a large number of records (e.g., Web pages, bitext pairs, or nodes in a graph) to generate partial results, which are then aggregated in some fashion. Naturally, the per-record computation and aggregation function vary according to task, but the basic structure remains fixed. Taking inspiration from higher-order functions in functional programming, MapReduce provides an abstraction at the point of these two operations. Specifically, the programmer defines a “mapper” and a “reducer” with the following signatures:

$$\begin{aligned} \text{map: } & (k_1, v_1) \rightarrow [(k_2, v_2)] \\ \text{reduce: } & (k_2, [v_2]) \rightarrow [(k_3, v_3)] \end{aligned}$$

Key-value pairs form the basic data structure in MapReduce. The mapper is applied to every input key-value pair to generate an arbitrary number of intermediate key-value pairs ($[\dots]$ is used to denote a list). The reducer is applied to all values associated with the same intermediate key to generate output key-value pairs. This two-stage processing structure is illustrated in Figure 1.

Under the framework, a programmer needs only to provide implementations of the mapper and reducer. On top of a distributed file system (Ghemawat et al., 2003), the runtime transparently handles all other aspects of execution, on clusters ranging from a few to a few thousand nodes. The runtime is responsible for scheduling map and reduce

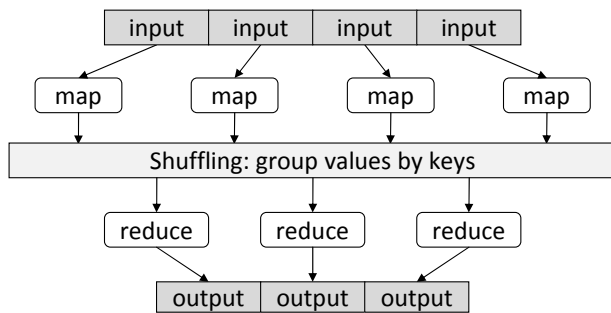


Figure 1: Illustration of the MapReduce framework: the “mapper” is applied to all input records, which generates results that are aggregated by the “reducer”. The runtime groups together values by keys.

workers on commodity hardware assumed to be unreliable, and thus is tolerant to various faults through a number of error recovery mechanisms. In the distributed file system, data blocks are stored on the local disks of machines in the cluster—the MapReduce runtime handles the scheduling of mappers on machines where the necessary data resides. It also manages the potentially very large sorting problem between the map and reduce phases whereby intermediate key-value pairs must be grouped by key.

As an optimization, MapReduce supports the use of “combiners”, which are similar to reducers except that they operate directly on the output of mappers (in memory, before intermediate output is written to disk). Combiners operate in isolation on each node in the cluster and cannot use partial results from other nodes. Since the output of mappers (i.e., the key-value pairs) must ultimately be shuffled to the appropriate reducer over a network, combiners allow a programmer to aggregate partial results, thus reducing network traffic. In cases where an operation is both associative and commutative, reducers can directly serve as combiners.

Google’s proprietary implementation of MapReduce is in C++ and not available to the public. However, the existence of Hadoop, an open-source implementation in Java spearheaded by Yahoo, allows anyone to take advantage of MapReduce. The growing popularity of this technology has stimulated a flurry of recent work, on applications in machine learning (Chu et al., 2006), machine translation (Dyer et al., 2008), and document retrieval (Elsayed et al., 2008).

3 Word Co-occurrence Matrices

To illustrate the arguments outlined above, I present a case study using MapReduce to build word co-occurrence matrices from large corpora, a common task in natural language processing. Formally, the co-occurrence matrix of a corpus is a square $N \times N$ matrix where N corresponds to the number of unique words in the corpus. A cell m_{ij} contains the number of times word w_i co-occurs with word w_j within a specific context—a natural unit such as a sentence or a certain window of m words (where m is an application-dependent parameter). Note that the upper and lower triangles of the matrix are identical since co-occurrence is a symmetric relation.

This task is quite common in corpus linguistics and provides the starting point to many other algorithms, e.g., for computing statistics such as pointwise mutual information (Church and Hanks, 1990), for unsupervised sense clustering (Schütze, 1998), and more generally, a large body of work in lexical semantics based on distributional profiles, dating back to Firth (1957) and Harris (1968). The task also has applications in information retrieval, e.g., (Schütze and Pedersen, 1998; Xu and Croft, 1998), and other related fields as well. More generally, this problem relates to the task of estimating distributions of discrete events from a large number of observations (more on this in Section 7).

It is obvious that the space requirement for this problem is $O(N^2)$, where N is the size of the vocabulary, which for real-world English corpora can be hundreds of thousands of words. The computation of the word co-occurrence matrix is quite simple if the entire matrix fits into memory—however, in the case where the matrix is too big to fit in memory, a naive implementation can be very slow as memory is paged to disk. For large corpora, one needs to optimize disk access and avoid costly seeks. As illustrated in the next section, MapReduce handles exactly these issues transparently, allowing the programmer to express the algorithm in a straightforward manner.

A bit more discussion of the task before moving on: in many applications, researchers have discovered that building the complete word co-occurrence matrix may not be necessary. For example, Schütze (1998) discusses feature selection

techniques in defining context vectors; Mohammad and Hirst (2006) present evidence that conceptual distance is better captured via distributional profiles mediated by thesaurus categories. These objections, however, miss the point—the focus of this paper is on practical cluster computing for academic researchers; this particular task serves merely as an illustrative example. In addition, for rapid prototyping, it may be useful to start with the complete co-occurrence matrix (especially if it can be built efficiently), and then explore how algorithms can be optimized for specific applications and tasks.

4 MapReduce Implementation

This section presents two MapReduce algorithms for building word co-occurrence matrices for large corpora. The goal is to illustrate how the problem can be concisely captured in the MapReduce programming model, and how the runtime hides many of the system-level details associated with distributed computing.

Pseudo-code for the first, more straightforward, algorithm is shown in Figure 2. Unique document ids and the corresponding texts make up the input key-value pairs. The mapper takes each input document and emits intermediate key-value pairs with each co-occurring word pair as the key and the integer one as the value. In the pseudo-code, EMIT denotes the creation of an intermediate key-value pair that is collected (and appropriately sorted) by the MapReduce runtime. The reducer simply sums up all the values associated with the same co-occurring word pair, arriving at the absolute counts of the joint event in the corpus (corresponding to each cell in the co-occurrence matrix).

For convenience, I refer to this algorithm as the “pairs” approach. Since co-occurrence is a symmetric relation, it suffices to compute half of the matrix. However, for conceptual clarity and to generalize to instances where the relation may not be symmetric, the algorithm computes the entire matrix.

The Java implementation of this algorithm is quite concise—less than fifty lines long. Notice the MapReduce runtime guarantees that all values associated with the same key will be gathered together at the reduce stage. Thus, the programmer does not need to explicitly manage the collection and distribution of

```

1: procedure MAP1( $n, d$ )
2:   for all  $w \in d$  do
3:     for all  $u \in \text{NEIGHBORS}(w)$  do
4:       EMIT( $(w, u), 1$ )

1: procedure REDUCE1( $p, [v_1, v_2, \dots]$ )
2:   for all  $v \in [v_1, v_2, \dots]$  do
3:      $sum \leftarrow sum + v$ 
4:   EMIT( $p, sum$ )

```

Figure 2: Pseudo-code for the “pairs” approach for computing word co-occurrence matrices.

```

1: procedure MAP2( $n, d$ )
2:   INITIALIZE( $H$ )
3:   for all  $w \in d$  do
4:     for all  $u \in \text{NEIGHBORS}(w)$  do
5:        $H\{u\} \leftarrow H\{u\} + 1$ 
6:   EMIT( $w, H$ )

1: procedure REDUCE2( $w, [H_1, H_2, H_3, \dots]$ )
2:   INITIALIZE( $H_f$ )
3:   for all  $H \in [H_1, H_2, H_3, \dots]$  do
4:     MERGE( $H_f, H$ )
5:   EMIT( $w, H_f$ )

```

Figure 3: Pseudo-code for the “stripes” approach for computing word co-occurrence matrices.

partial results across a cluster. In addition, the programmer does not need to explicitly partition the input data and schedule workers. This example shows the extent to which distributed processing can be dominated by system issues, and how an appropriate abstraction can significantly simplify development.

It is immediately obvious that Algorithm 1 generates an immense number of key-value pairs. Although this can be mitigated with the use of a combiner (since addition is commutative and associative), the approach still results in a large amount of network traffic. An alternative approach is presented in Figure 3, first reported in Dyer et al. (2008). The major difference is that counts of co-occurring words are first stored in an associative array (H). The output of the mapper is a number of key-value pairs with words as keys and the corresponding associative arrays as the values. The reducer performs an element-wise sum of all associative arrays with the same key (denoted by the function MERGE), thus ac-

cumulating counts that correspond to the same cell in the co-occurrence matrix. Once again, a combiner can be used to cut down on the network traffic by merging partial results. In the final output, each key-value pair corresponds to a row in the word co-occurrence matrix. For convenience, I refer to this as the “stripes” approach.

Compared to the “pairs” approach, the “stripes” approach results in far fewer intermediate key-value pairs, although each is significantly larger (and there is overhead in serializing and deserializing associative arrays). A critical assumption of the “stripes” approach is that at any point in time, each associative array is small enough to fit into memory (otherwise, memory paging may result in a serious loss of efficiency). This is true for most corpora, since the size of the associative array is bounded by the vocabulary size. Section 6 compares the efficiency of both algorithms.¹

5 Experimental Setup

Work reported in this paper used the English Gigaword corpus (version 3),² which consists of newswire documents from six separate sources, totaling 7.15 million documents (6.8 GB compressed, 19.4 GB uncompressed). Some experiments used only documents from the Associated Press Worldstream (APW), which contains 2.27 million documents (1.8 GB compressed, 5.7 GB uncompressed). By LDC’s count, the entire collection contains approximately 2.97 billion words.

Prior to working with Hadoop, the corpus was first preprocessed. All XML markup was removed, followed by tokenization and stopword removal using standard tools from the Lucene search engine. All tokens were replaced with unique integers for a more efficient encoding. The data was then packed into a Hadoop-specific binary file format. The entire Gigaword corpus took up 4.69 GB in this format; the APW sub-corpus, 1.32 GB.

Initial experiments used Hadoop version 0.16.0 running on a 20-machine cluster (1 master, 19 slaves). This cluster was made available to the Uni-

¹Implementations of both algorithms are included in Cloud⁹, an open source Hadoop library that I have been developing to support research and education, available from my homepage.

²LDC catalog number LDC2007T07

versity of Maryland as part of the Google/IBM Academic Cloud Computing Initiative. Each machine has two single-core processors (running at either 2.4 GHz or 2.8 GHz), 4 GB memory. The cluster has an aggregate storage capacity of 1.7 TB. Hadoop ran on top of a virtualization layer, which has a small but measurable impact on performance; see (Barham et al., 2003). Section 6 reports experimental results using this cluster; Section 8 explores an alternative model of computing based on “renting cycles”.

6 Results

First, I compared the running time of the “pairs” and “stripes” approaches discussed in Section 4. Running times on the 20-machine cluster are shown in Figure 4 for the APW section of the Gigaword corpus: the x -axis shows different percentages of the sub-corpus (arbitrarily selected) and the y -axis shows running time in seconds. For these experiments, the co-occurrence window was set to two, i.e., w_i is said to co-occur with w_j if they are no more than two words apart (after tokenization and stopword removal).

Results demonstrate that the stripes approach is far more efficient than the pairs approach: 666 seconds (11m 6s) compared to 3758 seconds (62m 38s) for the entire APW sub-corpus (improvement by a factor of 5.7). On the entire sub-corpus, the mappers in the pairs approach generated 2.6 billion intermediate key-value pairs totally 31.2 GB. After the combiners, this was reduced to 1.1 billion key-value pairs, which roughly quantifies the amount of data involved in the shuffling and sorting of the keys. On the other hand, the mappers in the stripes approach generated 653 million intermediate key-value pairs totally 48.1 GB; after the combiners, only 28.8 million key-value pairs were left. The stripes approach provides more opportunities for combiners to aggregate intermediate results, thus greatly reducing network traffic in the sort and shuffle phase.

Figure 4 also shows that both algorithms exhibit highly desirable scaling characteristics—linear in the corpus size. This is confirmed by a linear regression applied to the running time data, which yields R^2 values close to one. Given that the stripes algorithm is more efficient, it is used in the remainder of the experiments.

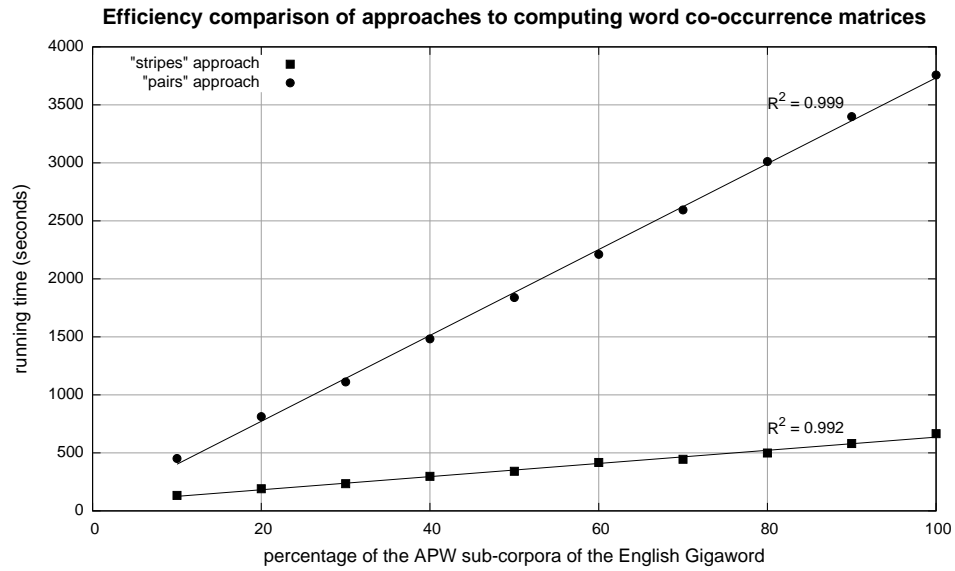


Figure 4: Running time of the two algorithms (“strips” vs. “pairs”) for computing word co-occurrence matrices on the APW section of the Gigaword corpus. The cluster used for this experiment contains 20 machines, each with two single-core processors.

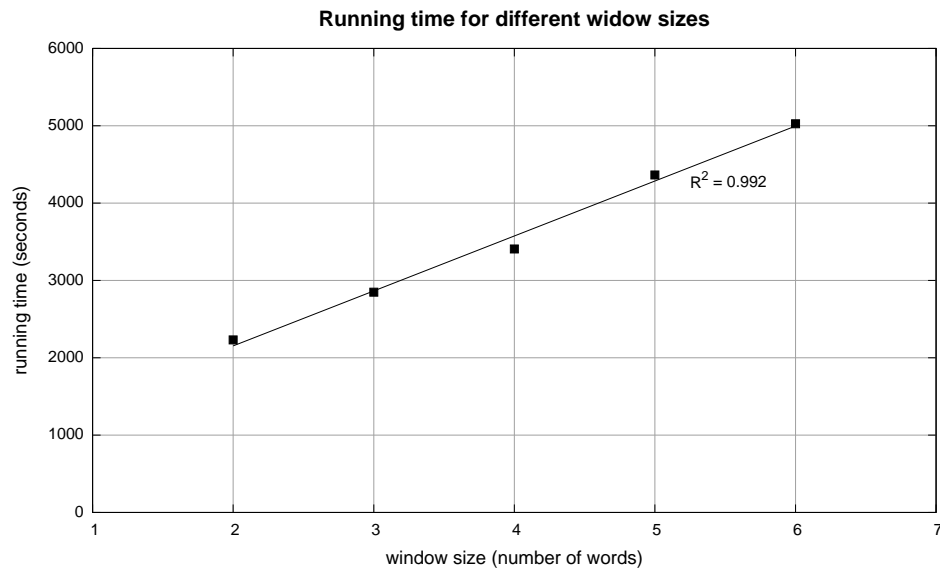


Figure 5: Running times for computing word co-occurrence matrices from the entire Gigaword corpus with varying window sizes. The cluster used for this experiment contains 20 machines, each with two single-core processors.

With a window size of two, computing the word co-occurrence matrix for the entire Gigaword corpus (7.15 million documents) takes 37m 11s on the 20-machine cluster. Figure 5 shows the running time as a function of window size. With a window of six words, running time on the complete Gigaword corpus rises to 1h 23m 45s. Once again, the stripes algorithm exhibits the highly desirable characteristic of linear scaling in terms of window size, as confirmed by the linear regression with an R^2 value very close to one.

7 Discussion

The elegance of the programming model and good scaling characteristics of resulting implementations make MapReduce a compelling tool for a variety of natural language processing tasks. In fact, MapReduce excels at a large class of problems in NLP that involves estimating probability distributions of discrete events from a large number of observations according to the maximum likelihood criterion:

$$P_{MLE}(B|A) = \frac{c(A, B)}{c(A)} = \frac{c(A, B)}{\sum_{B'} c(A, B')} \quad (1)$$

In practice, it matters little whether these events are words, syntactic categories, word alignment links, or any construct of interest to researchers. Absolute counts in the stripes algorithm presented in Section 4 can be easily converted into conditional probabilities by a final normalization step. Recently, Dyer et al. (2008) used this approach for word alignment and phrase extraction in statistical machine translation. Of course, many applications require smoothing of the estimated distributions—this problem also has known solutions in MapReduce (Brants et al., 2007).

Synchronization is perhaps the single largest bottleneck in distributed computing. In MapReduce, this is handled in the shuffling and sorting of key-value pairs between the map and reduce phases. Development of efficient MapReduce algorithms critically depends on careful control of intermediate output. Since the network link between different nodes in a cluster is by far the component with the largest latency, any reduction in the size of intermediate output or a reduction in the number of key-value pairs will have significant impact on efficiency.

8 Computing on Demand

The central theme of this paper is practical cluster computing for NLP researchers in the academic environment. I have identified two key aspects of what it means to be “practical”: the first is an appropriate programming model for simplifying concurrency management; the second is access to hardware resources. The Hadoop implementation of MapReduce addresses the first point and to a large extent the second point as well. The cluster used for experiments in Section 6 is modest by today’s standards and within the capabilities of many academic research groups. It is not even a requirement for the computers to be rack-mounted units in a machine room (although that is clearly preferable); there are plenty of descriptions on the Web about Hadoop clusters built from a handful of desktop machines connected by gigabit Ethernet.

Even without access to hardware, cluster computing remains within the reach of resource-constrained academics. “Utility computing” is an emerging concept whereby anyone can provision clusters on demand from a third-party provider. Instead of upfront capital investment to acquire a cluster and re-occurring maintenance and administration costs, one could “rent” computing cycles as they are needed—this is not a new idea (Rappa, 2004). One such service is provided by Amazon, called Elastic Compute Cloud (EC2).³ With EC2, researchers could dynamically create a Hadoop cluster on-the-fly and tear down the cluster once experiments are complete. To demonstrate the use of this technology, I replicated some of the previous experiments on EC2 to provide a case study of this emerging model of computing.

Virtualized computation units in EC2 are called instances. At the time of these experiments, the basic instance offers, according to Amazon, 1.7 GB of memory, 1 EC2 Compute Unit (1 virtual core with 1 EC2 Compute Unit), and 160 GB of instance storage. Each instance-hour costs \$0.10 (all prices given in USD). Computational resources are simply charged by the instance-hour, so that a ten-instance cluster for ten hours costs the same as a hundred-instance cluster for one hour (both \$10)—the Amazon infrastructure allows one to dynamically provision and release resources as necessary. This is at-

³<http://www.amazon.com/ec2>

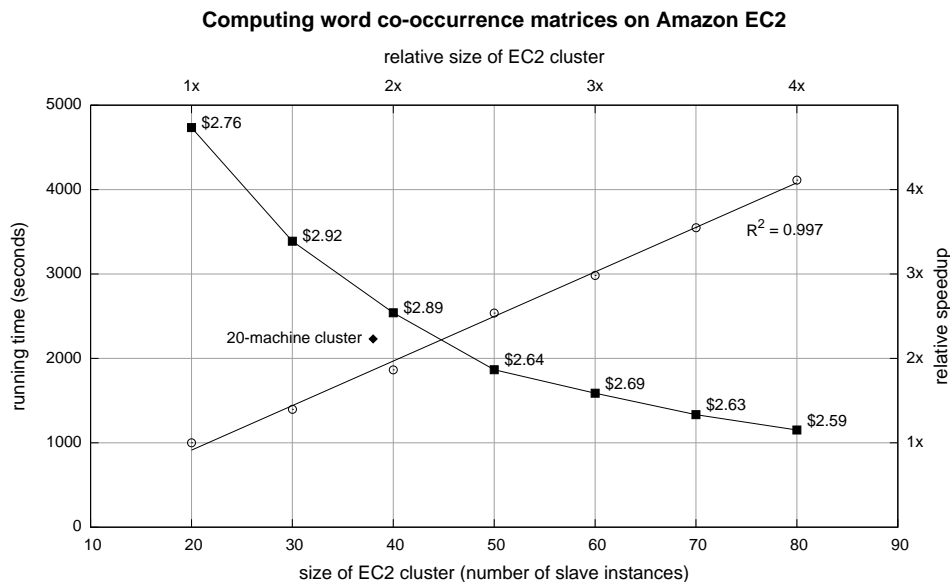


Figure 6: Running time analysis on Amazon EC2 with various cluster sizes; solid squares are annotated with the cost of each experiment. Alternate axes (circles) plot scaling characteristics in terms increasing cluster size.

tractive for researchers, who could on a limited basis allocate clusters much larger than they could otherwise afford if forced to purchase the hardware outright. Through virtualization technology, Amazon is able to parcel out allotments of processor cycles while maintaining high overall utilization across a data center and exploiting economies of scale.

Using EC2, I built word co-occurrence matrices from the entire English Gigaword corpus (window of two) on clusters of various sizes, ranging from 20 slave instances all the way up to 80 slave instances. The entire cluster consists of the slave instances plus a master controller instance that serves as the job submission queue; the clusters ran Hadoop version 0.17.0 (the latest release at the time these experiments were conducted). Running times are shown in Figure 6 (solid squares), with varying cluster sizes on the x -axis. Each data point is annotated with the cost of running the complete experiment.⁴ Results show that computing the complete word co-occurrence matrix costs, quite literally, a couple of dollars—certainly affordable by any academic researcher without access to hardware. For reference, Figure 6 also plots the running time of the same experiment on the 20-machine cluster used

⁴Note that Amazon bills in whole instance-hour increments; these figures assume fractional accounting.

in Section 6 (which contains 38 worker cores, each roughly comparable to an instance).

The alternate set of axes in Figure 6 shows the scaling characteristics of various cluster sizes. The circles plot the relative size and speedup of the EC2 experiments, with respect to the 20-slave cluster. The results show highly desirable linear scaling characteristics.

The above figures include only the cost of running the instances. One must additionally pay for bandwidth when transferring data in and out of EC2. At the time these experiments were conducted, Amazon charged \$0.10 per GB for data transferred in and \$0.17 per GB for data transferred out. To complement EC2, Amazon offers persistent storage via the Simple Storage Service (S3),⁵ at a cost of \$0.15 per GB per month. There is no charge for data transfers between EC2 and S3. The availability of this service means that one can choose between paying for data transfer or paying for persistent storage on a cyclic basis—the tradeoff naturally depends on the amount of data and its permanence.

The cost analysis presented above assumes optimally-efficient use of Amazon’s services; end-to-end cost might better quantify real-world usage conditions. In total, the experiments reported in this

⁵<http://www.amazon.com/s3>

section resulted in a bill of approximately thirty dollars. The figure includes all costs associated with instance usage and data transfer costs. It also includes time taken to learn the Amazon tools (I previously had no experience with either EC2 or S3) and to run preliminary experiments on smaller datasets (before scaling up to the complete corpus). The lack of fractional accounting on instance-hours contributed to the larger-than-expected costs, but such wastage would naturally be reduced with more experiments and higher sustained use. Overall, these costs appear to be very reasonable, considering that the largest cluster in these experiments (1 master + 80 slave instances) might be too expensive for most academic research groups to own and maintain.

Consider another example that illustrates the possibilities of utility computing. Brants et al. (2007) described experiments on building language models with increasingly-large corpora using MapReduce. Their paper reported experiments on a corpus containing 31 billion tokens (about an order of magnitude larger than the English Gigaword): on 400 machines, the model estimation took 8 hours.⁶ With EC2, such an experiment would cost a few hundred dollars—sufficiently affordable that availability of data becomes the limiting factor, not computational resources themselves.

The availability of “computing-on-demand” services and Hadoop make cluster computing practical for academic researchers. Although Amazon is currently the most prominent provider of such services, they are not the sole player in an emerging market—in the future there will be a vibrant market with many competing providers. Considering the tradeoffs between “buying” and “renting”, I would recommend the following model for an academic research group: purchase a modest cluster for development and for running smaller experiments; use a computing-on-demand service for scaling up and for running larger experiments (since it would be more difficult to economically justify a large cluster if it does not receive high sustained utilization).

If the concept of utility computing takes hold, it would have a significant impact on computer science research in general: the natural implication is

⁶Brants et al. were affiliated with Google, so access to hardware was not an issue.

that algorithms should not only be analyzed in traditional terms such as asymptotic complexity, but also in terms of monetary costs, in relationship to dataset and cluster size. One can argue that cost is a more direct and practical measure of algorithmic efficiency.

9 Conclusion

This paper addresses two challenges faced by academic research groups in scaling up natural language processing algorithms to large corpora: the lack of an appropriate programming model for expressing the problem and the difficulty in getting access to hardware. With this case study in building word co-occurrence matrices from large corpora, I demonstrate that MapReduce, via the open source Hadoop implementation, provides a compelling solution. A large class of algorithms in computational linguistics can be readily expressed in MapReduce, and the resulting code can be transparently distributed across commodity clusters. Finally, the “cycle-renting” model of computing makes access to large clusters affordable to researchers with limited resources. Together, these developments dramatically lower the entry barrier for academic researchers who wish to explore large-data issues.

Acknowledgments

This work was supported by the Intramural Research Program of the NIH, National Library of Medicine; NSF under awards IIS-0705832 and IIS-0836560; DARPA/IPTO Contract No. HR0011-06-2-0001 under the GALE program. Any opinions, findings, conclusions, or recommendations expressed in this paper are the author’s and do not necessarily reflect those of the sponsors. I would like to thank Yahoo! for leading the development of Hadoop, IBM and Google for hardware support via the Academic Cloud Computing Initiative (ACCI), and Amazon for EC2/S3 support. This paper provides a neutral evaluation of EC2 and S3, and should not be interpreted as endorsement for the commercial services offered by Amazon. I wish to thank Philip Resnik and Doug Oard for comments on earlier drafts of this paper, and Ben Shneiderman for helpful editing suggestions. I am, as always, grateful to Esther and Kiri for their kind support.

References

- Michele Banko and Eric Brill. 2001. Scaling to very very large corpora for natural language disambiguation. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics (ACL 2001)*, pages 26–33, Toulouse, France.
- Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, and Andrew Warfield. 2003. Xen and the art of virtualization. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP-03)*, pages 164–177, Bolton Landing, New York.
- Luiz André Barroso, Jeffrey Dean, and Urs Hölzle. 2003. Web search for a planet: The Google cluster architecture. *IEEE Micro*, 23(2):22–28.
- Thorsten Brants, Ashok C. Popat, Peng Xu, Franz J. Och, and Jeffrey Dean. 2007. Large language models in machine translation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 858–867, Prague, Czech Republic.
- Eric Brill, Jimmy Lin, Michele Banko, Susan Dumais, and Andrew Ng. 2001. Data-intensive question answering. In *Proceedings of the Tenth Text REtrieval Conference (TREC 2001)*, pages 393–400, Gaithersburg, Maryland.
- Cheng-Tao Chu, Sang Kyun Kim, Yi-An Lin, YuanYuan Yu, Gary Bradski, Andrew Ng, and Kunle Olukotun. 2006. Map-Reduce for machine learning on multicore. In *Advances in Neural Information Processing Systems 19 (NIPS 2006)*, pages 281–288, Vancouver, British Columbia, Canada.
- Kenneth W. Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29.
- Jeffrey Dean and Sanjay Ghemawat. 2004. MapReduce: Simplified data processing on large clusters. In *Proceedings of the 6th Symposium on Operating System Design and Implementation (OSDI 2004)*, pages 137–150, San Francisco, California.
- Susan Dumais, Michele Banko, Eric Brill, Jimmy Lin, and Andrew Ng. 2002. Web question answering: Is more always better? In *Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR 2002)*, pages 291–298, Tampere, Finland.
- Chris Dyer, Aaron Cordova, Alex Mont, and Jimmy Lin. 2008. Fast, easy, and cheap: Construction of statistical machine translation models with MapReduce. In *Proceedings of the Third Workshop on Statistical Machine Translation at ACL 2008*, pages 199–207, Columbus, Ohio.
- Tamer Elsayed, Jimmy Lin, and Douglas Oard. 2008. Pairwise document similarity in large collections with MapReduce. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics (ACL 2008), Companion Volume*, pages 265–268, Columbus, Ohio.
- John R. Firth. 1957. A synopsis of linguistic theory 1930–55. In *Studies in Linguistic Analysis, Special Volume of the Philological Society*, pages 1–32. Blackwell, Oxford.
- Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. 2003. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP-03)*, pages 29–43, Bolton Landing, New York.
- Zelig S. Harris. 1968. *Mathematical Structures of Language*. Wiley, New York.
- Saif Mohammad and Graeme Hirst. 2006. Distributional measures of concept-distance: A task-oriented evaluation. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP 2006)*, pages 35–43, Sydney, Australia.
- Michael A. Rappa. 2004. The utility business model and the future of computing services. *IBM Systems Journal*, 34(1):32–42.
- Hinrich Schütze and Jan O. Pedersen. 1998. A cooccurrence-based thesaurus and two applications to information retrieval. *Information Processing and Management*, 33(3):307–318.
- Hinrich Schütze. 1998. Automatic word sense discrimination. *Computational Linguistics*, 24(1):97–123.
- Jinxi Xu and W. Bruce Croft. 1998. Corpus-based stemming using cooccurrence of word variants. *ACM Transactions on Information Systems*, 16(1):61–81.