

# Graph Data Mining with Arabesque

Eslam Hussein<sup>△</sup>, Abdurrahman Ghanem<sup>△</sup>, Vinicius Vitor dos Santos Dias<sup>♣</sup>,  
Carlos H. C. Teixeira<sup>♣</sup>, Ghadeer AbuOda<sup>△</sup>,  
Marco Serafini<sup>△</sup>, Georgios Siganos<sup>△</sup>, Gianmarco De Francisci Morales<sup>△</sup>,  
Ashraf AboulNaga<sup>△</sup>, Mohammed Zaki<sup>♣</sup>  
<sup>△</sup>Qatar Computing Research Institute - HBKU, <sup>♣</sup>Universidade Federal de Minas Gerais,  
<sup>♣</sup>Rensselaer Polytechnic Institute

## ABSTRACT

Graph data mining is defined as searching in an input graph for all subgraphs that satisfy some property that makes them interesting to the user. Examples of graph data mining problems include frequent subgraph mining, counting motifs, and enumerating cliques. These problems differ from other graph processing problems such as PageRank or shortest path in that graph data mining requires searching through an exponential number of subgraphs. Most current parallel graph analytics systems do not provide good support for graph data mining. One notable exception is Arabesque, a system that was built specifically to support graph data mining. Arabesque provides a simple programming model to express graph data mining computations, and a highly scalable and efficient implementation of this model, scaling to billions of subgraphs on hundreds of cores. This demonstration will showcase the Arabesque system, focusing on the end-user experience and showing how Arabesque can be used to simply and efficiently solve practical graph data mining problems that would be difficult with other systems.

## 1. INTRODUCTION

Graph data is playing an increasingly important role in many fields such as biology, e-commerce, and social network analysis. Graph data appears in on-line operations, such as representing new “friend” relationships in a social network, and in analytics, such as predicting users who can become friends. The increase in the size of graph data and the complexity of workloads on this data have led to the development of parallel and distributed systems that support high throughput graph updates and retrieval, such as TAO [2], as well as systems that support large scale graph analytics, such as Pregel [10], GraphLab [9], and EmptyHeaded [1].

Most parallel graph analytics systems support computations that produce succinct properties of the graph or of individual vertices, or produce a small number of result subgraphs. Examples of these computations include PageRank, shortest path, and counting cliques. These systems do not provide good support for *graph data mining*, which we define as searching through the exponential number of subgraphs of an input graph to find subgraphs that satisfy some property that makes them interesting to the user. Examples of graph

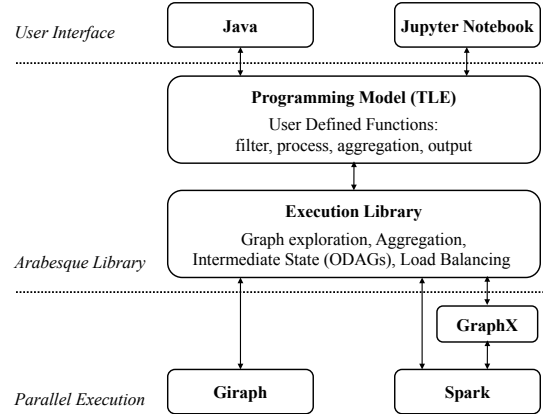


Figure 1: Overview of Arabesque.

data mining problems include frequent subgraph mining, counting motifs, and enumerating cliques or quasi-cliques (as opposed to only counting them). Some recent systems, such as Arabesque [15] and NScale [14], adopt a model where subgraphs (as opposed to vertices) are first class citizens in the computation, which enables better support for graph data mining.

This demonstration will showcase the Arabesque system, which was built with the specific goal of supporting efficient and scalable graph data mining. The technical details of Arabesque are presented elsewhere [15], and the code is available as open source<sup>1</sup>. This demonstration will focus on the end-user experience of Arabesque, and how it can be used to solve interesting and important graph data mining problems that would be difficult with other graph analytics systems. Participants in the demonstration will see how Arabesque fits within a typical data analytics toolchain, and will get a sense of the usability, programmability, and efficiency of Arabesque. The demonstration will be centered around three practical applications: finding frequent subgraphs in protein databases, analyzing motifs to measure the reaction to various events on Twitter, and analyzing cliques to detect communities of common interest among buyers on Amazon. The three applications are based on three different types of graph data mining problems supported by Arabesque: frequent subgraph mining, counting motifs, and enumerating cliques. In the next section, we present a brief overview of Arabesque, and in Section 3 we present the details of these applications.

## 2. OVERVIEW OF ARABESQUE

As mentioned earlier, graph data mining is characterized by enumerating the exponential number of subgraphs of an input graph

<sup>1</sup><http://arabesque.io>

and searching for patterns in these subgraphs. The Arabesque system [15] (Figure 1) is designed to support parallel graph data mining on hundreds of CPU cores in multiple servers (also referred to as worker nodes). A fundamental assumption made by Arabesque is that the input graph fits in the main memory of a single worker node, and can be replicated on all worker nodes. Today, the main memory of servers is typically in the 256GB to 2TB range, so this assumption covers a large fraction of graph data sets. Arabesque still needs to address the challenge of managing the exponentially sized intermediate state, which does *not* fit in the memory of a single worker node. Another challenge faced by Arabesque is distributing the computation to the CPU cores in a scalable, efficient, and load balanced way.

To address these challenges, Arabesque uses a programming model that can express graph data mining problems in a simple and succinct way, and is amenable to easy distribution on multiple cores. Arabesque also provides a scalable and efficient implementation of this programming model that works on top of parallel dataflow platforms such as Giraph<sup>2</sup> and Spark [16]. We present the Arabesque programming model and implementation next.

## 2.1 Programming Model

The Arabesque programming model is designed to support the automatic graph exploration required for graph data mining. It is based on a paradigm that we call *think like an embedding*, or *TLE*. An *embedding* is a subgraph representing an instance of a *pattern* of interest in the graph data mining problem, and a key characteristics of graph data mining is that we are interested in producing all output embeddings. For example, consider frequent subgraph mining, in which we want to find all instances of frequently occurring subgraph patterns. If subgraphs with, say, a vertex labeled  $A$  connected to a vertex labeled  $B$  connected to a vertex labeled  $C$  occur frequently in the input graph, we are interested not only in finding that  $A - B - C$  is frequent but also in producing all instances of  $A - B - C$ , say  $a_1 - b_1 - c_1, a_2 - b_2 - c_2, \dots, a_n - b_n - c_n$ . In this example,  $A - B - C$  is the *pattern*, and the instances  $a_i - b_i - c_i$  are the *embeddings* of this pattern, where lowercase letters indicate vertex ids of the input graph. A similar distinction between patterns and embeddings can be found in other graph data mining problems.

In the TLE programming model of Arabesque, the user provides two functions that accept one embedding as an argument: the *filter* function and the *process* function. The filter function is used to prune this search space: it takes an embedding as input and returns a boolean value indicating whether the embedding should be processed or not. The process function is used to analyze an embedding and generate the output required by the graph mining algorithm: it takes an embedding as input, processes the embedding as required by the graph data mining algorithm, and typically outputs a set of user-defined values to HDFS.

Arabesque explores the input graph in a series of bulk synchronous processing (BSP) steps, and maintains a set of candidate embeddings at each step. In the first step, the individual vertices of the input graph are the candidate embeddings, and in each subsequent step, each candidate embedding is expanded by adding its neighbors to it one by one to create larger candidate embeddings. In each step, Arabesque calls the filter function on all candidate embeddings, and discards the embeddings for which filter returns false from the candidate set. Arabesque then calls the process functions on embeddings remaining in the candidate set, and further expands these embeddings in subsequent BSP steps.

In addition to the filter and process functions, Arabesque allows

the user to specify other functions such as an aggregation filter function and an aggregation process function, which filter and process embeddings at the beginning of a BSP step based on aggregate information about all the embeddings of a pattern found in the previous step.

## 2.2 Implementation

A key characteristic of the TLE model is that there are no dependencies among embeddings. Each embedding can be filtered, processed, and expanded independently of other embeddings within a BSP step. Embeddings may be aggregated by pattern at the end of a BSP step, which introduces dependencies, but there are no dependencies within a step. The lack of dependencies enables Arabesque to utilize a coordination free strategy to avoid redundant work while exploring the graph based on the concept of *embedding canonicity*. Each worker thread is assigned a set of embeddings to expand in each step, without coordinating with other worker threads. It is possible that two worker threads can generate the same embedding independently. Without additional controls, both workers would call the filter and process functions on the embedding, which is wasteful. To avoid this situation, Arabesque defines a notion of canonicity for embeddings, and worker threads discard embeddings that they generate that are not canonical. To ensure coordination free exploration, canonicity in Arabesque is defined in a careful way that allows each thread to independently test the canonicity of embeddings that it generates.

The ability to expand embeddings independently enables Arabesque to balance load very well among the worker threads. Contrast this, for example, to the traditional *think like a vertex* (TLV) paradigm used by graph analytics systems such as Pregel. In TLV, computation and state are expressed at the level of a vertex in the input graph. One could use TLV for graph exploration by storing at each vertex all embeddings that this vertex is part of. The vertex function can expand embedding by adding each neighbor of a vertex to the embeddings that this neighbor is not already a part of. New embeddings would have to be sent to all vertices that these embeddings contain. This further multiplies the number of embeddings generated by the system, exacerbating the main bottleneck of graph mining algorithms. In addition, highly connected vertices generate a disproportionately large number of embeddings during expansion, leading to load imbalance. In our experiments, we have observed TLV to be up to two orders of magnitude slower than TLE.

As Arabesque explores the graph, it generates an exponential number of embeddings. To reduce the memory required for storing these embeddings, Arabesque uses a compact data structure called *Overapproximated Directed Acyclic Graph* (ODAG) that compresses the canonical embeddings generated in a BSP step. It also uses a *two-level aggregation* technique to speed up aggregation by pattern.

The full technical details of Arabesque are presented in [15]. In that paper, we show that Arabesque scales to billions of subgraphs on hundreds of cores on multiple worker nodes. From a software engineering perspective, Arabesque is implemented as a library that can easily be ported to any parallel dataflow execution engine. We currently have versions of Arabesque that runs on top of Giraph<sup>3</sup> and Spark<sup>4</sup>. Note that the Giraph version does not use the TLV programming model that Giraph implements. Instead, Arabesque uses Giraph to coordinate the BSP steps, and for resource allocation, job management, and fault tolerance. The Spark version of Arabesque, which we use in this demo, also implements graph exploration in BSP steps, taking advantage of Spark’s ability to keep

<sup>2</sup><http://giraph.apache.org>

<sup>3</sup><https://github.com/Qatar-Computing-Research-Institute/Arabesque>

<sup>4</sup><https://github.com/viniciusvdias/Arabesque>

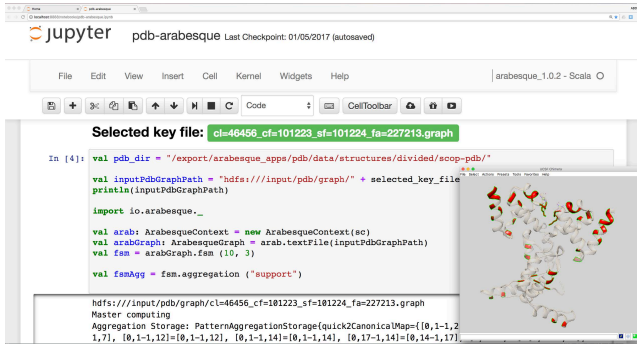


Figure 2: Mining protein structures with Arabesque.

results in memory between iterations through the use of resilient distributed datasets (RDDs). Arabesque on Spark has additional usability features such as a Jupyter notebook interface<sup>5</sup>, in addition to the Java interface, and the ability to pre-process the Arabesque input graphs (or post-process the output graphs) with GraphX [4]. Next, we describe the applications that we use to showcase these capabilities.

### 3. APPLICATIONS DEMONSTRATED

In this demonstration we will present three applications that use Arabesque: finding frequent structures in a protein database, modeling how users propagate information on Twitter in reaction to events, and finding communities of common interest among buyers from Amazon. The graph data mining algorithms used by the applications are, respectively, frequent subgraph mining, counting motifs, and enumerating cliques.

#### 3.1 Finding Frequent Structures in Proteins

The aim of this application is to identify and visualize frequently occurring patterns in the 3D structure of proteins. The input data for this application is the Protein Data Bank (PDB)<sup>6</sup>, which is an online repository containing the 3D structure of over 120K proteins.

The structure of a protein can be converted to a graph as follows. Each protein structure comprises a set of, say  $n$ , 3D coordinates, namely  $(x_i, y_i, z_i)$ , for  $i = 1, \dots, n$ . Each position or element  $i$ , also called an amino acid  $i$ , has a label. Let us denote  $a_i = (x_i, y_i, z_i)$ , and let  $l_i$  be its amino acid label. We can construct a graph for each protein, with a vertex for each amino acid, labeled with  $l_i$ . An edge exists between any two amino acids,  $a_i$  and  $a_j$ , if the Euclidean distance between them is below a given threshold, that is,  $\|a_i - a_j\|_2 \leq \theta$ , where  $\theta$  is the contact threshold (usually set to 7 angstroms, i.e.,  $7 \times 10^{-10}$  meters). The graphs for a set of proteins can be considered as disconnected components in an input graph to Arabesque.

Frequent subgraphs in these protein graphs represent frequently occurring patterns among the different protein structure. Identifying such frequently occurring patterns is important for many bioinformatics applications (e.g., [5, 6, 7, 11]). As a matter of fact, one of the ways to classify proteins in the PDB database is to group proteins by structure in a hierarchical organization, as in the structural classification of proteins (SCOP) project<sup>7</sup> [3], which uses manually identified, human-curated structural groupings. A scalable frequent subgraph mining implementation would be extremely helpful for bioinformatics applications on the PDB database, and in the demonstration we show how Arabesque can play this role.

<sup>5</sup><http://jupyter.org>

<sup>6</sup><http://www.rcsb.org/pdb>

<sup>7</sup><http://scop.berkeley.edu>

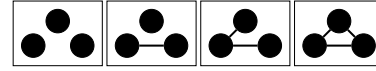


Figure 3: The 4 possible graphs on 3 vertices.

This application, like the other application in this demonstration, runs in the Jupyter notebook interface of Arabesque (Figure 2). The steps of the application are (1) extracting a relevant subset of the PDB database and constructing the input graph, (2) running frequent subgraph mining on the input graph, and (3) visualizing the frequent subgraphs found by Arabesque.

The goal of the application is to identify frequently occurring structures in a coherent subset of the PDB database. We use the SCOP classification to identify such a coherent subset. Thus, the first step of the application is for the user to choose one node of the SCOP hierarchy (referred to as a “SCOP key”), and to specify the threshold  $\theta$  for adding an edge between two amino acids. The application uses these user inputs to extract the relevant PDB data and construct the input graph.

Next, the application runs frequent subgraph mining on Arabesque. The user controls this step by specifying the required support and the maximum subgraph size explored. The output of this step is a set of subgraphs representing frequently occurring structures in the input protein data.

Finally, the user can visualize the frequently occurring structures identified by Arabesque. We use an external tool called UCSF Chimera<sup>8</sup> [13] for visualization. UCSF Chimera is a popular and powerful visualization tool that is widely used in the bioinformatics community, and users are able to take advantage of its full power for visualization.

#### 3.2 Twitter Reaction to Events

This application shows how to use subgraph counting to understand the reaction of social networks to exogenous events. In particular, we look at Twitter, and analyze the retweet network structure before and after the occurrence of an event. We demonstrate how  $k$ -profiles, i.e., the counts of all the subgraphs of size  $k$ , can detect the occurrence of an event.

Our dataset consist of several months of tweets related to an event in April 2016, when the Egyptian government announced that two islands in the Red Sea, which have a strategic location, were donated to Saudi Arabia. Online activists carried out a campaign on Twitter to call for demonstrations against the government’s decision.

For each week in the dataset, we build a retweet network  $G_i(V, E)$  such that the set of vertices  $V$  corresponds to the set of users active in our dataset, and there is an edge  $(u, v) \in E$  iff user  $u$  retweets user  $v$  (for the purposes of this demo, we ignore the direction of the edge). We filter out weeks where there is not enough activity (fewer than 500 edges in the graph).

$K$ -profiles are a useful generalization of triangle counting. For  $k = 3$ , there are 4 different possible subgraphs, as shown in Figure 3. For  $k = 4$ , there are 11 different ones. By counting the number of occurrences of each of these subgraph in the retweet network, we can extract a multi-dimensional vector that represents a *fingerprint* of the network. For simplicity, we only count connected subgraphs (e.g., only the last two graphs in Figure 3).

We extract the 4-profile for each retweet network, corresponding to each week in the dataset, and then compute the Euclidean distance between consecutive vectors. Figure 4 shows the results. The same pattern is clearly detectable also with 3-profiles (figure omitted for brevity).

<sup>8</sup><http://www.rbvi.ucsf.edu/chimera>

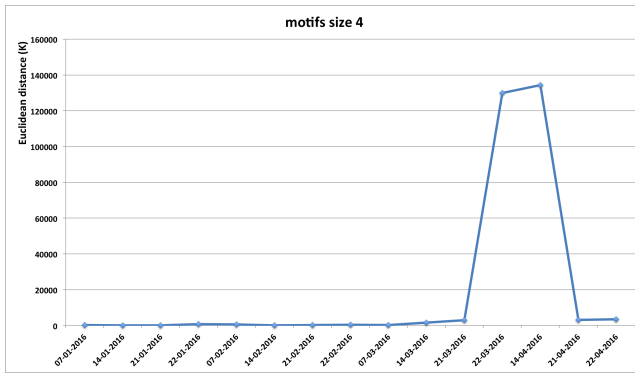


Figure 4: Event detection on Twitter with Arabesque.

The event is clearly distinguishable as a peak starting at the end of March and ending around April. This peak represents a change in the network activity that is easily to detect by monitoring the retweet patterns.

### 3.3 Finding Communities

In this use case, we use Arabesque to find communities in a co-purchase graph of the Amazon online shopping site. In particular, we use the  $k$ -clique percolation method [12], which is an established algorithm to find communities in a network. The method starts by identifying maximal cliques of size  $k$  or more. It then considers two cliques to be adjacent if they have  $k - 1$  common nodes. All cliques that are adjacent with each other, either directly or transitively, are considered to be part of the same community.

We applied the algorithm to find communities in the Amazon co-purchase graph of [8]. In this graph, vertices correspond to items on sale, and edges connect items that are frequently bought together in the same order. Therefore, communities represent similar items that are often bought together. We ran the method with a value  $k = 10$ . The steps of the application are: (1) building the input graph, (2) using Arabesque to find maximal cliques of size  $k$ , (3) identifying adjacent cliques and communities, and (4) visualizing some of the cliques.

Figure 5 shows some of the communities we found, which will be shown during the demo. Communities are enclosed in circles, and vertices of overlapping communities have the same color. We can see that some communities with large overlaps are similar. For example, the two communities in blue have large overlaps and they both contain essays, mainly related to personal wellness, spirituality, and home care. Other communities have a smaller overlap and are less related. The two communities in green have a small overlap and they are more diverse. The one on the left mainly contains essays on biology, history, traveling, and sociology, while the one on the right contains mainly books related to children. The overlap is only on two generic self-help books.

## 4. CONCLUSION

Arabesque is a scalable and efficient parallel system for graph data mining. It enables users to solve problems not easily solvable by other systems. This demonstration will present the user/programmer experience with Arabesque, focusing on complete applications that showcase different capabilities of the system.

## References

- [1] C. R. Aberger, S. Tu, K. Olukotun, and C. Ré. EmptyHeaded: A relational engine for graph processing. In *SIGMOD*, 2016.
- [2] N. Bronson et al. TAO: Facebook’s distributed data store for the social graph. In *USENIX Annual Technical Conference (ATC)*, 2013.

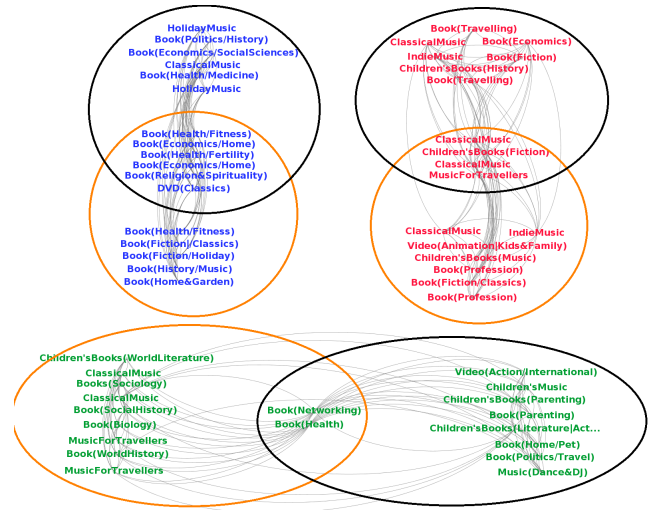


Figure 5: Screenshot of the community finding application on Amazon.

- [3] N. K. Fox, S. E. Brenner, and J.-M. Chandonia. SCOPe: Structural Classification of Proteins – extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Research*, 42(D1), 2014.
- [4] J. E. Gonzalez, R. S. Xin, A. Dave, D. Crankshaw, M. J. Franklin, and I. Stoica. GraphX: Graph processing in a distributed dataflow framework. In *OSDI*, 2014.
- [5] J. Hu, X. Shen, Y. Shao, C. Bystroff, and M. J. Zaki. Mining protein contact maps. In *BIOKDD*, 2002.
- [6] J. Huan, W. Wang, D. Bandyopadhyay, J. Snoeyink, J. Prins, and A. Tropsha. Mining protein family specific residue packing patterns from protein structure graphs. In *Proc. Int. Conf. on Research in Computational Molecular Biology*, 2004.
- [7] M. Jambon, A. Imbert, G. Deléage, and C. Geourjon. A new bioinformatic approach to detect common 3D sites in protein structures. *Proteins: Structure, Function, and Bioinformatics*, 52(2), 2003.
- [8] J. Leskovec, L. A. Adamic, and B. A. Huberman. The dynamics of viral marketing. *ACM Transactions on the Web (TWEB)*, 1(1):5, 2007.
- [9] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, and J. M. Hellerstein. Distributed GraphLab: A framework for machine learning in the cloud. *PVLDB*, 5(8), 2012.
- [10] G. Malewicz, M. H. Austern, A. J. C. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. In *SIGMOD*, 2010.
- [11] P. Meysman, C. Zhou, B. Cule, B. Goethals, and K. Laukens. Mining the entire protein databank for frequent spatially cohesive amino acid patterns. *BioData Mining*, 8, 2015.
- [12] G. Palla, I. Derényi, I. Farkas, and T. Vicsek. Uncovering the overlapping community structure of complex networks in nature and society. *Nature*, 435(7043), 2005.
- [13] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, and T. E. Ferrin. UCSF Chimera – A visualization system for exploratory research and analysis. *J. Comp. Chemistry*, 25(13), 2004.
- [14] A. Quamar, A. Deshpande, and J. J. Lin. NScale: Neighborhood-centric large-scale graph analytics in the cloud. *Vldb J.*, 25(2), 2016.
- [15] C. H. C. Teixeira, A. J. Fonseca, M. Serafini, G. Siganos, M. J. Zaki, and A. Aboulmaga. Arabesque: A system for distributed graph mining. In *SOSP*, 2015.
- [16] M. Zaharia et al. Apache Spark: A unified engine for big data processing. *Comm. ACM*, 59(11), 2016.