

# ADAM<sub>pro</sub>: Database Support for Big Multimedia Retrieval

Ivan Giangreco · Heiko Schuldt

**Abstract** For supporting retrieval tasks within large multimedia collections, not only the sheer size of data but also the complexity of data and their associated metadata pose a challenge. Applications that have to deal with big multimedia collections need to manage the sheer volume of data and to effectively and efficiently search within these data. When providing similarity search, a multimedia retrieval system has to consider the actual multimedia content, the corresponding structured metadata (e.g., content author, creation date, etc.) and – for providing similarity queries – the extracted low-level features stored as densely populated high-dimensional feature vectors. In this paper, we present ADAM<sub>pro</sub>, a combined database and information retrieval system that is particularly tailored to big multimedia collections. ADAM<sub>pro</sub> follows a modular architecture for storing structured metadata, as well as the extracted feature vectors and it provides various index structures, i.e., Locality-Sensitive Hashing, Spectral Hashing, and the VA-File, for a fast retrieval in the context of a similarity search. Since similarity queries are often long-running queries, ADAM<sub>pro</sub> supports progressive queries that provide the user with streaming result lists by returning (possibly imprecise) results as soon as they become available. We provide the results of an evaluation of ADAM<sub>pro</sub> on the basis of several collection sizes up to 50 million entries and feature vectors with different numbers of dimensions.

**Keywords** Databases, Multimedia retrieval systems, Big data

---

Ivan Giangreco · Heiko Schuldt  
Databases and Information Systems Group,  
University of Basel,  
Basel, Switzerland  
E-mail: [ivan.giangreco@unibas.ch](mailto:ivan.giangreco@unibas.ch)

## 1 Introduction

Multimedia is *Big Data*, both in terms of their volume and their heterogeneity. Many applications that have to deal with such big multimedia collections need support for managing the sheer volume of data and for effectively and efficiently searching within these data – based on annotated (structured) metadata and/or based on intrinsic features of the multimedia objects. Consider, for instance, the following applications: a TV station is looking for videos and scenes with some specific visual content to enrich a news report; a medical researcher is looking for all mammograms showing a certain visual characteristic which might indicate a case of breast cancer; or a user is looking for a piece of music that she remembers without, however, knowing its name.

Obviously, all these applications have in common that large multimedia collections need to be searched on the basis of their content. In the last years, this has successfully spurred research in the field of Machine Learning in order to detect and possibly learn features for characterising the content of multimedia objects and thus to serve as basis for retrieving results when comparing the object features to a query (e.g., [9, 10]). However, these feature extractors only form one side of the coin. The other side of the coin is formed by the organisation and storage of feature data (in general: any form of metadata on multimedia objects), and the support for various types of queries using these metadata.

In this paper, we address multimedia queries that combine Boolean retrieval based on structured metadata (e.g., content author, creation date, etc.) with the vector space retrieval model to support similarity queries on the objects' content. By this, we consider both the information retrieval and the database approach: In-

formation retrieval systems traditionally focus on high dimensional feature spaces and support nearest neighbour queries (similarity search). Database systems, on the other hand, traditionally only come with very limited support for similarity searches and partial matches, but are very good at organising data and retrieving exact matches. For multimedia content and its associated metadata, the ‘one size fits all’ approach generally does not work. As a consequence, multimedia retrieval systems often store feature data in file-based structures. From a database perspective, this has various drawbacks: On the one hand, the absence of *physical and logical data independence* makes the organisation of data a difficult and tedious task that is prone to errors. On the other hand, by letting the retrieval system take care of storing data, the principle of *separation of concerns* is violated in that application logic (i.e., on *what* feature data to store) and storage logic (i.e., on *how* to organise and store feature data) are not well separated. The missing separation of concerns leads to the re-engineering of components necessary for the sole task of managing and organising data.

In this paper, we introduce  $ADAM_{pro}$ , a database and information retrieval system that jointly supports Boolean and vector space retrieval and that is particularly tailored to very large multimedia collections.  $ADAM_{pro}$  is an extension of the ADAM system [6, 7] and focuses on storage support for big multimedia data. It follows a modular architecture: based on the nature of the data to be managed (structured or unstructured), individual modules can be replaced to increase the overall query efficiency and to reduce response time. In addition,  $ADAM_{pro}$  jointly supports various index structures. By combining index structures that quickly produce (approximate) results with index structures that take longer to produce (correct) results, so-called *progressive* queries can be supported. In the context of the retrieval process, results are presented to a user in a streaming fashion as soon as they become available. The query results will be continuously updated as more (and also more precise) results will be available. While this is not necessary for small collections, *progressive* queries in  $ADAM_{pro}$  allow to provide fast query results especially for big multimedia collections.

The contribution of the paper is threefold:

1. We present the architecture of the  $ADAM_{pro}$  system and the interplay between its components at query time to be able to efficiently process a query. In this paper, we focus on the vector space retrieval model in which feature data reflecting the objects’ content is represented by means of (high-dimensional) feature vectors and a  $k$  nearest neighbour (kNN) search to find the most similar objects.
2. We introduce the concept of progressive query evaluation which aims at reducing retrieval time.  $ADAM_{pro}$  seamlessly combines various index structures, i.e., Locality-Sensitive Hashing (LSH), Spectral Hashing, and the Vector Approximation File (VA-File), to answer kNN searches efficiently; the former two are used to generate early, albeit pre-mature (approximate) results while the latter guarantees exact results, but at the price of higher retrieval times.
3. We present results of the evaluation of  $ADAM_{pro}$  under varying configurations and particularly address the scalability of  $ADAM_{pro}$ . This includes both the sheer size of collections (no. of objects) as well as the complexity of the feature data (no. of dimensions of the feature space).

The remainder of the paper is organised as follows: in Section 2, we present a sample application. Section 3 discusses details of the system architecture of  $ADAM_{pro}$  and introduces the concept of progressive queries. Section 4 reports on the evaluation of  $ADAM_{pro}$ . Section 5 discusses related work and Section 6 concludes.

## 2 IMOTION: A Sample Application

As an example for the use of  $ADAM_{pro}$ , consider IMOTION (Intelligent Multi-Modal Augmented Video Motion Retrieval System)<sup>1</sup>, a system for large-scale video retrieval applications [14]. The objective of IMOTION is to provide a rich variety of different query paradigms for searching in video collections. This includes traditional *keyword search* on the basis of automatically collected metadata (e.g., content author, video length, etc.) or manually added tags describing the content of a video or a shot. In addition, IMOTION also supports a large variety of similarity search-based queries, e.g., *Query-by-Sketch (QbS)*, *Query-by-Example (QbE)*, querying by motion and querying by audio. For the retrieval, this means that users can specify either an existing image or a video snippet as a query, or provide the system with a hand-drawn sketch of the most relevant object(s) (e.g., using a color sketch as depicted in Figure 1), or draw motion or record an audio snippet to search for. All the query options can be seamlessly combined – either in a pipelined fashion (e.g., start with a keyword query and use one of the results as query object for QbE), or combined in a single query by superimposing for instance a sketch and a query image (by adding, via a sketch, an object that is not visible on the query image).

<sup>1</sup> <http://www.imotion-project.eu>

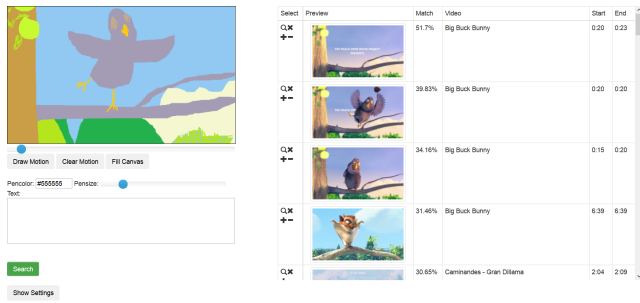


Fig. 1: Sketch-based video retrieval in IMOTION (from [14]).

For supporting retrieval, IMOTION makes use of a large number of feature extractors [14] which produce comparatively high-dimensional dense feature vectors (up to over several hundred dimensions) that are used at query time for comparison to the query object. Given the large size of the video collections the IMOTION query engine is supposed to handle this exceeds the capabilities of current storage systems that are not adapted to the use case at hand (as will be shown in Section 4). In the IMOTION system, ADAM<sub>pro</sub> takes over the task of storing and organising both the structured metadata and the extracted feature vectors, and supports the retrieval logic for retrieving exact or approximate matches.

### 3 Architecture

In retrieval systems such as the IMOTION system, two phases can generally be distinguished: In the *off-line* phase, the retrieval engine extracts features from the given multimedia objects. The feature extraction phase serves two purposes: First, it allows the adaption of a (multimedia) document to the retrieval framework used; second, it reduces the search complexity, since it avoids the full inspection of objects, but only considers the extracted features for comparison. In the *on-line* phase, i.e., the time-sensitive query phase, on the other hand, the extracted features are compared to the query object, and ranked by similarity.

Figure 2 shows ADAM<sub>pro</sub> in the described setting: in the off-line phase, ADAM<sub>pro</sub> takes over the task of storing the feature vectors as they are inserted into the system. In the on-line phase, ADAM<sub>pro</sub> is responsible for a fast response time given a query.

To this end, ADAM<sub>pro</sub> combines various storage subsystems depending on the data and the queries at hand; the combination of various systems shows the advantage of each system over the others (and over a monolithic system) in one specific phase of the retrieval:

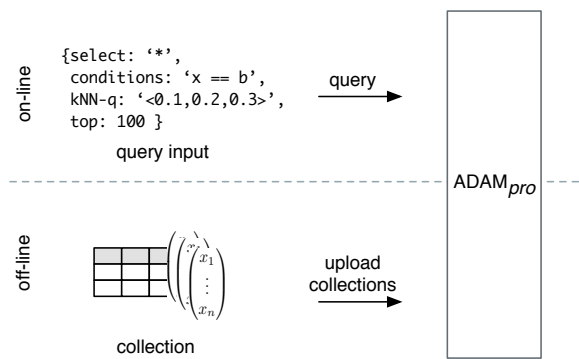


Fig. 2: On-line and off-line phase of a retrieval system.

- for structured metadata, ADAM<sub>pro</sub> uses a relational DBMS, as it allows querying for all attributes in a very elegant manner with good performance; these data are used for a pre-filtering,
- the index structures for a  $k$  nearest neighbour (kNN) search are stored in a file-based format, as they can be well distributed over multiple workers; these data are used to further filter the tuples to retrieve,
- for the feature vector data, we use a key-value store which allows to retrieve the full feature vectors (using the keys filtered by the index) quickly and efficiently for further computation; these data are used for the full distance computation.

The overall architecture of ADAM<sub>pro</sub> is depicted in Figure 3: On the left hand side, the orchestrator is depicted, which takes care of incoming requests (such as insert operations or queries) by calling the corresponding components. The *metadata storage component* is responsible for storing and retrieving structured metadata in a relational database. The *index storage component* builds the index structures and stores them in index files, separated from the actual content and metadata. Finally, the *feature storage component* stores the full feature vectors. At query time, the metadata storage component will filter results based on Boolean predicates (i.e., on the metadata, for instance `date = 15/03/2015`). The results will then be processed by the index storage component that further prunes the result list, by using the index structures available for retrieving the  $k$  nearest neighbours and performing a similarity search on the basis of the given query vector. Finally, the remaining elements are collected from the feature storage component that retrieves the full feature vectors and performs the exact computation of distances.

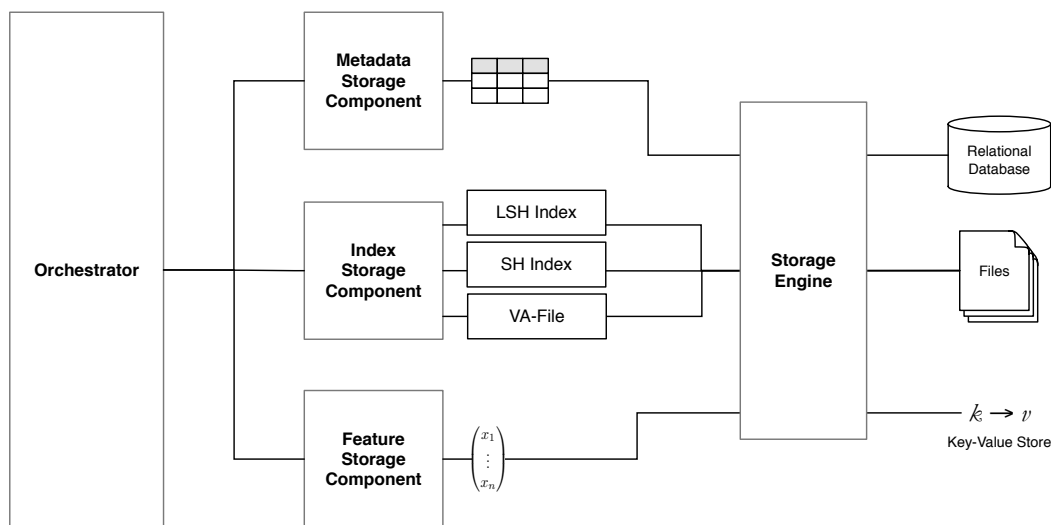


Fig. 3: Architecture of the  $ADAM_{pro}$  system.

### 3.1 Schema Definition and Data Insertion

In the data definition step, the user specifies the logical schema for the data. Given the various data types,  $ADAM_{pro}$  distributes the physical data to different systems:  $ADAM_{pro}$  stores feature data in a key-value store and indexes the data using various index structures built into the system. The metadata, on the other hand, is stored in a relational DBMS. In this phase,  $ADAM_{pro}$  takes care of creating the appropriate schema in the various systems involved.

### 3.2 Retrieval

Consider as an example the following query: A user is looking for the top 100 images that were taken on the 15<sup>th</sup> of March 2015 and that are similar to the given query image. This query uses a Boolean predicate (e.g.,  $date = 15/03/2015$ ) and involves a similarity query (i.e., all objects similar to the query image or all feature vectors similar to the query vector, respectively). The results should obviously be ordered by similarity and pruned at 100 results. As shown in this example, a query in this setting may ask for all objects similar to the given query object while at the same time fulfilling all Boolean predicates.

$ADAM_{pro}$  supports both Boolean retrieval and kNN similarity search. Retrieval based on Boolean predicates can be applied on all structured fields (i.e., the structured metadata). The similarity search is applied on the feature vector. To increase the retrieval performance, various index structures for high-dimensional data are used (Section 3.3). To increase the retrieval performance,

these index structures prune results from the final result set (as they take into account the limiting factor  $k$  of a kNN search), rather than only ordering the results. Therefore, it is crucial that the Boolean retrieval is always performed before the  $k$ NN search, as otherwise results would get lost. Furthermore, by first performing the (fast) Boolean retrieval, the similarity search can avoid to consider results that do not adhere to the Boolean predicate and by that improve the system's performance.

In Figure 4, we display a high-level execution plan for a query in the  $ADAM_{pro}$  system: (1) If Boolean predicates are available in the query, the query is first sent to the relational database that returns a result set fulfilling the Boolean predicates (TID list). (2) Using the built-in index structures, the query vector and the TID result set from step 1, in the second step, the nearest result candidates are retrieved. This result set contains possibly more than  $k$  elements and is not yet sorted. Furthermore, the result elements do not yet contain the exact distance values. (3) In the last step, using the TID set of step 2, the full feature vectors are retrieved and the exact distances are computed. The result is then returned.

### 3.3 Index Structures

To support efficient  $k$  nearest neighbour retrieval, the  $ADAM_{pro}$  system implements Locality-Sensitive Hashing (LSH) [8], Spectral Hashing (SH) [18] and the Vector Approximation-File (VA-File) [17]. We detail the implemented index structures in the following.

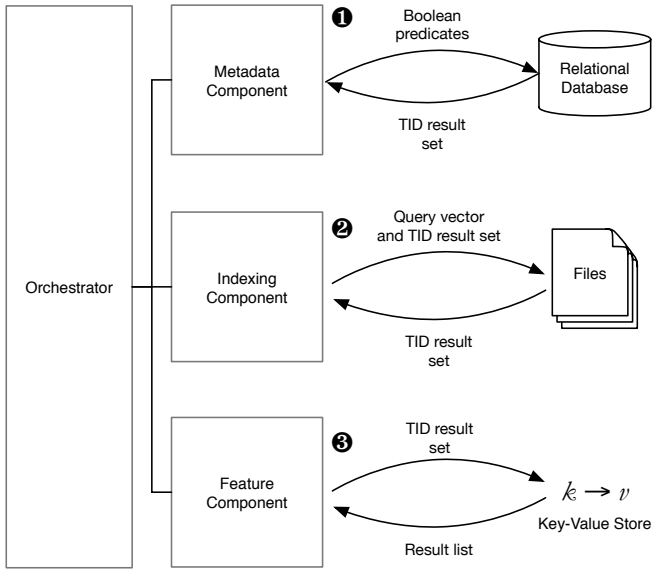


Fig. 4: General query plan in the ADAM<sub>pro</sub> system.

### 3.3.1 Locality-Sensitive Hashing

The main idea of Locality-Sensitive Hashing [8] is to hash each feature vector  $f$  using several hash functions and use the hashes for finding near neighbours. Formally, using a hash function, a  $d$  dimensional vector is mapped onto the space of integers ( $h_a(f) = \mathbb{R}^d \mapsto \mathbb{N}$ ). A family  $\mathcal{H}$  of hash functions is called locality-sensitive or more specifically  $(R, cR, P_1, P_2)$ -sensitive, if for any two feature vectors  $p, q \in \mathbb{R}$ :

- if  $\|p - q\| \leq R$  then  $P_{\mathcal{H}}[h(q) = h(p)] \geq P_1$
- if  $\|p - q\| \geq cR$  then  $P_{\mathcal{H}}[h(q) = h(p)] \leq P_2$

with  $P_1 > P_2$ . Intuitively, if two vectors are close, the probability that they collide in their hash should be high; vice-versa if two vectors are far apart, the probability that they collide in their hash should be small.

ADAM<sub>pro</sub> currently supports the Minkowski distance; therefore, for LSH, we use the family of hash functions proposed in [4] for  $l_p$  norms, based on  $p$ -stable distributions. For the hash function, we pick a random projection  $a \in \mathbb{R}^d$  with entries from a  $p$ -stable distribution, chop the line into equi-width segments ( $w$ ) and shift by a random value  $b \in [0, w)$ . Formally, the hashing function is given as  $h_{a,b}(v) = \lfloor \frac{av+b}{w} \rfloor$ .

While LSH can be very efficient, it has to be noted that the hashing approach yields both false positive (irrelevant items in the result list) and false negatives (missing relevant items).

### 3.3.2 Spectral Hashing

Spectral Hashing [18] belongs also to the family of hash-based indexing methods, however, falls into the category of ‘learning to hash’, i.e., the hash functions are generated based on the data at hand. The idea of Spectral Hashing [18] is to find a hash function such that similar items are mapped to similar hash codes, i.e., small distances in the feature space should result in small Hamming distances between the codes. The embedding is done using the eigenfunctions computed along the PCA directions. In essence, the algorithm will

1. find the principal components of the data using a principal component analysis (PCA)
2. compute the Laplacian eigenfunctions with the smallest eigenvalues along every PCA direction
3. threshold the eigenfunctions at zero to obtain the binary codes

As with LSH, Spectral Hashing may yield both false positive and false negatives.

### 3.3.3 Vector Approximation (VA) File

Behind Vector Approximation-File (VA-File) [17] lies the idea to build an index that yields exact results and at the same time is very performant. The authors argue that for increasing dimensionality, any tree-based index structure degenerates to a sequential scan. Therefore, the authors suggest to compress the feature vectors in a quantization step to a short signature that can quickly be scan in a sequential manner and which allows to early prune the result list. This is achieved, by using the signatures to compute at query time upper and lower bounds to the distance and by that early exclude items that are too distant. The upper and a lower bound of the distance can be calculated with very few simple calculations and the computation is therefore computationally less complex than a full distant computation on the vectors. Furthermore, by only reading the signatures, less page accesses and therefore I/O accesses are necessary, which would largely increase the retrieval time. While VA-File may degenerate to a sequential scan, it will always return all true positives.

To create a VA signature, a fixed-length bit string for each data point is generated. For that purpose, the data space is divided in  $2^{b_{tot}}$  cells, where  $b_{tot}$  denotes the total length of the bit signature, and the cells are enumerated in a binary way. Each dimension  $d$  receives  $b_d$  bits that are finally concatenated to create the full bit mask [17].

### 3.4 Progressive Query Results

Queries in multimedia retrieval systems are often comparably long-running queries, as they cannot profit from tree-based index structures that significantly reduce the search time complexity. The reason for this is that feature vectors, on which the queries are performed, do not have an absolute ordering, but the ordering results only based on the given query. [17] argues that with increasing dimensionality, tree-based index structures degenerate to a sequential scan of the data. Therefore, predicting the query time for indexes for high-dimensional data is a very difficult task: Traditional database systems consider the number of index and data pages, the height of the index tree, the length of TID lists in the leaf nodes of the tree, etc. For the indexes used in *ADAM<sub>pro</sub>*, these parameters are not appropriate predictors for the retrieval time. Consider, for example, a VA-File index: The algorithm behind the VA-File will scan all database elements. However, by decreasing the size of the signature to be scanned, the number of operations to be performed for computing the final distance and consequently the number of data pages to be loaded, can be significantly decreased. Nevertheless, using the number of index pages containing the signature for estimating the retrieval time, will actually not correctly estimate the retrieval time, as in the most degenerate case, the VA-File has to consider all vectors stored on the data pages. In particular, as this is not an inherent property of the data only, but the data in combination with the query at hand, it is a hard task to achieve to predict the query time. Finally, these predictions do not consider whether the indexes will return exact or only approximate results.

For this reason, *ADAM<sub>pro</sub>* supports so-called *progressive querying* that results in streamed result lists. Using this approach, *ADAM<sub>pro</sub>* runs at the same time all physical plans to execute the same logical plan and returns the (possibly approximate) answer(s) as soon as they become available to the user. Starting the same query using different query plans at the same time may decrease the efficiency of the entire system, but it allows to trade computation (which is not a bottleneck in modern environments) with query response time. For a user who is waiting a short time for her results, the results from the database may only be very approximate; if the user, on the other hand, waits a bit longer, she may get an answer that is for sure precise. In any case, she will get the first possible answer as soon as it is available.

Consider, for clarification, the following example: A query for the  $k$  nearest neighbours is started on all available indexes: the Spectral Hashing index may return

first, due to its low query complexity, however the result may contain false positives or lack true positives. Only in the next step, the exact results from the VA-File index may arrive. On the other hand, in very degenerate cases, a sequential scan may return even before the VA-File index, ensuring that the user gets the results as fast as possible.

### 3.5 Distribution

In a distributed setting, there is no obvious partitioning scheme that can be generally applied to the feature data and that allows to prune nodes at query time from the retrieval without possibly losing result elements (this is also a consequence of the fact that tree-based methods do not work well for feature vector data). While, for instance, text retrieval systems can choose to query only specific nodes that are responsible for a certain keyword appearing in the query, this approach is not easily adaptable to kNN retrieval in the multimedia context. As a consequence, for a query, all data items have to be considered, i.e., a query has to be processed by all nodes and the sub-results of each node have finally to be merged. The distribution of the VA-File, for instance, has already been discussed in [16]. We have implemented the same ideas for all index structures and we perform the retrieval in *ADAM<sub>pro</sub>* in a map/reduce fashion. At query time, the index structures are partitioned to be processed by multiple workers (possibly residing on the same node). Each node takes care of filtering out the  $k$  nearest neighbours and sends the partial result list to the master node that merges the result lists to a global result list of nearest neighbours.

## 4 Performance Evaluation

### 4.1 Implementation

*ADAM<sub>pro</sub>* is implemented in Java/Scala using Apache Spark 1.5<sup>2</sup>. For storing the metadata, we use PostgreSQL<sup>3</sup>, the Apache Parquet columnar file format on the Hadoop file system (HDFS) for the index, and Apache Cassandra 2.1<sup>4</sup> as key-value store for storing the feature vectors. *ADAM<sub>pro</sub>* can be accessed via a REST interface. LSH is implemented based on E2LSH<sup>5</sup> and Spectral Hashing is based on the Matlab code provided by the authors<sup>6</sup>.

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

<sup>3</sup> <http://www.postgresql.org/>

<sup>4</sup> <http://cassandra.apache.org/>

<sup>5</sup> <http://www.mit.edu/~andoni/LSH/>

<sup>6</sup> <http://www.cs.huji.ac.il/~yweiss/SpectralHashing/>

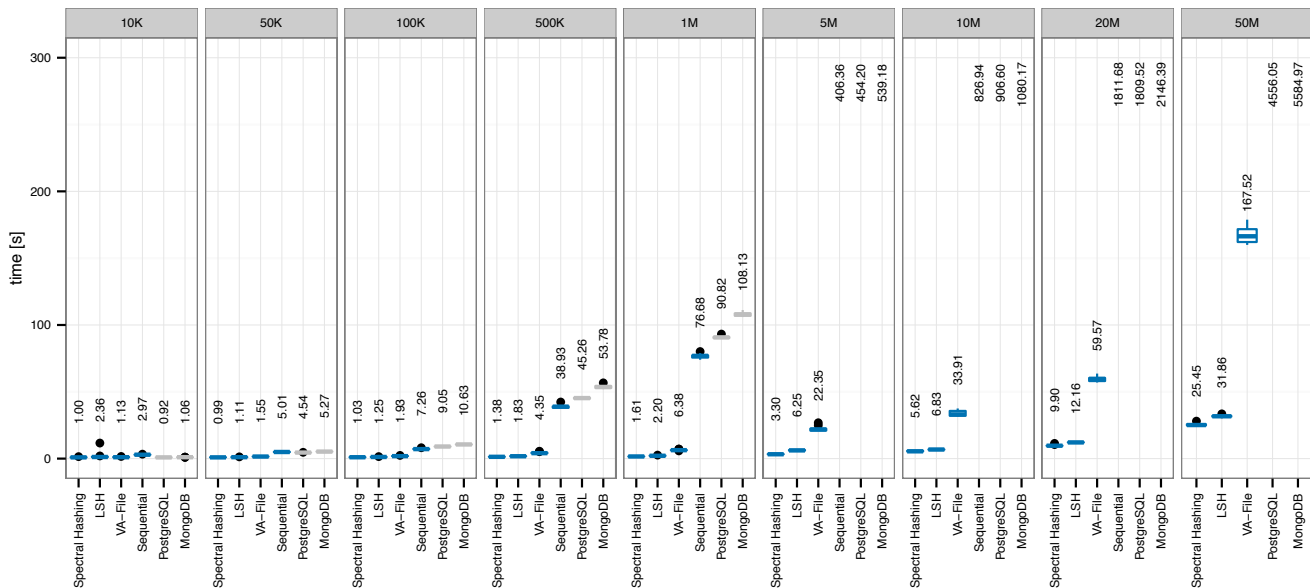


Fig. 5: Performance evaluation showing the retrieval time at varying collection sizes for the various methods implemented in ADAM<sub>pro</sub> with dimensionality of the feature vectors at 100. As a baseline, both PostgreSQL and MongoDB have been added to the plot. The experimental runs are summarised by the mean time for each method. The plot is cut at 300 seconds; for all values above 300 seconds only the mean time is displayed.

## 4.2 Experimental Setup

We have evaluated ADAM<sub>pro</sub> using artificially generated data at various numbers of dimensions and collection sizes. The evaluation setup involves the following parameters:

- collection size: 10K, 50K, 100K, 500K, 1M, 5M, 10M, 20M, 50M
- dimensions: 10, 50, 100, 200, 500
- execution plans: sequential scan, LSH scan, Spectral Hashing scan, VA-File scan

The vectors added to the collection and the query vectors are composed of uniformly distributed float values  $\in (0, 1)$ . To avoid anomalies in the results, we have run each experimental setting five times and average over the different runs for the same parameter setting. We run ADAM<sub>pro</sub> (and the systems to compare it to) on Microsoft Windows Azure using one DS13 instance running Ubuntu 14.04 with 8 cores and 56 GB memory. For all experiments, the parameters for generating the index (number of bits for signature, etc.) have been set beforehand to a general value independent of the given data.

As a baseline, we use PostgreSQL 9.4<sup>7</sup> and MongoDB 3.0<sup>8</sup>. In PostgreSQL, we use a custom function

<sup>7</sup> <http://www.postgresql.org/>

<sup>8</sup> <http://www.mongodb.com>

to compute the distance between two float arrays and use `ORDER BY` and `LIMIT` to model a  $k$  nearest neighbor search. In MongoDB, on the other hand, we make use of a server-based script that computes in a map/reduce fashion the distance between a query array and the vectors stored in the collection: the map function emits the distance between the query vector and the vector from the collection, whereas the reduce function sorts and slices the results. Both baselines do not make use of any indexing structure that could improve the kNN retrieval.

## 4.3 Evaluation of single execution plans

We first consider the results of the evaluation for every single execution plan: Figure 5 shows a box plot that compares the retrieval time at varying collection sizes with a fixed dimension of 100 for the feature vectors. Note that the plot has been cut at 300 seconds and for all values above this threshold, only the mean time is displayed. It can be seen that when increasing the number of items, Locality Sensitive Hashing and Spectral Hashing retrieve the results faster than the other scanning methods (i.e., than the VA-File scan and the sequential scan). Furthermore, it can be seen that in our experiments the VA-File always performs better than the sequential scan. This means that given the data and the queries used in the evaluation, we never

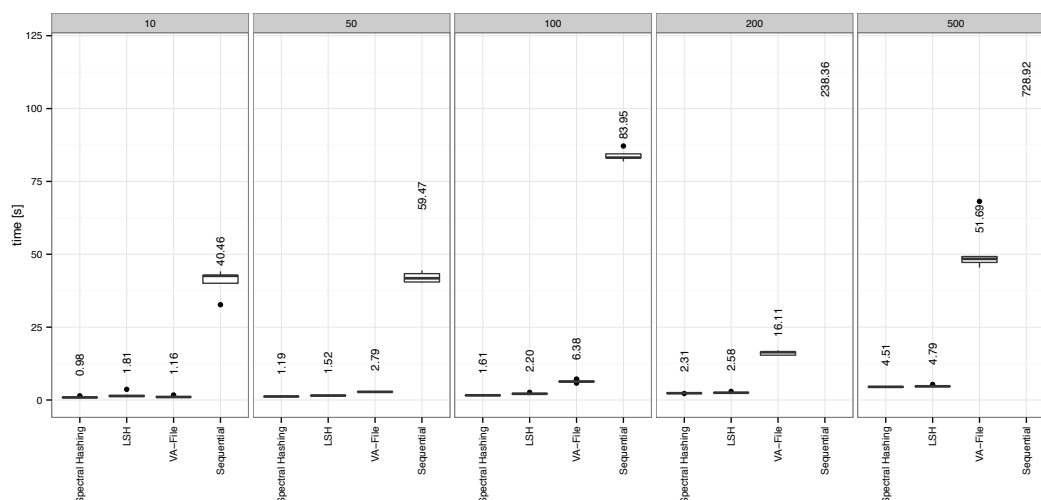


Fig. 6: Performance evaluation showing the retrieval time at varying number of dimensions of the feature vectors stored in  $ADAM_{pro}$  using the various methods implemented at a collection size of 1M elements. The plot is again cut at 120 seconds and for all values above this threshold only the mean time is displayed.

get into the degenerate case in which a sequential scan is truly necessary.

This behaviour is as expected, as LSH and SH allow for a simple lookup, whereas the VA-File has to scan all signatures; however, in exchange, the VA-File returns precise results. For increasing collection size, in particular for collections that contain more than 50K elements, our system performs in any case better than performing a sequential scan in MongoDB or PostgreSQL, the baseline to our evaluation.

As can be seen from Figure 5, the retrieval time is obviously dependent on the collection size. Moreover, as it can be seen from Figure 6, it also depends on the dimensionality. Figure 6 shows the retrieval time for increasing numbers of dimensions of the feature vectors (at a fixed collection size of 1M elements). As can be seen, the retrieval time increases with the dimensionality of the feature vector.

#### 4.4 Evaluation of progressive querying

Given these observations, we evaluate the behaviour of progressive querying. In particular, we show the results exemplified at collection sizes of 100K (Figure 7a) and 1M (Figure 7b) elements, respectively. In Figure 7, we show at which time after starting the query  $ADAM_{pro}$  presents its results to the user. In the original  $ADAM_{pro}$  implementation, the query execution is normally cancelled as soon as exact results (i.e., results from the sequential scan or from VA-File) are retrieved; for this evaluation, to be able to show the times of the various scans, we have adjusted the implementation not to

abort the execution (nevertheless, we have marked the time at which  $ADAM_{pro}$  would stop the further execution in its normal setup with a red line). This means that even after finding the final and exact results, we continue to execute the query and measure the query time for the remaining execution plans.

In Figure 7, it can be seen that predicting the query time for the various indexing structures is a difficult task: it cannot be clearly said which index structure performs better under which conditions. Particularly for small dimensionalities, it is not obvious whether LSH, Spectral Hashing or VA-File will return first its results. With increasing dimensionality, both hashing based methods perform clearly better than VA-File and the sequential scan. Particularly for increasing dimensionality (and collection size), the progressive querying approach becomes more and more important: If, in a collection of 1M elements and feature vectors with a dimensionality of 200 elements, a user realises that the results received after about 20 seconds are good enough or not worthy to further consider, she may as well abort the further execution of the query; on the other hand, she may wait to get precise results after about 60 seconds. As expected, in all runs, the sequential scan is the last execution plan that returns its query results. Since our query handler is aware of the fact that VA-File returns precise results, it would abort the further execution after VA-File has answered.



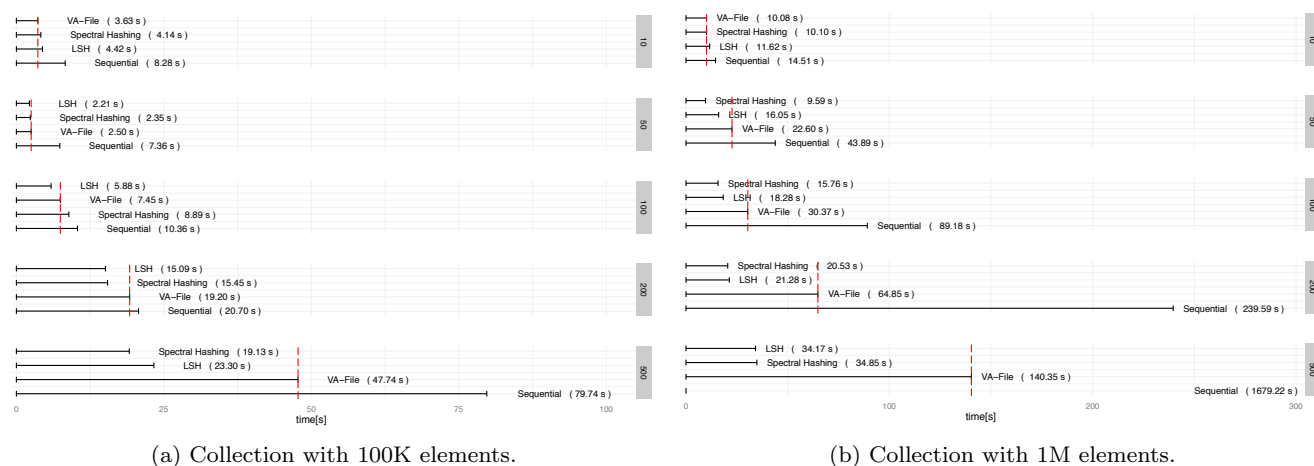


Fig. 7: Timeline displaying the mean time at which the results using the various scans become available when using progressive querying for various dimensions at a collection size of 1M. The red line denotes when ADAM<sub>pro</sub> would stop the further execution of the retrieval, as precise results have been found.

## 5 Related Work

One of the earliest works attempting to integrate an information retrieval system and a database management system can be found in [15]. In this early work from 1980, the authors note especially the lack of support for information retrieval queries in the query language and in the internal indexing techniques of a database management system (DBMS). The authors motivate the integration of the two systems into a database management and information retrieval system (DBMIRS) by the need for queries supporting both formatted (structured) and unformatted (unstructured) data retrieval.

An early integration of multimedia data in databases can be found in the IBM project Garlic [3]. In here, the authors integrate multiple federated databases (some of those from the IBM project QBIC) into one distributed system for multimedia data. The system is based on an object-oriented database model. For query formulation, the authors extend the object-oriented query language (OQL).

Similarly, [12] introduces Chabot, an information retrieval system using an underlying PostgreSQL database for storing extracted features. The authors implement complex types, user-defined indices and user-defined functions into the database to support information retrieval data types and queries. Chabot not only supports text-based queries, but also allows to improve results by providing content-based queries (e.g., “images with some orange in it”). Both Garlic and Chabot support Boolean predicates rather than similarity-based queries.

[5] presents Mirror, a database supporting content-based multimedia retrieval. The authors describe the

engineering factors for creating a distributed multimedia IR-DBMS that uses Moad, a new relational algebraic framework based on the non-first normal form (NF<sup>2</sup>). Mirror is implemented on top of the object-relational DBMS Monet [2].

Further, DISIMA DBMS [13] is a DBMS that allows to store syntactic features, i.e., color, shape, texture, and semantic features, i.e., real world objects or concepts, in an object-oriented data model. The system supports content-based searches and searches on image semantics. The authors implement an extended version of OQL for multimedia objects (MOQL) and VisualMOQL, a visual counterpart to MOQL. To increase the performance of the system, the authors use three-dimensional extendible hashing (3DEH) that allows to pre-filter images based, for instance, on the average color.

[1] introduces a system that combines low-level (syntactic) features with semantic features in a commercial object-relational database. The database is extended by several User Defined Types following the MPEG-7 standard descriptors, and operations implemented in PL/SQL, e.g., to evaluate similarity measures.

In [11], the authors make use of the map/reduce paradigm for querying large sets of image data in a cloud environment. The authors use an indexing method called extended Cluster Pruning (eCP) for indexing the feature data and port it to the map/reduce paradigm on the Hadoop platform.

## 6 Conclusion

In this paper, we have presented  $ADAM_{pro}$ , a modular system that manages large multimedia collections.  $ADAM_{pro}$  flexibly supports various storage systems and indexing structures (LSH, Spectral Hashing, VA-File) to increase the overall query efficiency and reduce response time. Furthermore, with  $ADAM_{pro}$ , we have presented the concept of progressive queries that embraces the idea of returning a result stream to the user: depending on how long the user waits, the system will refine the results and produce correct instead of only approximate results.

In our future work, we plan to consider different distribution scenarios for  $ADAM_{pro}$  and we plan to further increase the collection size and dimensionality.

**Acknowledgements** This work was partly supported by the Swiss National Science Foundation in the context of the CHIST-ERA project IMOTION, contract no. 20CH21.151571. Furthermore, the authors would like to thank the reviewers for their time in reviewing our manuscript and for their helpful comments to improve this paper.

## References

1. Carlos E. Alvez and Aldo R. Vecchietti. Combining Semantic and Content Based Image Retrieval in ORDBMS. In Rossitza Setchi, Ivan Jordanov, Robert J. Howlett, and Lakhmi C. Jain, editors, *Knowledge-Based and Intelligent Information and Engineering Systems, 14th International Conference, KES 2010, Cardiff, UK, September 8-10, 2010, Proceedings, Part II*, volume 6277 of *Lecture Notes in Computer Science*, pages 43–55. Springer Berlin Heidelberg, 2010.
2. Peter A Boncz and Martin L Kersten. Monet. an impressionist sketch of an advanced database system. In *In Proc. IEEE BIWIT workshop*, San Sebastian, Spain, 1994.
3. Michael J. Carey, Laura M. Haas, Peter M. Schwarz, Manish Arya, William F. Cody, Ronald Fagin, Myron Flickner, Allen Luniewski, Wayne Niblack, Dragutin Petkovic, Joachim Thomas II, John H. Williams, and Edward L. Wimmers. Towards Heterogeneous Multimedia Information Systems: The Garlic Approach. In *RIDEDOM 1995: International Workshop on Research Issues in Data Engineering - Distributed Object Management*, pages 124–131, Taipei, Taiwan, 1995.
4. Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S. Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the 20th ACM Symposium on Computational Geometry, Brooklyn, New York, USA, June 8-11, 2004*, SCG '04, pages 253–262, 2004.
5. Arjen P. de Vries and H. M. Blanken. Database technology and the management of multimedia data in the mirror project. In *Proc. SPIE*, volume 3527, pages 443–453. International Society for Optics and Photonics, 1998.
6. Ivan Giangreco, Ihab Al Kabary, and Heiko Schuldt. ADAM - A database and information retrieval system for big multimedia collections. In *Proceedings of the 2014 IEEE International Congress on Big Data*, pages 406–413, Anchorage, AK, USA, June/July 2014. IEEE.
7. Ivan Giangreco, Ihab Al Kabary, and Heiko Schuldt. ADAM: a system for jointly providing IR and database queries in large-scale multimedia retrieval. In *Proceedings of the 37th International ACM Conference on Research and Development in Information Retrieval (SIGIR'14)*, pages 1257–1258, Gold Coast, Australia, July 2014. ACM.
8. Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the 13th Annual ACM Symposium on the Theory of Computing*, pages 604–613, Dallas, Texas, USA, 1998.
9. Andrej Karpathy and Fei-Fei Li. Deep visual-semantic alignments for generating image descriptions. *CoRR*, abs/1412.2306, 2014.
10. Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
11. Diana Moise, Denis Shestakov, Gylfi Thór Gudmundsson, and Laurent Amsaleg. Indexing and searching 100m images with map-reduce. In *International Conference on Multimedia Retrieval, ICMR'13, Dallas, TX, USA, April 16-19, 2013*, pages 17–24, 2013.
12. Virginia Ogle and Michael Stonebraker. Chabot: Retrieval from a Relational Database of Images. *Computer*, 28(9):40–48, 1995.
13. Vincent Oria, M. Tamer Özsu, and Paul Iglinski. Querying Images in the DISIMA DBMS. In *MIS 2001: Workshop on Multimedia Information Systems*, pages 89–98, Capri, Italy, 2001.
14. Luca Rossetto, Ivan Giangreco, Heiko Schuldt, Stéphane Dupont, Omar Seddati, Metin Sezgin, and Yusuf Sahillioğlu. IMOTION – a content-based video retrieval engine. In *Proceedings of the 21st International Conference on MultiMedia Modeling (MMM'15), Part II*, pages 255–260, Sydney, Australia, January 2015. Springer.
15. Hans-Jörg Schek. Methods for the Administration of Textual Data in Database Systems. In Robert N. Oddy, Stephen E. Robertson, C. J. van Rijsbergen, and P. W. Williams, editors, *SIGIR 1980: International Conference on Research and Development in Information Retrieval*, pages 218–235, Cambridge, England, 1980. Butterworth & Co.
16. Roger Weber, Klemens Böhm, and Hans-Jörg Schek. Interactive-time similarity search for large image collections using parallel va-files. In *Research and Advanced Technology for Digital Libraries, 4th European Conference, ECDL 2000, Lisbon, Portugal, September 18-20, 2000, Proceedings*, pages 83–92, 2000.
17. Roger Weber, Hans-Jörg Schek, and Stephen Blott. A Quantitative Analysis and Performance Study for Similarity-Search Methods in High-Dimensional Spaces. In *VLDB 1998: International Conference on Very Large Data Bases*, pages 194–205, New York, USA, 1998.
18. Yair Weiss, Antonio Torralba, and Robert Fergus. Spectral hashing. In *Advances in Neural Information Processing Systems 21, Proceedings of the Twenty-Second Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 8-11, 2008*, pages 1753–1760, 2008.