

Automatic Extraction of Figures from Scientific Publications in High-Energy Physics

Piotr Adam Praczyk (piotr.praczyk@gmail.com) is PhD student at Universidad de Zaragoza, Spain, and research grantholder at the Scientific Information Service of CERN, Geneva, Switzerland

Javier Nogueras-Iso (jnog@unizar.es) is Associate Professor, Computer Science and Systems Engineering Department, Universidad de Zaragoza, Spain

Salvatore Mele (Salvatore.Mele@cern.ch) is the leader of the Open Access section at the Scientific Information Service of CERN, Geneva, Switzerland

Abstract:

Plots and figures play an important role in the process of understanding a scientific publication, providing overviews of large amounts of data or ideas that are difficult to intuitively present using only the text. State of art in digital libraries, serving as gateways to knowledge encoded in scholarly writings, does not take full advantage of the graphical content of documents. Enabling machines to automatically unlock the meaning of scientific illustrations would allow immense improvements in the way scientists work and the knowledge is being processed. In this paper we present a novel solution for the initial problem of processing graphical content, obtaining figures from scholarly publications stored in PDF format. Our method relies on vector properties of documents and as such, does not introduce additional errors, characteristic for methods based on raster image processing. Emphasis has been placed on correctly processing documents in High Energy Physics. The described approach makes distinction between different classes of objects appearing in PDF documents and uses spatial clustering techniques to group objects into larger logical entities. A number of heuristics allow the rejection of incorrect figure candidates and the extraction of different types of metadata.

Keywords:

digital library, information retrieval, scholarly publication, scholarly information retrieval, plot, figure, figure extraction, PDF, High-Energy Physics

1 Introduction

Notwithstanding the technological advances of large-scale digital libraries and novel technologies to package, store and exchange scientific information, scientists communication pattern has changed little in decades if not centuries. The key information of scientific articles is still packaged in a form of text and, for several scientific disciplines, in a form of figures.

New semantic text-mining technologies are unlocking the information in scientific discourse and there exist some remarkable examples of attempts to extract figures from scientific publications (Kataria, On Utilization of Information Extracted From Graph Images in Digital Documents. 2008) but current attempts do not provide sufficient level of generality to deal with figures from High-Energy Physics (HEP) and cannot be applied in a digital library like INSPIRE, which is the main point of our interest. Publications being the

main area of our interest tend to contain highly specific types of figures, which we understand as any type of graphical content illustrating the text and referenced from within of it. In particular they contain high volume of plots which are line-art images illustrating a dependency of a certain quality on a parameter.

The graphical content of scholarly publications allows much more efficient access to the most important results presented in a publication (Hearst, et al. 2007) (Johnston 2011). The human brain perceives the graphical content much faster than reading an equivalent block of text. Presenting figures together with the publication summary, when displaying search results, would allow more accurate assessment of the article content and in turn lead to a better usage of researchers' time. Enabling users to search for figures describing similar quantities or phenomena to a given one could become a very powerful tool for finding publications describing similar results. Combined with additional metadata, it could provide knowledge about evolution of certain measurement or idea over time.

These and many more applications created an incentive to research possible ways integration of figures in INSPIRE. INSPIRE is a digital library for HEP (Holtkamp, et al. 2010), the application field of this work. It provides a large scale digital library service (1 Million records, 50'000 users), which is starting to explore new mechanisms of using figures in articles of the field to index, retrieve and present information (Praczyk, Nogueras-Iso and Dallmeier-Tiessen, et al. 2012) (Praczyk, Nogueras-Iso and Kaplun, et al. 2011). As a first step, direct access to graphical content before accessing the text of a publication can be provided. Secondly, a description of graphics ("blue-band plot", "the yellow shape region") could be used in addition to metadata and/or full-text queries to retrieve a piece of information. Finally, articles could be aggregated in clusters containing the same or similar plots, in a possible alternative automated answer to a standing issue in information management.

The indispensable step to realise this vision is an automated, resilient and high-efficiency extraction of figures from scientific publications. In this paper, we present an approach that we have developed to address this challenge. The focus has been put on developing a general method allowing the extraction of data from documents stored in Portable Document Format (PDF). The results of the algorithm consist of metadata, raster images of a figure but also vector graphics, which allows easier further processing.

The PDF format has been chosen as the input of the algorithm because it is a de facto standard in scientific communication. In the case of HEP, Mathematics and other exact sciences, the majority of publications are prepared using the Latex document formatting system and later compiled into a PDF file. The electronic versions of publications from outstanding scientific journals are also provided in PDF format. The internal structure of PDF files does not always reveal the location of graphics. In some cases images are included as external entities and easily distinguishable from the rest of documents content, but other times they are mixed with the rest of the content. Therefore, in order not to miss any figure, the low-level structure of a PDF must be analysed. The work described in this paper focuses on the area of High Energy Physics. However, with minor variations, described methods could be applicable in the case of a different area of knowledge.

The rest of the paper is organised as follows. Section 2 presents the state of the art in the area of PDF document analysis and figure extraction. Section 3 describes every step of the extraction method in detail. Section 4 presents the results of the evaluation of the presented method on a test-bed of HEP articles. Finally, the paper ends with some conclusions and outlook on future work.

2 Related Work

Over years of development of Digital Libraries and document processing, researchers developed several methods of automatically extracting and processing graphics appearing in PDF documents. Based on properties of the processed content, these methods can be divided into two groups. The attempts of the first category deal with PDF documents in general, not making any assumptions about the content of encoded graphics or document type. The methods from the second group are more specific to figures from scientific publications. Our approach belongs to the second group.

General tools include command line programs like *pettifogs*¹ or web-based applications like *pdftoword*². These solutions are useful in a general case of documents, but all suffer from the same difficulties when processing scientific publications: Graphics that are recognised by such tools have to be marked as graphics inside PDF documents. This is the case with raster graphics and some other internally stored types objects. In the case of scholarly documents, most graphics are constructed internally using PDF primitives and thus cannot be correctly processed by tools from the first group. Moreover, general tools do not have the necessary knowledge to produce metadata describing the extracted content.

With respect to specific tools for scientific publications it must be noted first that important scientific publishers like Springer or Elsevier have created services to allow access to figures present in scientific publications: the improvement of SciVerse Science Direct site for searching images in the case of Elsevier³ (Elsevier 2012); and the SpringerImages service in the case of Springer⁴ (Eichhorn 2011). These services allow searches triggered from a text box, where the user can introduce a description of the required content. It is also possible to browse images by categories such as types of graphics (Image, Table, Line art, Video and so on). The search engines are limited to searches based on figure captions. In this sense, there is little difference between the image search and text search implemented in a typical digital library.

Most of existing works aiming at the retrieval and analysis of figures use the rasterised graphical representation of source documents as its basis. Browuer et al. (Browuer, et al. 2008) (Kataria, Browuer, et al. 2008) describe a method of detecting plots by means of wavelet analysis. In their work they focus on the extraction of data points from identified figures. In particular, they address the challenge of correctly identifying overlapping points of data in plots. This problem would not manifest itself often in the case of vector graphics, which is the scenario proposed in our extraction method. Vector graphics preserve much more information about the documents content than simple values of pixel colours. In particular, vector graphics describe overlapping objects separately. Raster methods are also much more prone to additional errors being introduced during the recognition/extraction phase. The methods described in this paper could be used along with the method of Kataria (Kataria, On Utilization of Information Extracted From Graph Images in Digital Documents. 2008) for the case of documents being output of a digitisation process.

Liu et al. (Liu. Y, et al. 2007) present a page box-cutting algorithm for the extraction of tables from PDF documents. Their approach is not directly applicable but their ideas of geometrical clustering of PDF primitives are similar to the ones proposed in our work. However, our experiments with their

¹<http://sourceforge.net/projects/pdf-images/> (last access: 17.12.2012)

²<http://www.pdftoword.com/> (last access: 17.12.2012)

³<http://www.sciencedirect.com/> (last access: 17.12.2012)

⁴<http://www.springerimages.com/> (last access: 17.12.2012)

implementation and HEP publications have shown that the heuristics used in their work cannot be directly applied to the case of HEP, showing the need for an adapted approach, even in the case of tables.

A different category of work, not directly related to graphics extraction, but being useful when designing algorithms, has been devoted to the analysis of graph usage in scientific publications. The results presented by Cleveland et al. (Cleveland 1984) describe a more general case than publications of High Energy Physics. Even if the data presented in the work came from scientific publications before 1984, included observations, as for example typical sizes of graphs, were useful with respect to general properties of figures and were taken into account when adjusting parameters of the presented algorithm.

Finally, there exist attempts to extract layout information from PDF documents. The knowledge of page layout is useful to distinguish completely independent parts of the content. The approach of layout and content extraction presented by Chao et al. (Chao and Fan 2004) is the closest to the one we propose in this paper. The difference lies in the fact that we are focusing on the extraction of plots and figures from scientific documents, which usually follow stricter conventions. Therefore, we can make more assumptions about their content and extract more precise data. For instance, our method emphasises the role of detected captions and permits them to modify the way in which graphics are treated. We also extract portions of information that are difficult to be extracted using more general methods, such as captions of figures.

3 The Method

PDF files have a complex internal structure allowing to embed various external objects and to include various types of metadata. However, the central part of every PDF file consists of a visual description of the subsequent pages. The imaging model of PDF uses a language based on a subset of the PostScript language. PostScript is a complete programming language containing instructions (called also operators) which allow to render text and images on a virtual canvas. The canvas can correspond to a computer screen or to another, possibly virtual, device used to visualise the file. The subset of PostScript, which was used to describe content of PDFs had been stripped from all the flow control operations (like loops and conditional executions), which makes it much simpler to interpret than the original PostScript. Additionally, the state of the renderer is not preserved between subsequent pages, making their interpretation independent.

In order to avoid many technical details, which are irrelevant in this context, we will consider a PDF document as a sequence of operators (also called the content stream). Every operator can trigger a modification of the graphical state of the PDF interpreter, which might be drawing a graphical primitive, rendering an external attached object, or modifying a position of the graphical pointer⁵ or a transformation matrix⁶. The outcome of an atomic operation encoded in the content stream depends not only on parameters of the operation, but also on the way previous operators modified the state of the interpreter. Such a design makes a PDF file easy to render but not necessarily easy to analyse.

⁵ At every moment of the execution of a PostScript program, the interpreter maintains a number of variables. Some of them encode current positions within the rendering canvas. Such positions are used to locate the subsequent character or to define the starting point of the subsequent graphical primitive.

⁶ Transformation matrices are encoded inside the interpreters' state. If an operator requires arguments indicating coordinates, these matrices are used to translate the provided coordinates to the coordinate system of the canvas.

Figure 1 provides an overview of the proposed extraction method. At the very first stage, the document is preprocessed and operators are extracted (see Section 3.1). Later, graphical⁷ and textual⁸ operators are clustered using different criteria (see Sections 3.4 and 3.5) and the first round of heuristics rejects regions which cannot be considered figures. In the next phase, the clusters of graphical operators are merged with text operators representing fragments of text to be included inside a figure (see Section 3.4). The second round of heuristics detects clusters which are unlikely to be figures. Text areas detected by the means of clustering text operations are searched for possible figure captions (see Section 3.5). Captions are matched with corresponding figure candidates and geometrical properties of captions are used to refine the detected graphics. The last step generates data in a format convenient for further processing (see Section 3.6).

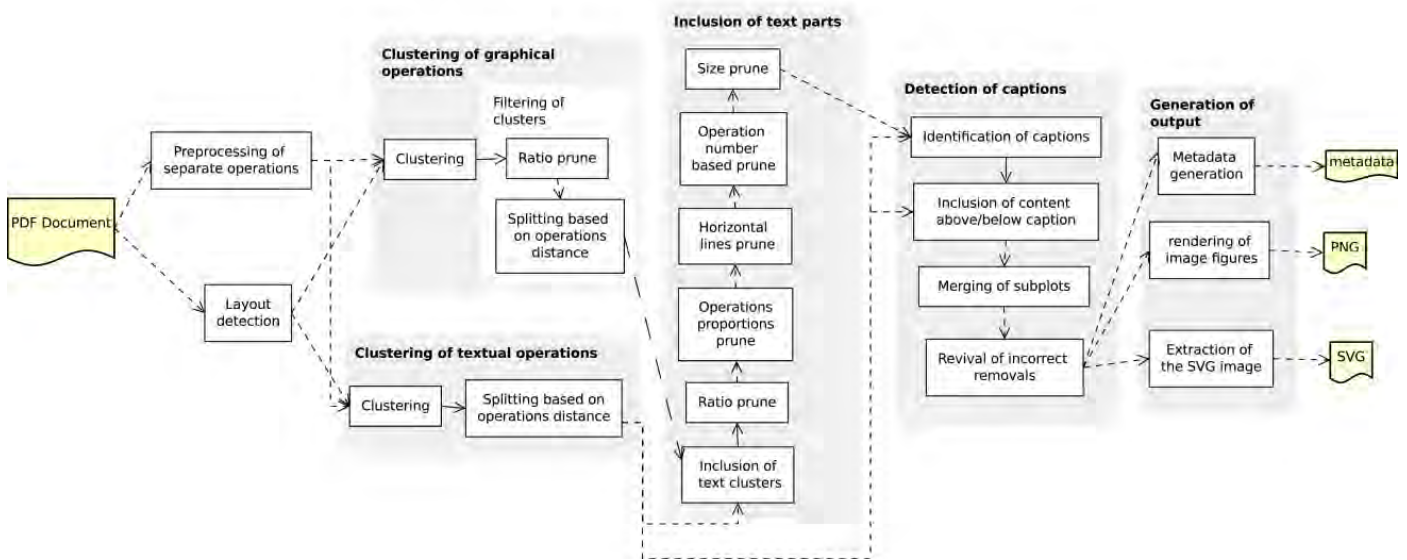


Figure 1: Overview of the figure extraction method

Additionally, it must be noted that another important preprocessing step of the method consists of the layout detection. Section 3.7 discusses an algorithm segmenting pages into layout elements called page divisions. This considerably improves the accuracy of the extraction method because elements from different page divisions can no longer be considered belonging to the same cluster (and subsequently figure). This allows to apply the method separately to different columns of a document page.

3.1 Preprocessing of Operators

The proposed algorithm considers only certain properties of a PDF operator rather than trying to completely understand its effect. Considered properties consist of the operators' type, the region of the page where the operator produces output and, in the case of textual operations, the string representation of the result. For simplicity, we suppress the notion of coordinate system transformation, inherent for the PDF rendering, and describe all operators in a single coordinate system of a virtual 2-dimensional canvas where operations take effect. Transformation operators⁹ are assigned an empty operation region as they do not modify the result directly but affect subsequent operations.

⁷ Graphical operators are those which trigger the rendering of a graphical primitive.

⁸ Textual operations are the PDF instructions which cause the rendering of the text. Textual operations receive the string representation of the desired text and use the current font which is saved in the interpreters' state.

⁹ Operations which do not produce any visible output, but solely modify the interpreters' state.

In our implementation, an existing PDF rendering library has been used to determine boundaries of operators. Rather than trying to understand all possible types of operators, we check the area of the canvas that has been affected by an operation. If the area is empty, we consider the operation to be a transformation. If there exists a non-empty area that has been changed, we check if the operator belongs to a maintained list of textual operators. This list is created based on the PDF specification. If so, the operators argument list is scanned searching for a string and the operation is considered to be textual. An operation that is neither a transformation nor a textual operation is considered to be graphical. It might happen that text is generated using a graphical operator. However, such a situation is unusual. In the case of operators triggering rendering of other operators, which is the case for example when rendering text using type-3 fonts, we consider only the top-level operation.

In most cases, separate operations are not equivalent to logical entities considered by a human reader (such as a paragraph, a figure, a heading and so on...). Graphical operators are usually responsible for displaying lines or curve segments while humans think in terms of illustrations, data lines and so on. Similarly, in the case of text, operators do not have to represent complete or separate words or paragraphs. They usually render parts of words and sometimes parts of more than one word.

The only assumption we make about the relation between operators and logical entities is that a single operator does not trigger rendering of elements from different detected entities (figures, captions). This is usually true because logical entities tend to be separated by a modification of the context (There is a distance between text paragraphs or an empty space between curves).

3.2 Clustering of Graphical Operators

3.2.1 The Clustering Algorithm

The representation of a document as a stream of rectangles allows the calculation of more abstract elements of the document. In our model, every logical entity of the document is equivalent to a set of operators. The set of all operators of the document is divided into disjoint subsets in the process called clustering. Operators are decided to belong to the same cluster based on the position of their boundaries. The criteria for the clustering is based on a simple but important observation: operations forming a logical entity have boundaries lying close to each other. Groups of operations forming different entities are separated by empty spaces.

```

1: Input: OperationSet input_operations {Set of operators of the same type}
2: Output: Map<Rectangle, OperationSet> {Spatial clusters of operators}
3: IntervalTree tx ← IntervalTree()
4: IntervalTree ty ← IntervalTree()
5: Map<Operation, Operation> parent ← Map()
6: for all Operation op ∈ input_operations do
7:   Rectangle boundary ← extendByMargins(op.boundary)
8:   repeat
9:     OperationSet int_opsx ← tx.getIntersectingOps(boundary)
10:    OperationSet int_opsy ← ty.getIntersectingOps(boundary)
11:    OperationSet int_ops ← int_opsx ∩ int_opsy
12:    for all Operation int_op ∈ int_ops do
13:      Rectangle bd ← tx[int_op] × ty[int_op]
14:      boundary ← smallestEnclosing(bd, boundary)
15:      Parent[int_op] ← op
16:      tx.remove(int_op); ty.remove(int_op)
17:    end for
18:  until int_ops = ∅
19:  tx.add(boundary, op); ty.add(boundary, op)
20: end for
21: Map<Rectangle, OperationSet> results ← Map()
22: for all Operation op ∈ input_operations do
23:   Operation root_ob ← getRoot(parent, op)
24:   Rectangle rec ← tx[int_ob] × ty[int_ob]
25:   if not results.has_key(rec) then
26:     results[rec] ← List()
27:   end if
28:   results[rec].add(op)
29: end for
30: return results

```

Algorithm 1: The clustering algorithm

The clustering of textual operations yields text paragraphs and smaller objects like section headings. However, in the case of graphical operations, we can obtain consistent parts of images, but usually not complete figures yet. Outcomes of the clustering are utilised during the process of figures detection.

Algorithm 1 shows the pseudo-code of the clustering algorithm. The input of the algorithm consists of a set of preprocessed operators annotated with their affected area. The output is a division of the input set into disjoint clusters. Every cluster is assigned a boundary equal to the smallest rectangle containing boundaries of all included operations.

In the first stage of the algorithm (lines 6-20), we organise all input operations in a data structure of forest of trees. Every tree describes a separate cluster of operations. The second stage (lines 21-29) converts the results (clusters) into a more suitable format.

The clustering of operations is based on the relation of their rectangles being close to each other. Definition 1 formalises the notion of being close, making it useful for the algorithm.

Definition 1: Two rectangles are considered to be located close to each other if they are intersecting after expanding their boundaries in every direction by a margin.

The value by which rectangles should be extended is a parameter of the algorithm and might be different in various situations. In order to detect if rectangles are close to each other, we needed a data structure allowing to store a set of rectangles. This data structure was required to allow retrieving all stored rectangles that intersect a given one.

We have constructed the necessary structure using an important observation about the operation result areas. In our model all bounding rectangles have their edges parallel to the edges of the reference canvas on which the output of the operators is rendered. This allowed us to reduce our problem from the case of 2-dimensional rectangles to the case of 1-dimensional intervals. We can assume that edges of the rectangular canvas define the coordinates system. It is easy to prove that two rectangles of edges parallel to the axis of the coordinates system intersect only if both their projections in the directions of axis intersect. The projection of a rectangle into an axis is always an interval.

The observation made above has allowed us to build the required 2-dimensional data structure by remembering 2 one-dimensional data structures that allow to remember a number of intervals and for a given interval return the set of intersecting ones. Such a one-dimensional data structure has been provided by interval-trees (Edelsbrunner and Maurer 1981). Every interval inside the tree has an arbitrary object assigned to it, which in this case is a representation of the PDF operator. This object can be treated as an identifier of the interval. The data structure also implements a dictionary interface, mapping objects to actual intervals.

At the beginning, the algorithm initialises two empty interval trees representing projections on X and Y axis respectively. Those trees store values about projections of the biggest so-far calculated areas rather than about particular operators. Each cluster is represented by the most recently discovered operation belonging to it.

During the algorithm execution, each operator from the input set is considered only once. The order of processing is not important. The processing of a single operator proceeds as follows (the interior of the outermost “for all” loop of the algorithm).

1. Firstly, the boundary of the operation is extended by the width of margins. The spatial data structure described earlier is utilised to retrieve boundaries of all already detected clusters (lines 9-10)
2. The forest of trees representing clusters is updated. The currently processed operation is added without a parent. Roots of all trees representing intersecting clusters (retrieved in previous step) are attached as children of the new operation.
3. The boundary of the processed operation is extended to become the smallest rectangle containing all boundaries of intersecting clusters and the original boundary. Finally, all intersecting clusters are removed from the spatial data structure.
4. Lines 9-17 of the algorithm are repeated as long as there exist areas intersecting the current boundary. In some special cases, more than one iteration may be necessary.
5. Finally, the calculated boundary is inserted into the spatial data structure as a boundary of a new cluster. The currently processed operation is designed to represent the cluster and so, is remembered as a representant of the cluster.

After processing all available operations, the post-processing phase begins. All the trees are transformed into lists. The resulting data structure is a dictionary having boundaries of detected clusters as keys and lists of belonging operations as values. This is achieved in lines 21-29. During the process of retrieving the cluster to which a given operation belongs, we use a technique called path compression, known from the Find&Union data structure (Cormen, Leiserson and Rivest 1990).

3.2.2 Filtering of Clusters

Graphical areas detected by a simple clustering usually do not directly correspond to figures. The main reason for this is that figures may contain not only graphics, but also portions of text. Moreover, not all graphics present in the document must be part of a figure. For instance, common graphical elements not belonging to a figure include logos of institutions and text separators like lines and boxes; various parts of mathematical formulas usually include graphical operations; and in the case of slides from presentations, the graphical layout should not be considered part of a figure.

The above shows that the clustering algorithm described earlier is not sufficient for the purpose of figures detection and it yields a results set wider than expected. In order to take into account the aforementioned characteristics, precalculated graphical areas are subject to further refinement. This part of the processing is highly domain-dependent as it is based on properties of scientific publications in a particular domain, in this case publications of HEP. In the course of the refinement process, previously computed clusters can be completely discarded, extended with new elements, or some of their parts might be removed. In this subsection we discuss the heuristics applied for rejecting and splitting clusters of graphical operators.

There are two main reasons for rejecting a cluster. The first of them is a size being too small compared to a page size. The second is the figure candidate having its aspect ratio outside a desired interval of values.

The first heuristic is designed to remove small graphical elements appearing for example inside mathematical formulas, but also small logos and other decorations. The second one discards text separators and different parts of mathematical equations, such as a line separating numerator from a denominator inside a fraction. The thresholds used for filtering are provided as configurable properties of the algorithm and their values are assigned experimentally in a way maximising the accuracy of figures detection.

Additionally, the analysis of the order of operations forming the content stream of a PDF document may help to split clusters that were incorrectly joined by Algorithm 1. Parts of the stream corresponding to logical parts of the document usually form a consistent subsequence. This observation allows to construct a method of splitting elements incorrectly clustered together. We can assign content streams not only to entire PDF documents or pages, but also to every cluster of operations. The clustering algorithm presented in Algorithm 1 returns a set of areas with a list of operations assigned to each of them. The content stream of a cluster consists of all operations from such a set ordered in the same manner as in the original content stream of the PDF document. The usage of the original content stream allows us to define a distance in the content stream as follows:

Definition 2 If o_1 and o_2 are two operations appearing in the content stream of the PDF document, by the distance between these operations we understand the number of textual and graphical operations appearing after the first of them and before the second of them.

In order to detect situations when a figure candidate contains unnecessary parts, the content stream of a figure candidate is read from the first to the last operation. For every two subsequent operations, the distance between them in the sense of the original content stream is calculated. If the value is larger than a given threshold, the content stream is split into two parts which become separate figure candidates. For both candidates, a new boundary is calculated.

This heuristic is especially important in the case of less formal publications such as slides from presentations at conferences. Presentation slides tend to have a certain amount of graphics appearing on

every page and not carrying any meaning. Simple geometrical clustering would connect elements of page style with all the rest of the document content. Measuring the distance in the content stream and defining a threshold on the distance facilitates the distinction between the layout and the rest of the page. This technique might be also useful in order to automatically extract the template used for a presentation, although this transcends the scope of this publication.

3.3 Clustering of Textual Operators

The same algorithm that is applied to cluster graphical elements can be used to cluster parts of text. Detecting larger logically consistent parts of text is important because they should be treated as single entities during subsequent processing. This comprises inclusion inside a figure candidate (captions of axes, parts of a legend etc...), classification of a text paragraph as a figure caption and so on.

3.4 Inclusion of Text Parts

The next step in figures extraction involves the inclusion of lost text parts inside figure candidates. At the stage of operations clustering, only the operations of the same type (graphical or textual) were considered. The results of those initial steps become subsequently the input to the clustering algorithm that will detect relations between previously detected entities. By doing this, we move one level farther in the process of abstracting from operations. Initially we start from basic meaningless operations. Later we detect parts of graphics and text and finally we are able to see the relations between both.

Not all clusters detected at this stage are interesting because some clusters might consist uniquely of text areas. Only those results that include at least one graphical cluster may be subsequently considered figure candidates.

Another round of heuristics allows to mark unnecessary intermediate results as deleted. Applied methods are very similar to those described in Section 3.2.2, only thresholds deciding on the rejections must change because we operate on geometrically much larger entities. Also the way of application is different - candidates rejected at this stage can be later restored to the status of a figure. Instead of permanently removing, heuristics of this stage only mark figure candidates as rejected. This happens in the case of the candidates having incorrect aspect ratio, incorrect sizes or consisting only of horizontal lines (which is usually the case with mathematical formulas but also tables).

In addition to using the aforementioned heuristics, having clusters consisting of a mixture of textual and graphical operations allows to apply new ones. During the next phase, we analyse the type of operations rather than their relative location. In some cases, steps described earlier might detect objects that should not be considered a figure, such as text surrounded by a frame. This situation can be recognised by the calculation of a ratio between the number of graphical and textual operations in the content stream of a figure candidate. In our approach we have defined a threshold which indicates which figure candidates should be rejected because they contain too few graphics. This allows to remove for instance blocks of text decorated with graphics for aesthetic reasons. The ratio between numbers of graphical and textual operations is smaller in the case of tables than in the case of figures so extending the heuristic with an additional threshold could improve the table/figure distinction. Another heuristic analyses ratio between the total area of graphical operations and the area of the entire figure candidate.

Subsequently, we mark as deleted figure candidates containing horizontal lines as the only graphical operations. These candidates describe tables or mathematical formulas which have survived previous steps of the algorithm. Tables can be reverted to the status of figure candidates in later stages of processing.

Figure candidates that survive all the phases of filtering are finally considered to be figures. Figure 2 shows a fragment of a publication page with indicated text areas and final figure candidates detected by the algorithm.

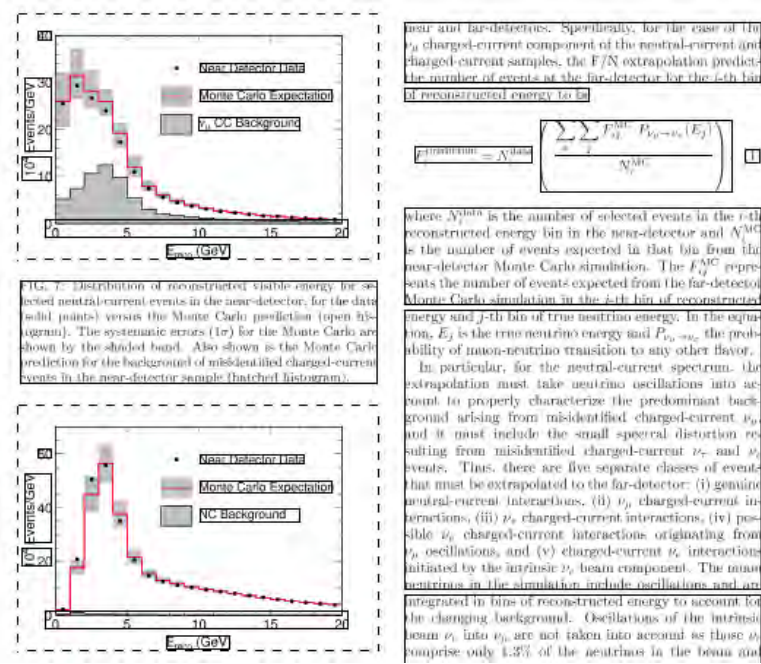


Figure 2: A fragment of the PDF page with boxes around every detected text area and each figure candidate. Dashed rectangles indicate figure candidates. Solid rectangles indicate text areas.

3.5 Detection and Matching of Captions

The input of the part of the algorithm responsible for detecting figure captions consists of previously determined figures and all text clusters. The observation of scientific publications shows that typically captions of figures start with a figure identifier (for instance see the grammar for figure captions proposed by Bathia et al. (Bhatia, Lahiri and Mitra 2009)).

The identifier usually starts with a word describing a figure type and is followed by a number or some other unique identifier. In the case of more complex documents, the figure number might have a hierarchical structure reflecting for example the chapter number. The set of possible figure types is very limited. In the case of HEP publications, the most usual combinations include words “Figure”, “Plot” and different variations of their spelling and abbreviating.

During the first step of the caption detection, all text clusters from the publication page are tested for the possibility of being a caption. This consists of matching the beginning of the text contained in a textual cluster with a regular expression determining what is a figure caption. The role of the regular expression is to elect strings starting with one of the predefined words, followed by an identifier or beginning of a sentence. The identifier is subsequently extracted and included in the metadata of a caption. The caption detection has to be designed to reject paragraphs of the type “Figure 1 presents results of (...)”. In order to achieve this, we reject the possibility of having any lower case text after the figure identifier.

Having the set of all the captions, we start searching for corresponding figures. All previous steps of the algorithm take into account the division of a page into text columns (see Section 3.7 about the layout detection). When matching captions with figure candidates, we do not take into account the page layout.

Matching between figure candidates and captions happens at every document page separately. We consider every detected caption once, starting with those located at the top of the page and moving downwards towards the end. For every caption we search figure candidates lying nearby. First we search above the caption and in the case of failure, we move below the caption. We take into account all figure candidates, including those rejected by heuristics.

In the case of finding multiple figure candidates corresponding to a caption, we merge them into a single figure, treating previous candidates as subfigures of a larger figure. We also include small portions of text and graphics previously rejected from figure candidates which lie between figure and caption and between different parts of a figure. These parts of text usually contain identifiers of the subfigures. The amount of unclustered content that can be included in a figure is a parameter of the extraction algorithm and is expressed as a percentage of the height of the document page.

It might happen that captions are located in a completely different location, but this case is rare and tends to appear in older publications. The distance from the figure is calculated based on the page geometry. The captions should not be too distant from the figure.

3.6 Generation of the Output

The choice of the format in which data should be saved at the output of the extraction process should take into account further requirements.

The most obvious use case of displaying figures to end users in response to text-based search queries does not yield very sophisticated constraints. A simple raster graphics annotated with captions and possibly some extracted portions of metadata would be sufficient. Unfortunately, the process of generating raster representations of figures might lose many important pieces of information that could be used in the future for an automatic analysis.

In order to store as much data as possible, apart from storing the extracted figures in a raster format (e.g., PNG), we also decided to preserve their original vector character. Vector graphics formats, similarly to PDF documents, contain information about graphical primitives. Primitives can be organised in larger logical entities. Sometimes rendering of different primitives leads to a modification of the same pixel of resulting image. Such a situation might happen for example when circles are used to draw data points lying nearby on the same plot. In order to avoid such issues, we convert figures into Scalable Vector Graphics (Ferraiolo 2001) format.

On the implementation level, the extraction of vector representation of a figure proceeds in a manner similar to regular rendering of a PDF document. The interpreter preserves the same elements of the state and allows their modification by transformation operations. A virtual canvas is created for every detected figure. The content stream of the document is processed and all the transformation operations are executed modifying the interpreters state. The textual and graphical operators are also interpreted, but they affect only the appropriate canvas of the figure to which the operation belongs. If a particular operation does not belong to any figure, no canvas is affected. The behaviour of graphical canvases used during the SVG generation is different from the case of raster rendering. Instead of creating graphical output, every operation is transformed into a corresponding primitive and saved within a SVG file.

The PDF format was designed in such a manner that the number of external dependencies of a file is minimised. This design decision led to the inclusion of the majority of fonts in the document itself. It would

be possible to embed font glyphs in the SVG file and use them in order to render strings. However, for the sake of simplicity, we decided to omit font definitions in the SVG output.

A text representation is extracted from every text operation and the operation is replaced by a SVG text primitive with a standard font value. This simplification affects how the output looks like, but the amount of formatting information that is lost is minimal. Moreover, this does not pose a problem as vector representations are intended to be used during automatic analysis of figures rather than for displaying purposes. A possible extension of the presented method could involve embedding complete information about used glyphs.

Finally, the generation of the output is completed with some metadata elements. An exhaustive categorisation of the metadata that can be compiled for figures could be the customisation of the one proposed by Liu et al (Liu. Y, et al. 2007) for table metadata. In the case of figures, the following categories could be distinguished: 1) environment/geography metadata (information of the document where the figure is located); 2) affiliated metadata (e.g., captions, references, or footnotes); 3) layout metadata (information about the original visualisation of the figure); 4) content data; 5) and figure type metadata. For the moment, we compile only environment/geography metadata and affiliated metadata.

The geography/environment metadata consists of the document title, the document authors, the document date (creation and publication), and the exact location of a figure inside a publication (page and boundary). Most of these elements are provided by simply making a reference to the original publication in the Inspire repository. The affiliated metadata consists of the text caption and the exact location of the caption in the publication (page and boundary). In the future, metadata from other categories will be annotated for each figure.

3.7 Detection of the Page Layout

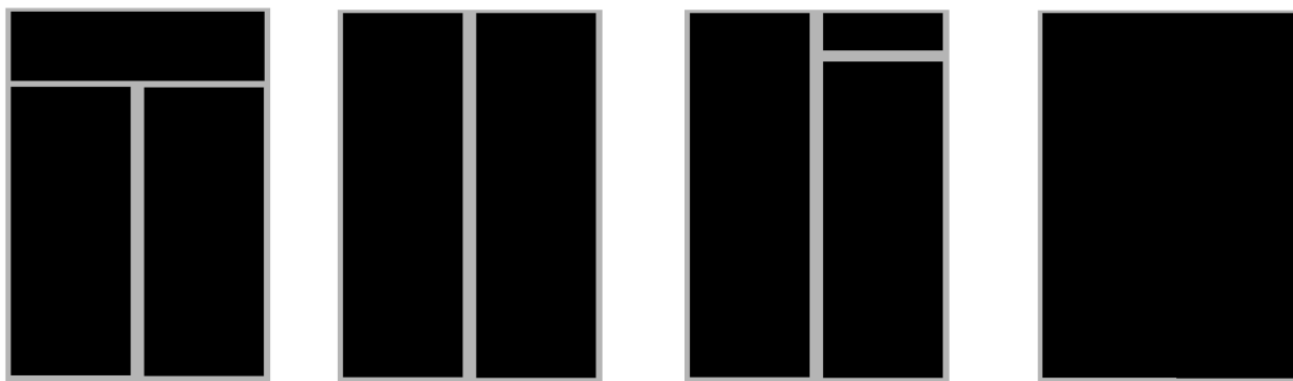


Figure 3: Sample page layouts that might appear in a scientific publication. The black colour indicates areas where content is present.

In this section we discuss how to detect the page layout, an issue which has been omitted in the main description of the extraction algorithm, but which is essential for an efficient detection of figures. Figure 3 depicts several possibilities of organising content on the page. As mentioned in previous sections, the method of clustering operations based on their geometrical position may fail in the case of documents having a complex page layout. The content appearing in different columns should never be considered belonging to the same figure. This cannot be assured without enforcing additional constraints during the clustering phase.

In order to address this difficulty, we enhanced the figure extractor with a preprocessing phase of detecting the page layout. Being able to identify how the document page is divided into columns, enables us to execute the clustering within every column separately. It is intuitively obvious, what can be understood as a page layout, although in order to provide a method of calculating such, we need a more formal definition, which we provide below.

By the layout of a page, we understand a particular division of a page into areas called columns. Each area is a sum of disjoint rectangles. The division of a page into areas must satisfy a set of conditions summarised in Definition 3.

Definition 3: Let P be a rectangle representing the page. The set D containing subareas of a page is called a page division if and only if

$$\begin{aligned}
 &Q = P \\
 &\forall_{Q \in D} \\
 &\forall_{x, y \in D} x \cap y = \emptyset \\
 &\forall_{Q \in D} Q \neq \emptyset \\
 &\forall_{Q \in D} \exists_{R = \{x : x \text{ is a rectangle}, \forall_{y \in R \setminus \{x\}} y \cap x = \emptyset\}} Q = \bigcup_{x \in R} x
 \end{aligned}$$

Every element of a division is called a page area.

In order to be considered a page layout, borders of areas from the division must not intersect the content of the page. Definition 3 does not guarantee that the layout is unique. A single page might be assigned different divisions satisfying the definition. Additionally, not all valid page layouts are interesting from the point of view of figures detection. The segmentation algorithm calculates one of such divisions, imposing additional constraints on the detected areas. The layout-calculation procedure utilises the notion of separators, introduced by Definition 4.

Definition 4: A vertical (or horizontal) line inside a page or on its borders is called a separator if its' horizontal (vertical) distance from the page content is larger than a given constant value.

The algorithm consists of two stages. First, the vertical separators of a sufficient length are detected and used to divide the page into disjoint rectangular areas. Each area is delimited by two vertical lines each of which forms a consistent interval inside of one of the detected vertical separators. At this stage, horizontal separators are completely ignored. Figure 4 shows a fragment of a publication page processed by the first stage of the layout-detection. The upper horizontal edge of one of the areas lies too close too close to two text lines. With the constant of the Definition 4 chosen to be sufficiently large, this edge would not be a horizontal separator and thus the generated division of the page would require additional processing to become a valid page layout. The second stage of the algorithm transforms the previously detected rectangles into a valid page layout by splitting rectangles into smaller parts and by joining appropriate rectangles to form a single area.

Mg	$3s4s^3S_1 \rightarrow 3s3p^3P_0$	516.73	517.11	this	$-(0.09 \pm 0.01)$	-0.017	this
	$3s4s^3S_1 \rightarrow 3s3p^3P_1$	517.27	517.51	work	$-(0.06 \pm 0.01)$	-0.012	work
	$3s4s^3S_1 \rightarrow 3s3p^3P_2$	518.36	518.52		$-(0.06 \pm 0.01)$	-0.012	

Table 1. Electronic transitions of various elements measured at an increased helium pressure. The table includes the free atomic transitions, the wavelength of the transitions in superfluid helium under saturated vapour pressure and the pressure line shifts. Also mentioned is the change of the wavelength because of a pressure increase relative to the wavelength of the transitions at saturated vapour pressure.

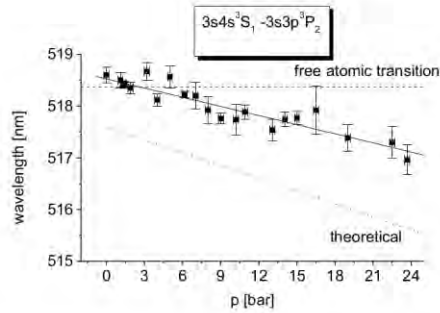


Fig. 9. Emission wavelength of the $3s4s^3S_1 \rightarrow 3s3p^3P_2$ transition of the magnesium atom as a function of the helium pressure. The dotted line corresponds to the wavelength calculated by use of the bubble model. The dashed line corresponds to the free atomic transition at 518.37 nm.

deviation is about 14 nm [13].

The quality of the agreement of the calculated and measured pressure shifts for all three lines can be tested with a statistical hypothesis test, the students test. The deviation of the three values is compatible with statistical fluctuations. Therefore a mean pressure line shift of $(0.07 \pm 0.01 \text{ nm/bar})$ can be derived. This very good consistency between the experimental and the theoretical values allows the conclusion that the magnesium atom seem to maintain a bubble like structure under increased helium pressures. The pressure shift is monotonous.

5 Discussion

As a consequence of the higher pressure the bubble like defect shrinks, i.e. the equilibrium radius decreases. The repulsive part of the pair potential energies due to Pauli forces rises in the upper P state already at larger radii than for the lower S state which implies a smaller wavelength for emitted radiation.

Figure 4: Example of intermediate layout-detection results requiring the refinement

Algorithm 2 shows the pseudo-code of the detection of vertical separators. The input of the algorithm consists of the image of the publication page. The output is a list of vertical separators aggregated by their x-coordinates. Every element of this list consists of two elements: an integer indicating the x-coordinate and the list of y-coordinates describing the separators. The first element of this list indicates the y-coordinate of the beginning of the first separator. The second element is the y-coordinate of the end of the same separator. The third and fourth elements describe the second separator and the same mechanism is used for the remaining separators (if they exist).

The algorithm proceeds according to the sweeping principle (Cormen, Leiserson and Rivest 1990) known from the computational geometry. The algorithm reads the publication page starting from the left. For every x-coordinate value, a set of corresponding vertical separators is detected (lines 9 – 18). Vertical separators are searched as consistent sequences of blank points. A point is considered blank if all the points in its horizontal surrounding of the radius defined by the constant from Definition 5 are of the background colour. Not all blank vertical lines can be considered separators. Short empty spaces usually delimit lines of text or different small units of the content. In line 11 we test detected vertical separators for being long enough.

If a separator has been detected in a particular column of a publication page, the adjacent columns also tend to contain similar separators. Lines 19-31 of the algorithm are responsible for electing the longest candidate among the adjacent columns of the page. The maximisation is performed across a set of adjacent columns for which at least one separator exists.


```

1: Input: the page image
2: Output: vertical separators of the input page
3: List<Pair<int, List<int>>> separators  $\leftarrow \emptyset$ 
4: int max_weight  $\leftarrow 0$ ;
5: boolean maximizing  $\leftarrow$  false
6: for all  $x \in \{\min_x \dots \max_x\}$  do
7:   emptyb  $\leftarrow 0$ , current_eval  $\leftarrow 0$ 
8:   empty_areas  $\leftarrow$  List()
9:   for all  $y \in \{0 \dots \text{page\_height}\}$  do
10:    if point at (x, y) is not blank then
11:      if  $y - \text{empty}_b - 1 > \text{height}_{\min}$  then
12:        empty_areas.append(emptyb)
13:        empty_areas.append( $y = \text{page\_height} ? y : y-1$ )
14:        current_eval  $\leftarrow$  current_eval +  $y - \text{empty}_b$ 
15:      end if
16:      emptyb  $\leftarrow y + 1$ 
17:    end if
18:  end for
  {We have already processed the entire column. Now we are comparing with adjacent already processed columns}
19:  if max_weight < current_eval then
20:    max_weight  $\leftarrow$  current_eval
21:    max_separators  $\leftarrow$  empty_areas
22:    maxx  $\leftarrow x$ 
23:  end if
24:  if maximizing then
25:    if empty_areas =  $\emptyset$  then
26:      separators.add(<maxx, max_separators>)
27:      maximizing  $\leftarrow$  false, max_weight  $\leftarrow 0$ 
28:    end if
29:  else
30:    maximizing  $\leftarrow$  (empty_areas  $\neq \emptyset$ )
31:  end if
32: end for
33: return separators

```

Algorithm 2: Detecting vertical separators

The detected separators are used to create the preliminary division of the page (Similar to the one from the example of Figure 44). Similarly to the previous step, separators are considered one by one in the order of increasing x coordinate. At every moment of the execution, the algorithm maintains a division of the page into rectangles. This division corresponds only to the already detected vertical separators. Updating the previously considered division is facilitated by processing separators in a particular well-defined order.

Before presenting the final outcome, the algorithm must refine the previously-calculated division. This happens in the second phase of the execution. All the horizontal borders of the division are then moved along adjacent vertical separators until they become horizontal separators in the sense of Definition 4. Typically, moving the horizontal borders result in dividing already existing rectangles into smaller ones. If such a situation happens, both newly created parts are assigned to different page layout areas. Sometimes when moving separators is not possible, different areas are combined together, forming a larger one.

4 Tuning and Testing

The extraction algorithm described here has been implemented in Java and tested on a random set of scientific articles coming from the Inspire repository. The testing procedure has been used to evaluate the quality of the method, but also allowed to tweak the parameters of the algorithm in order to maximise the outcomes.

4.1 Preparation of the testing set

In order to prepare the testing set, we randomly selected 207 documents stored in INSPIRE. In total, these documents consisted of 3728 pages which contained 1697 figures altogether.

The records have been selected according to a uniform probability distribution across the entire record space. This way, we have created a collection that is representative for the entire INSPIRE including historical entries.

Currently, INSPIRE consists of: 1,140 records describing publications written before 1950; 4,695 between 1950 and 1960; 32,379 between 1960 and 1970; 108,525 between 1970 and 1980; 167,240 between 1980 and 1990; 251,133 between 1990 and 2000; and 333,864 in the first decade of XXIst century. In total, up to July 2012, INSPIRE manages 952,026 records. It can be seen that the rate of growth has increased with time and most of INSPIRE documents come from the last decade.

The results on such a testing set should accurately estimate the efficiency of extraction for existing documents but not necessarily for new documents, being ingested into INSPIRE. This is because INSPIRE contains entries describing old articles which were created using obsolete technologies or scanned and encoded in PDF. The extraction algorithm is optimised for born-digital objects. In order to test the hypothesis that the extractors provides better results for newer papers, the testing set has been split into several subsets. The first set consists of publications published before 1980. The rest of the testing set has been split into subsets corresponding to decades of publication.

In order to simplify the counting of correct figure detections and to provide a more reliable execution and measurement environment, every testing document has been split into a number of PDF documents consisting of a single page. Subsequently, every single page document has been manually annotated with the number of figures appearing inside.

4.2 Execution of the Tests

The efficient execution of the testing was possible thanks to a special script executing the plots extractor on every single page separately and then computing the total number of successes and failures. The script allows the execution of tests in a distributed heterogeneous environment and allows dynamic connection and disconnection of computing nodes. In the case of a software failure, the extraction request is resubmitted to a different computation node, allowing to avoid problems related to a worker node configuration rather than to the algorithm implementation itself.

During the preparation of the testing set, we manually annotated all the expected extraction results. Subsequently, the script compared these metadata with the output of the extractor. Using aggregated numbers from all extracted pages allowed us to calculate efficiency measures of the extraction algorithm. As quality measures, we used recall and precision (Baeza-Yates and Ribeiro-Neto 1999). Their definitions are included in the following equations:

$$recall = \frac{\#correctly\ extracted\ figures}{\#figures\ present\ in\ the\ testset}$$

$$precision = \frac{\#correctly\ extracted\ figures}{\#extracted\ figures}$$

At every place where we needed a single comparable quality measure rather than two semi-independent numbers, we have used a harmonic average of the precision and the recall (Baeza-Yates and Ribeiro-Neto 1999).

$$harmonic\ average = \frac{2}{\frac{1}{precision} + \frac{1}{recall}}$$

Table 1 summarises the results obtained during the test execution for every subset of our testing set. Figure 5 shows the dependency of recall and precision on the time of publication. The extractor parameters used in this test execution were chosen based on intuition and small number of manually triggered trials. In the next section we describe an automatic tuning procedure we have used to find the most optimal algorithm arguments.

Table 1: Results of the test execution

	-1980	1980-90	1990-2000	2000-10	2010-12
Number of existent figures	114	60	170	783	570
Number of correctly detected figures	59	53	164	703	489
Number of incorrectly detected figures	26	78	65	40	73
Total number of pages	85	136	760	1919	828
Number of correctly processed pages	20	44	712	1816	743

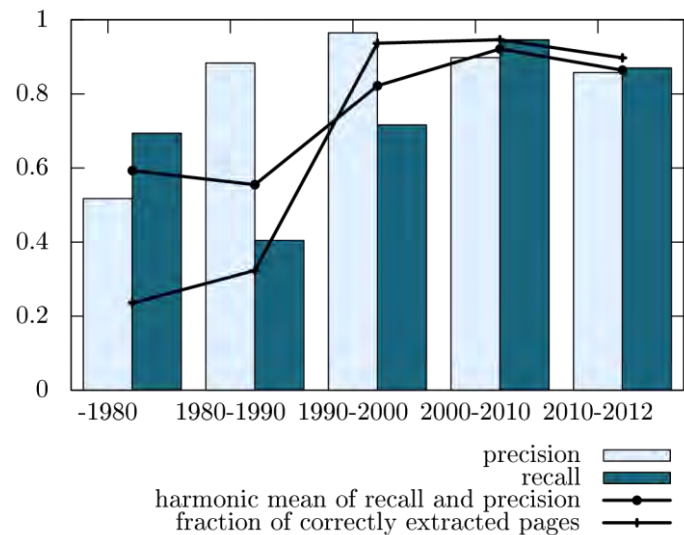


Figure 5: Recall and Precision as functions of decade of the date of the publication.

It can be seen that, as expected, the efficiency increases with the increasing time of publication. A total recall and precision for all samples since 1990, which constitutes a majority of the INSPIRE corpus, were both 88%.

Precision and recall based on the correctly detected figures do not give a full image of the algorithm efficiency because the extraction has been executed on a number of pages not containing any figures. The correctly extracted pages not having any figures do not appear in the recall and precision statistics because in their case the expected and detected number of figures are both equal to 0.

Besides recall and precision, Figure 5 depicts also the fraction of pages which have been extracted correctly. Taking into account the samples since 1990, 3,271 pages out of 3,507 have been detected completely correctly, which makes 93% success rate counted by number of pages. As it can be seen, this measure is higher than both the precision and the recall.

The analysis of the extractor results in the case of failure shows that in many cases, even if results are not completely correct, they are not far from the expectation. There are different reasons of the algorithm failing. Some of them may result from non-optimal choice of algorithm parameters, others from document layout being too far from the assumed one. In some rare cases, even manual inspection of the document does not allow an obvious identification of figures.

4.3 The Automatic Tuning of Parameters

In previous section we have shown the results obtained by executing the extraction algorithm on a sample set. During this execution we were using extractor arguments which seemed to be the most correct based on our observation but also on other research (Cleveland 1984)(typical sizes of figures, margin sizes etc.). This way of algorithm configuration was useful during the development, but is not likely to yield the best possible results. In order to find better parameters, we have implemented a method of automatic tuning. Metrics described in the previous section (recall and precision) provided a good method of measuring the efficiency of the algorithm running based on given parameters.

The choice of optimal parameters can be relative to the choice of documents on which the extraction is to be performed. The way in which the testing set has been selected, allowed us to use it as representative for the HEP publications. In order to tune the algorithm, we have used a described subset of testing set from the previous step as a reference. The subset consisted of all entries created after 1990. This allowed us to minimise the presence of scanned documents which, by design, cannot be correctly processed by our method.

The adjustment of parameters has been performed by a dedicated script which has executed the extraction using various parameter values and has read results. The script has been configured with a list of tuneable parameters together with their type and allowed values range. Additionally, the script had the knowledge of the believed best value, which was the one used in previous testing.

In order to decrease the complexity of training, we have made several assumptions about the parameters. These assumptions are only an approximation of real nature of parameters but the practice has shown that they are good enough to permit the optimisation:

- We assume that the precision and recall are continuous with respect to the parameters. This allows us to assume that efficiency of the algorithm for parameter values close to a given one will be close. The Optimisation has proceeded by sampling the parametric space in a number of points and executing tests using the selected points as parameter values. Having N parameters to optimise and dividing the space of every parameter into M regions leads to the execution of M^N tests. Execution of every test is a timely operation due to the size of the training set.
- We assume that parameters are independent from each other. This means that we can divide the problem of finding an optimal solution in the N -dimensional space of N configuration arguments into finding N solutions in 1-dimensional subspaces. Such an assumption seems to be intuitive and considerably reduces the number of necessary tests from $O(M^N)$ to $O(M \cdot N)$, where M is the number of samples taken from a single dimension.

In our tests, the parametric space has been divided into 10 equal intervals in every direction. In addition to checking the extraction quality in those points, we have executed one test for the so-far best argument. In order to increase the level of fine-tuning of the algorithm, each test has been re-executed in the region, where chances of finding a good solution were considered the highest. This consisted of a region centred around the highest result and having a radius of 10% of the parameter space.

Figure 6 and Figure 7 show the dependency of the recall and the precision on an algorithm parameter. The parameter depicted in Figure 6 indicates what minimal aspect ratio the figure candidate must have in order to be considered a correct figure. It can be seen that tuning this heuristic increases the efficiency of the extraction. Moreover, the dependency of recall and precision on the parameter is monotonic which is the most compatible with the chosen optimisation method.

The parameter of Figure 7 specifies which fraction of the area of the entire figure candidate has to be occupied by graphical operations. This parameter has a lower influence on the extraction efficiency. Such a situation can happen when more than one heuristic influences the same aspect of the. This is contradictory with the assumption of parameter independence, but we have decided to use the present model for the simplicity.

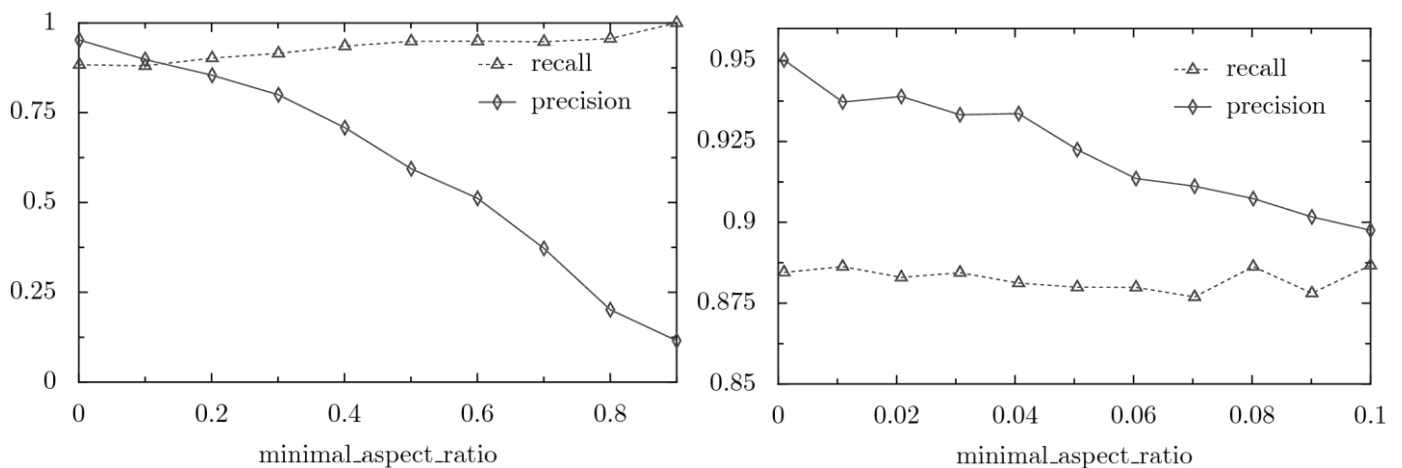


Figure 6: Effect of the minimal aspect ratio on precision and recall

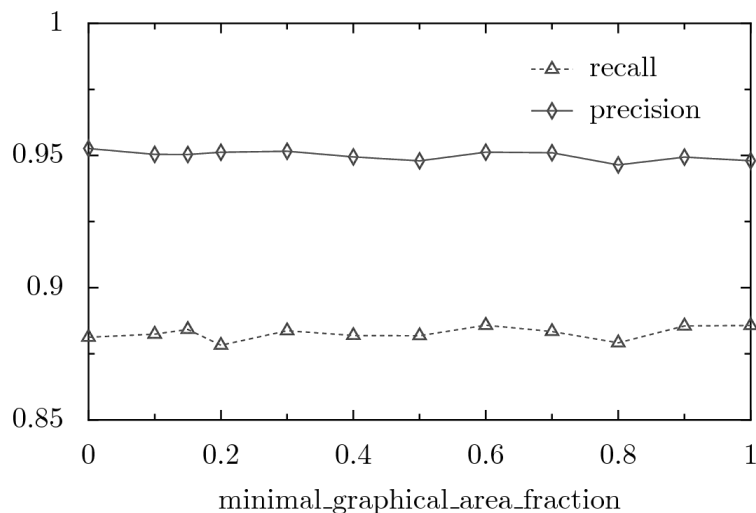


Figure 7: Effect on the precision and recall of the area fraction occupied by graphical operations

After executing the optimisation algorithm, we have managed to achieve the recall of 94.11% and the precision of 96.6% which is a considerable improvement with respect to previous results of 88%.

5 Conclusions and Future Work

This work has presented a method for extracting figures from scientific publications in a machine readable format, which is the main step towards the development of services enabling access and search of images stored in scientific digital libraries. In recent years, figures have been gaining increasing attention in the Digital Libraries community. However, little has been done to decipher the semantics of these graphical representations and to bridge the semantic gap between content which can be understood by machines and this which is managed by digital libraries. Extracting figures and storing then in the uniform and machine readable format constitutes the first step towards the extraction and the description of the internal semantics of figures. Storing semantically described and indexed figures would open completely new possibilities of accessing the data and discovering connections between different types of publishing artefacts and different resources describing related knowledge (Praczyk, Nogueras-Iso and Kaplun, et al. 2011).

Our method of detecting fragments of PDF documents that correspond to figures is based on a series of observations of the character of publications. However, tests have shown that additional work is needed to improve the correctness of the detection. Also the performance should be re-evaluated after we have a large set of correctly annotated figures, confirmed by users of our system. The heuristics used by the algorithm are based on a number of numeric parameters which we have tried to optimise using automatic techniques. The tuning procedure has made several arbitrary assumptions on the nature of the dependency between parameters and extraction results. A future approach to the parameter optimisation, requiring much more processing, could involve the execution of a genetic algorithm (Theodoridis and Koutroumbas 2006) which would treat the parameters as gene samples. This could potentially allow a discovery of a better parameter set because a smaller set of assumptions would be imposed on the parameters. A vector of algorithm parameters could play the role of a gene and random mutations could be introduced to previously considered and subsequently crossed genes. The evaluation and selection of surviving genes could be performed by the usage of the metrics described in section 4.2. Another approach to improving the quality of the tuning could involve extending the present algorithm by a discovery of mutually dependent parameters and usage of special techniques (relaxing the assumptions) to fine-tune in subspaces spanned by these parameters.

All of our experiments have been performed using a corpus of publications from HEP. The usage of the extraction algorithm on a different corpus would require tuning the parameters for the specific domain of application. For the area of HEP, we can also consider preparing several sets of execution parameters varying by decade of document publication or by other easy to determine characteristics. Subsequently, we could decide which extraction method to run, based on those metrics.

In addition to a better tuning of the existing heuristics, there are improvements which can be made at the level of the algorithm. As an example, we could mention extending the process of clustering text parts. In the current implementation, the margins by which textual operations are extended during the clustering process are fixed as algorithm parameters. This approach proved to be robust in most cases. In fact, distances between text lines tend to be different depending on the currently utilised style. Every text portion tends to have one style that dominates. An improved version of the text-clustering algorithm could use local rather than global properties of the content. This would not only allow to correctly handle the

entire document written using different text styles, but also help to manage cases of single paragraphs differing from the rest of the content.

Another important, not implemented yet, improvement related to figure metadata is the automatic extraction of figure references from the text content. Important information about figure content might be stored in the surroundings of the place where publication text refers to a figure. Furthermore, the meta-data could be extended by the usage of some type of classifier which would assign a graphics type to the extracted result. Currently, we are only distinguishing between tables and figures based on simple heuristics involving number and type of graphical areas and the text inside of the detected caption. In the future, we could detect line-plots from photos, histograms and so on. Such a classifier could be implemented using Artificial Intelligence techniques such as Support Vector Machines (Theodoridis and Koutroumbas 2006).

Finally, partial results of the figures extraction algorithm might be useful in performing other PDF analyses:

- The usage of clustered text areas could allow a better interpretation and indexing of textual content stored in digital libraries with full text access. Clusters of text tend to describe logical parts like paragraphs, section and chapter titles and so on. A simple extension of the current schema could allow the extraction of predominant formatting style of the text encoded in a page area. Text parts written in different styles could be indexed in a different manner giving for instance more importance to segments written with larger font.
- In section 3.4 we mentioned that the algorithm detects not only figures, but also tables. A heuristic is being used in order to distinguish tables from different types of figures. Our present effort concentrates on correct treatment of figures, but a useful extension could allow extraction of different types of entities. For instance, another common type of content ubiquitous in HEP documents are mathematical formulas. Thus, in addition to figures, it would be important to extract tables and formulas in structured format allowing a further processing.

The internal architecture of the implemented prototype of the figure extractor allows easy implementation of extension modules which can compute other properties of PDF documents.

Acknowledgements

This work has been partially supported by CERN, and the Spanish Government through the project TIN2012-37826-C02-01.

6 Bibliography

Baeza-Yates, R., and B. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press Books. Addison-Wesley, 1999.

Bhatia, S., S. Lahiri, and P. Mitra. "Generating Synopses for Document-Element Search." *Proceedings of the 18th ACM Conference on Information and Knowledge Management*. New York, 2009. 2003-2006.

Browner, W., S. Kataria, S. Das, P. Mitra, and C. L. Giles. "Segregating and Extracting Overlapping Data Points in Two-dimensional Plots." *Proceedings of the 8th ACM/IEEE-CS joint conference on Digital Libraries*. New York, 2008. 276-279.

- Chao, H., and J. Fan. "Layout and Content Extraction for PDF Documents." *Document Analysis Systems*, 2004: 213-224.
- Cleveland, W. S. "Graphs in Scientific Publications." *The American Statistician* 38, no. 4 (1984): 261-269.
- Cormen, T. H., C. E. Leiserson, and R. L. Rivest. *Introduction to Algorithms*. Cambridge: MIT Electrical Engineering and Computer Science Series, 1990.
- Edelsbrunner, H., and H. A. Maurer. "On the Intersection of Orthogonal Objects." *Information Processing Letters*, 1981: 13.
- Eichhorn, G. "Trends in Scientific Publishing at Springer." *Future professional communication in astronomy II*. Cambridge, MA, 2011.
- Elsevier. *SciVerse Science Direct: Image Search*. 2012.
<http://www.info.sciverse.com/sciencedirect/using/searching-linking/image> (accessed 5 24, 2013).
- Ferraiolo, J. *Scalable Vector Graphics (SVG) 1.0 Specification*. Iuniverse Inc., 2001.
- Hearst, M. A., A. Divoli, J. Ye, and M. A. Wooldridge. "Exploring the Efficacy of Caption Search for Bioscene Journal Search interfaces." *Proceedings of the Workshop on Bio NLP 2007: Biological, Translational and Clinical Language Processing*. 2007. 73-80.
- Holtkamp, A., S. Mele, T. Simko, and T. Smith. "INSPIRE: Realizing the Dream of a Global Digital Library in High-Energy Physics." *3rd Workshop Conference: Towards a Digital Mathematics Library*. Paris, 2010. 83-92.
- Johnston, L. "Web Reviews: See the Science: Scitech image databases." *Sci-Tech News* 65 (2011).
- Kataria, S. "On Utilization of Information Extracted From Graph Images in Digital Documents." *Bulletin of IEEE Technical Comittee on Digital Libraries* 4 (2008).
- Kataria, S., W. Browuer, P. Mitra, and C. L. Giles. "Automatic Extraction of Data Points and Text Blocks from 2-dimensional plots in digital documents." *Proceedings of the 23rd National Conference on Artificial Intelligence*. 2008. 1169-1174.
- Liu, Y, K. Bai, P. Mitra, and C. L. Giles. "Tableseer: Automatic Table Metadata Extraction and Searching in Digital Libraries." *JCDL'07*. Vancouver, 2007.
- Praczyk, P. A., J. Nogueras-Iso, S. Dallmeier-Tiessen, and M. Whalley. "Integrating Scholarly Publications and Research Data - Preparing for Open Science, a Case Study from High-Energy Physics with Special Emphasis on (Meta)data Models." *Metadata and Semantics Research - CCIS 343* (2012): 146-157.
- Praczyk, P. A., J. Nogueras-Iso, S. Kaplun, and T Simko. "A Storage Model for Supporting Figures and Other Artefacts in Scientific Libraries: the Case Study of Invenio." *4th Workshop on very Large Digital Libraries*. Berlin, 2011.
- Russell, S., and P. Norvig. *Artificial Intelligence: A Modern Approach*. Prentice Hall, 2009.
- Theodiridis, S., and K. Koutroumbas. *Pattern Recognition, third Edition*. Academic Press, 2006.

