



HAL
open science

A Hybrid Segmentation of Web Pages for Vibro-Tactile Access on Touch-Screen Devices

Waseem Safi, Fabrice Maurel, Jean-Marc Routoure, Pierre Beust, Gaël Dias

► **To cite this version:**

Waseem Safi, Fabrice Maurel, Jean-Marc Routoure, Pierre Beust, Gaël Dias. A Hybrid Segmentation of Web Pages for Vibro-Tactile Access on Touch-Screen Devices. 3rd Workshop on Vision and Language (VL 2014) associated to 25th International Conference on Computational Linguistics (COLING 2014), Aug 2014, dublin, Ireland. pp.95 - 102. hal-01076613

HAL Id: hal-01076613

<https://hal.science/hal-01076613>

Submitted on 22 Oct 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Hybrid Segmentation of Web Pages for Vibro-Tactile Access on Touch-Screen Devices

Waseem Safi¹ Fabrice Maurel¹ Jean-Marc Routoure^{1,2} Pierre Beust¹ Gaël Dias¹

¹ University of Caen Basse-Normandie - UNICAEN

² National Superior Engineering School of Caen - ENSICAEN

14032 Caen - France

firstName.lastName@unicaen.fr

Abstract

Navigating the Web is one of important missions in the field of computer accessibility. Many specialized techniques for Visually Impaired People (VIP) succeed to extract the visual and textual information displayed on digital screens and transform it in a linear way: either through a written format on special Braille devices or a vocal output using text-to-speech synthesizers. However, many researches confirm that perception of the layout of web pages enhances web navigation and memorization. But, most existing screen readers still fail to transform the 2-dimension structures of web pages into higher orders. In this paper, we propose a new framework to enhance VIP web accessibility by affording a “first glance” web page overview, and by suggesting a hybrid segmentation algorithm to afford nested and easy navigation of web pages. In particular, the web page layout is transformed into a coarse grain structure, which is then converted into vibrating pages using a graphical vibro-tactile language. First experiments with blind users show interesting issues on touch-screen devices.

1 Introduction

In October 2013, the world health organization estimated that the number of Visually Impaired People (VIP) in the world is 285 million: 39 million of them are blind and 246 million have low vision. In particular, the organization defined four levels of visual functions depending on the international classification of diseases: normal vision, moderate visual impairment, severe visual impairment and blindness.

VIP depend on screen readers in order to deal with computer operating systems and computational programs. One of most important and desired targets by VIP is navigating the Web, considering the increased importance and expansion of web-based computational programs. Screen readers present some solutions to navigate the textual and graphical contents of web pages, either by transforming a web page into a written Braille, or into a vocal output. In addition to these solutions, some screen readers installed on touch devices transform a web page into a vocal-tactile output.

But, there are some drawbacks for these proposed solutions. On the one hand, Braille techniques are costly and only few number of VIP have learned Braille (in France in 2011, there were about 77,000 visually impaired people and only 15,000 of them had learned Braille). On the other hand, transforming the information of a web page into a vocal format might not be suitable in public and noisy environments. Finally most of Braille solutions are not suitable for mobile devices [Maurel et al, 2012].

In addition to these drawbacks, the most important one is the failure to transform the 2-dimension web page structure. Indeed, as reported by many authors, perceiving the 2D structure of web documents greatly improves navigation efficiency and memorization as it allows high level text reading strategies such as: rapid or cursory reading, finding or locating information, to name but a few [Maurel et al, 2003].

Our work focuses on developing and evaluating a sensory substitution system based on a vibro-tactile solution, which may solve the mentioned drawbacks. In particular, we study how to increase the VIP perception of a 2D web page structure and how to enhance their techniques to navigate the contents of

This work is licensed under a Creative Commons Attribution 4.0 International Licence. License details: <http://creativecommons.org/licenses/by/4.0/>

web pages on touch-screen devices. The suggested solution is very cheap compared to Braille devices and may be efficient in noisy and public environments compared to vocal-tactile solutions.

Our contribution is three-fold: (1) designing a Tactile Vision Sensory System (TVSS) represented by an electronic circuit and an Android program in order to transform light contrasts of touch-screen devices into low-frequencies tactile vibrations; (2) designing an algorithm for segmenting web pages in order to support the visually impaired persons by a way which may enhance their ability to navigate the textual and graphical contents of web pages and (3) analyzing the effects of the suggested segmentation method on navigation models and tactics of blind persons, and its effect on enhancing their strategies for text reading and looking for textual information.

The paper is organized as follows. First, in section 2, we review most advanced VIP targeted technologies. Then, in section 3, we describe the new proposed framework. In section 4, we view the state of the art for web pages segmentation methods. In the fifth section, our hybrid segmentation method is presented and how this method could be integrated in our framework. In section 6, we enumerate the desired effects of the proposed segmentation method on navigation models and tactics of blind persons, and how it may enhance their strategies for text reading and searching of textual information. Finally, in the seventh section, perspectives and conclusions are presented.

2 VIP targeted technologies

Current products for VIP such as screen readers mainly depend on speech synthesis or Braille solutions, e.g. ChromeVox ^[3], Windows-Eyes ^[4], or JAWS (Job Access With Speech) ^[5]. Braille displays are complex and expensive electromechanical devices that connect to a computer and display Braille characters. Speech synthesis engines convert texts into artificial speech, where the text is analyzed and transformed into phonemes. These phonemes are then processed using signal processing techniques.

Some screen readers can also support tactile feedback when working on touch-screen devices, such as Mobile Accessibility ^[6] and Talkback ^[7] for Android, or VoiceOver ^[8] for iPad. Many of these products propose shortcuts for blind users to display a menu of HTML elements existing in the web page, for example headers, links and images. But, the main drawback of all these products is the fact that they transfer the information of web pages into a linear way i.e. without any indication of the 2-dimension global structure.

Many researches tried to enhance the way by which VIP interact with web pages, such as [Alaeldin et al, 2011], who proposed a tactile web navigator to enable blind people access the Web. This navigator extracts texts from web pages and sends them to a microcontroller responsible of displaying the text in Braille language using an array of solenoids. A tactile web browser for hypertext documents has been proposed by [Rotard et al, 2005]. This browser renders texts and graphics for VIP on a tactile graphics display and supports also a voice output to read textual paragraphs and to provide a vocal feedback. The authors implemented two exploration modes, one for bitmap graphics and another one for scalable vector graphics. A pin matrix device is then used to produce the output signal for blind users. The main drawback of these two proposed systems is that they need specific devices (solenoids and pin matrix), which are expensive and cannot be integrated to handled devices such as PDAs or Tablet PCs. Another interesting model called MAP-RDF (Model of Architecture of Web Pages) has been proposed by [Boulssa et al, 2011]. This model allows representing the structure of a web page and provides blind users with an overview of the web page layout and the document structure semantics. Tactos is a perceptual interaction system, which has been suggested by [Lenay et al, 2003] and consists of three elements: (1) tactile simulators (two Braille cells with 8 pins) represent a tactile feedback system, (2) a graphics tablet with a stylus represents an input device and (3) the computer. More than 30 prototypes of Tactos have been released to serve a lot of users in many domains. Tactos has been successfully used to recognize simple and complex shapes. The device has been also used in geometry teaching domain in an institution for visually impaired and blind children. Tactos also allowed psychology researchers to propose and develop new paradigms for studying perceptions and mediated communication of blind persons [Tixier et al, 2013]. However, it shows the same drawback as the previous systems, which are expensive and need specific devices. Moreover, the blind user can only explore the web page with a stylus and both hands are occupied by the system. Moreover, it is unemployable for a large set of environments, for example in public.

3 Proposed Framework

The “first glance” can be defined as the ability to understand the document layout and its structural semantics in a blink of an eye [Maurel et al, 2012]. In this work, we aim to increase the ability of visually impaired persons to understand the 2-dimension web page layout in order to enhance their tactics to navigate the Web with a vibro-tactile feedback.

The first phase in our model is to extract visual structures in the navigated web page and convert these “visual” blocks into zones (or segments) to facilitate the navigation in later phases. We achieve this phase depending on a hybrid segmentation method. Then the system represents the extracted visual elements as symbols using a graphical language. The third phase is to browse these graphical symbols depending on the size of the used touched-screen device; and in the fourth phase, our system provides a vibro-tactile feedback when the blind user touches the tablet by giving the user a vibro-tactile feedback by transforming light contrasts of touch-screen devices into low-frequencies tactile vibrations. A tablet (Asus Model TF101 with Android operating system) has being used for our tests.

To achieve the desired system, we have designed an electronic circuit, which controls two micro-vibrators placed on two fingers. A Bluetooth connection with an android tablet allows controlling the vibration intensity (i.e. amplitude) of vibrators. An Android dedicated program on the tablet views an image on the screen and detects where the user touches the tablet screen (the viewed image represents the result of web page segmentation). The intensity of the light emitted by the tablet at touched points is then transmitted to the embedded device in order to control the vibration intensity. In this paper, we focus only on the first phase (extracting visual structures in the navigated web page, and convert them into zones), with considering that detailed description of hardware components of the system, and results of pre-tests are described in [Maurel et al, 2012] and [Maurel et al, 2013].

4 Related Works

Segmenting a web page is a fundamental phase for understanding its global structure. Extracting the global structure of web pages is useful in many domains such as information retrieval, data extraction, and similarity of web pages.

Many approaches have been suggested for segmenting web pages, such as:

- 1) DOM-based segmentation: it depends on analyzing the DOM tree (Document Object Model), and extracting the main structure of web pages depending on HTML tags. An example of this approach is the work of [Sanoja et al, 2013], which determines firstly the layout template of a web page, and then it divides the page into minimum blocks, and finally collects these minimum blocks into content blocks.
- 2) Vision-based segmentation: this method divides the web page depending on the visual view of web page contents on a web browser. The most famous tool depends on this approach is VIPS (VIsion based Page Segmentation) [Deng et al, 2003].
- 3) Image processing based segmentation: this approach captures an image for the visual view of a web page, and then depends on image processing techniques to divide the captured image into sub blocks [Cai et al, 2004] [Cao et al, 2010].
- 4) Text-based Segmentation: this approach focuses on extracting only information about texts existed in a web page. After dividing the web page into blocks of texts, it could be possible to find the semantic relations between these textual blocks. This method is useful in many information retrieval domains such as question answering applications [Foucault e al, 2013].
- 5) Fixed-length segmentation: this approach divides the web pages into fixed length blocks (passages), after removing all HTML tags, where each passage contains a fixed number of words [Callan, 1994].
- 6) Densitometric analysis based segmentation: this approach depends on methods applied in quantitative linguistics, where text-density refers to a measure for identifying important textual segments of a web page [Kohlschütter et al, 2008].
- 7) Graph-based segmentation: This approach depends on transforming the visual segments of a web page into graph nodes, then applying many common graph methods on these nodes for combining

them into blocks, or for making a clustering for these nodes. Some common works which depend on this approach are [Chakrabarti et al, 2008] [Liu et al, 2011].

-8) and Hybrid-based segmentation: This approach combines many approaches indicated previously.

5 Suggested Hybrid Segmentation Algorithm

Most of segmentation algorithms render firstly the web page using a web browser, and then segments the HTML elements into many blocks depending on the visual layout. The constructed hybrid segmentation algorithm has been tested on 154 pages collected manually from many newspaper and e-commerce sites (www.leparisien.fr, www.lefigaro.fr, www.liberation.fr, www.amazon.fr, www.materiel.net), and the results have been integrated with our under-development Android program. The obtained results are promised because the segmentation algorithm can extract well the web page blocks depending on the visual structure, and the algorithm can also convert correctly these blocks into zones (clustering the blocks). Our algorithm blends three segmentation approaches, DOM-based segmentation, vision-based segmentation, and graph-based segmentation.

Proposed Corpora

To achieve the previous mentioned model, we construct two corpora, one for training, and another for testing. We selected many criteria for crawling web pages, such as, the type of crawled pages (information web sites, and e-commerce web sites), the size (about 10,000 pages), the language (French), the version of web site (Classic, Mobile), and the technology used to build the crawled web site (framework JavaScript: JQuery, mootools, ... / CMS: Prestashop, Drupal, Joomla...).

5.1 Vision-Based Approach

In this phase, we render the web page using Mozilla FireFox browser, and getting its visual structure by injection Java-script code inside the HTML source code of the rendered web page. The obtained visual structure indicates a global hierarchy of the rendered web page, and assigns a bounding box for each HTML element. Figure 1.a represents a part of a web page, and the result of its vision-based segmentation is presented in figure 1.b.



(a) A part of a web page (b) Vision-based segmentation
Figure 1. A part of a web page (leparisien.fr) and its vision-based segmentation

The input of this phase is a web page HTML source code, and its output is injected information about bounding boxes for each HTML element. In next sections, we refer to bounding boxes by blocks (i.e. each bounding box represents an HTML element, and may contain other bounding boxes.).

5.2 DOM-Based Approach

After segmenting a web page depending on its visual structure, we analyze its DOM structure by applying filters and re-organization rules for enhancing results of next phases. Dead-Nodes filter is an example of these filters: it deletes all HTML nodes that do not affect on the appearance, for example nodes with height or width equals to "0px" (zero pixel); or nodes with style properties ("display :

none" or "visibility:hidden"). An example of re-organization rules is Paragraph-Reorganization rule, where this rule re-constructs all paragraph child-nodes in one node contains the extracted text; we made this rule after analyzing many DOM structures, and observing that some paragraph nodes contain child-nodes which affect negatively on extracting the text, such as <i>, , etc..., and these child-nodes contain important texts. We made many filters and re-organization rules, and integrated them with our framework, and then we tested applying these rules and filters on the vision-based segmented web pages (154 pages mentioned previously). As a result of applying the two approaches (vision-based and DOM-based), we succeeded to get the first glance visual structure for many pages.

The result of this phase is a filtered DOM-tree, each of its nodes is visible and contains a bounding box information. Figure 1.b represents a hierarchy of some HTML nodes, the first level contains 3 main blocks (B1, B2, and B3), and each one contains many sub blocks, for example B3 contains B3-1 and B3-2.

To illustrate results of applying previous two mentioned approaches, we represented an obtained filtered DOM-tree on the used tablet. Figure 2 views a graphical representation for a page web, since each rectangle represents a block in the analyzed web page. Red rectangles represent images (tags), green rectangles represent links (<a> tags), blue rectangles represent list of items (or tags), and finally, black rectangles represent paragraphs (<p> tags).



Figure 2. A graphical representation for a page web (leparisien.fr)

5.3 Graph-Based Approach

After segmenting the web page depending on its visual structures and analyzing its DOM-structure, we apply a new graph-based segmentation algorithm which called “Blocks2Zones Clustering” in order to group many similar blocks together in one zone. Clustering many blocks together is necessary in order to decrease the number of viewed blocks in some interfaces (instead of viewing many blocks, we view one zone represents these blocks and then the user can navigate intra-elements inside the zone by double clicking on the graphical element of the chosen zone.), and to group closed blocks in one zone (here, closeness depends on distances between blocks, this will be described next sections in details). The pseudo-code of the proposed algorithm is:

Blocks2Zones Clustering Algorithm

Input (Blocks, N° of desired Zones)

Output: Graph of N nodes (N Zones)

1- Transform the blocks into a graph (Non-Directed graph)

1.1. Blocks \rightarrow Nodes,

1.2. Make relations between the nodes, and assign weights for these relations.

2- If number of zones \geq number of blocks

end the algorithm,

Else

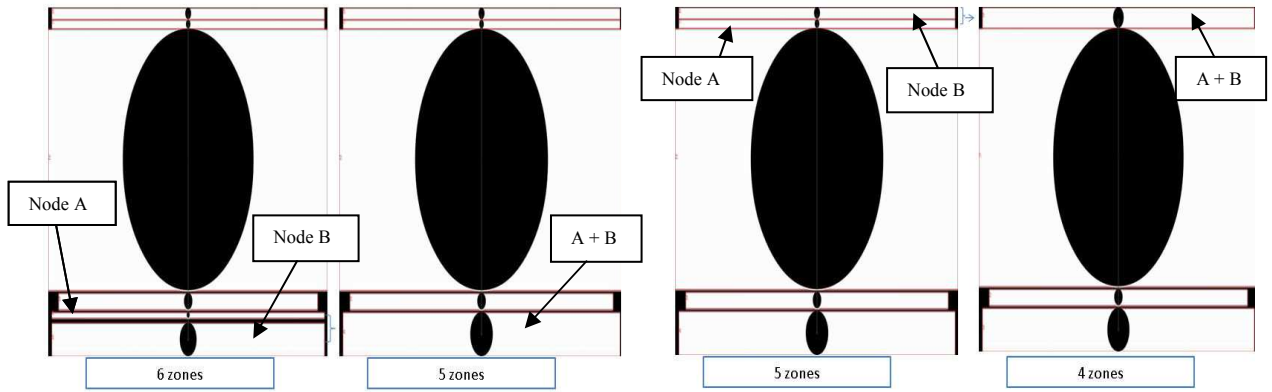
3- Find the node with the smallest size (node A) (Figure 3.a (6 zones), Figure 3.b (5 zones))

4- For node A, find the relation which has the largest weight (node B). (Figure 3.a (6 zones), Figure 3.b (5 zones))

5- Group the nodes A, and B (A+B). (Figure 3.a (5 zones), Figure 3.b (4 zones))

6- Repeat steps 3-4-5 till number of blocks == number of zones

Figure 3 represents some examples of applying this algorithm, where each rectangle represents a zone (a block or a collection of blocks), and the center of each ellipse represents the zone center.



(a) Converting 6 zones to 5 zones

(b) Converting 5 zones to 4 zones

Figure 3. Examples of applying Blocks2Zones clustering Algorithm

To calculate weights between nodes, we tested 2 relations of distances: the first one is Minkowski Manhattan distance ($d(p, q) = d(q, p) = ||p - q|| = \sum_{i=1}^n |p_i - q_i|$), and the second is Minkowski Euclidian distance ($d(p, q) = d(q, p) = \sqrt{(p_1 - q_1)^2 + (p_2 - q_2)^2 + \dots + (p_n - q_n)^2}$). To ensure which distance should be used, we applied an internal quality criterion for the two used distances; the applied criterion is Sum of Squared Error (SSE) ($SSE = \sum_{K=1}^n \sum_{x_i \in C_K} ||x_i - \alpha_i||^2$ Where C_K is the set of instances in cluster K , and $\alpha_{k,j} = \frac{1}{N_k} \sum_{x_i \in C_K} x_{i,j}$).

Results of applying SSE measure on the two distances (Minkowski Manhattan, and Minkowski Euclidian) were the same; which means that using either Minkowski Manhattan distance or Minkowski Euclidian distance is equal in our algorithm to calculate weights between nodes.

Applying this hybrid segmentation algorithm on a filtered DOM-tree (obtained from applying Vision-based approach then Dom-based approach) converts a web page to a set of zones, each zone contains many other zones or blocks, and each block represents a visual structure and may contain many other blocks. The purpose of the proposed vibro-tactile access protocol is then to transform the semantic of symbols in these zones, or blocks, or HTML elements into vibrations with different frequencies and amplitudes.

6 Desired Effects of Suggested Algorithm on Web Navigation Models of VIP

We had finished designing the suggested algorithm, and integrated it with our framework. However, practical experiments with VIP are the most important criteria to select success or failing this algorithm in enhancing VIP web navigation models, and this is our next step to be achieved. But, depending on our previous experiments in dealing with VIP targeted techniques, we can expect some desired effects of the proposed segmentation method on navigation models and tactics of blind persons, and how it may enhance their strategies for text reading and searching of textual information.

Firstly, this segmentation model can give VIP an impression of the layout of navigated web page (first glance layout), and (as indicated previously) perceiving the 2D structure of web documents greatly improves navigation efficiency and memorization.

Secondly, this suggested segmentation method can group together many closed blocks in one zone, and in this way the user can select easily if these zone contents are important for him/her or not, for example collecting all header blocks in one zone, or collecting all footer blocks in one zone.

Thirdly, this model can allow high level text reading strategies such as rapid or cursory reading or locating information, this can be achieved by ignoring navigating any zone does not contain textual information. By the way if one zone contains textual information, the user can navigate it and decide if it contains important information for him/her or not.

7 CONCLUSION AND PERSPECTIVES

In this paper, we summarized our current work which aims to design an approach for non-visual access to web pages on touch-screen devices, and we focused on the suggested hybrid segmentation algorithm. We expect that integrating this method of segmentation with the designed vibro-tactile protocol can give VIP an impression of the first glance layout of web pages.

In the same way that the environment enables a blind person to move in space with sidewalks and textures which will be explored by his/her white cane, we hope giving the blind user an ability to navigate documents depending on "textual sidewalks" and "graphical paths" which will be discovered by his/her finger.

Next steps in this research will be 1) making real experiments to study effects of suggested segmentation algorithm on VIP web navigation models, 2) adding advances techniques in text summarization to facilitate navigating textual information, 3) adding elements to the graphical vibro-tactile language in order to represent more HTML elements such links, buttons, input fields, and other elements, 4) we plan also to add thermic actuators for translating the notion of colors. This may be very useful and hopeful for blind users to transfer information about colors.

8 References

- [1] Maurel, F., Dias, G., Routoure, J-M., Vautier, M., Beust, P., Molina, M., Sann, C., « *Haptic Perception of Document Structure for Visually Impaired People on Handled Devices* », *Procedia Computer Science*, Volume 14, Pages 319-329, ISSN : 1877-0509, 2012.
DOI=<http://dx.doi.org/10.1016/j.procs.2012.10.036>
- [2] Maurel, F., Vigouroux, N., Raynal, M., Oriola, B., « *Contribution of the Transmodality Concept to Improve Web Accessibility* ». In *Assistive Technology Research Series*, Volume 12, 2003, Pages 186-193. International conference; 1st, Smart homes and health telematics; Independent living for persons with disabilities and elderly people. ISSN : 1383-813X. 2003.
- [3] <http://www.chromevox.com/> [Access 24/5/2014]
- [4] <http://www.synapseadaptive.com/gw/wineyes.htm> [Access 24/5/2014]
- [5] <http://www.freedomscientific.com/> [Access 24/5/2014]
- [6] <https://play.google.com/store/apps/details?id=es.codefactory.android.app.ma.vocalizerfrfdemo&hl=fr> [Access 24/5/2014]
- [7] <https://play.google.com/store/apps/details?id=com.google.android.marvin.talkback&hl=fr> [Access 24/5/2014]
- [8] <http://www.apple.com/fr/accessibility/> [Access 24/5/2014]
- [9] Alaidin, A., Mustafa, Y., Sharief, B., 2012. « *Tactile WebNavigator Device for Blind and Visually Impaired People* ». In *Proceedings of the 2011 Jordan Conference on Applied Electrical Engineering and Computing Technologies*, Jordan, 2011.
DOI=<http://dx.doi.org/10.1109/AEECT.2011.6132519>
- [10] Rotard, M., Knödler, S., Ertl, T., « *A Tactile Web Browser for the Visually Disabled* ». In *Proceedings of the sixteenth ACM Conference on Hypertext and Hypermedia*. ACM, New York, NY, USA, 2005, pages 15-22, 2005.
DOI= <http://dx.doi.org/10.1145/1083356.1083361>
- [11] Boulssa, Y., Mojahid, M., Oriola, B., Vigouroux, N., « *Accessibility for the Blind, an Automated Audio/Tactile Description of Pictures in Digital Documents* ». *IEEE International Conference on Advances in Computational Tools for Engineering Applications*, 2009, Pages: 591 – 594, 2009.
DOI=<http://dx.doi.org/10.1109/ACTEA.2009.5227855>

- [12] Lenay, C., Gapenne, O., Hanne-ton, S., Marque, C., Genouëlle, C., “*Sensory Substitution, Limits and Perspectives*». In *Touch for Knowing Cognitive psychology of haptic manual perception*, Amsterdam, Pages: 275-292, 2003,
- [13] Tixier, M., Lenay, C., Le-Bihan, G., Gapenne, O., Aubert, D., « *Designing Interactive Content with Blind Users for a Perceptual Supplementation System* », TEI 2013, 2013, in *Proceedings of the 7th International Conference on Tangible, Embedded and Embodied Interaction*, Barcelona, Spain, Pages 229-236, 2013.
DOI= <http://dx.doi.org/10.1145/2460625.2460663>
- [14] Maurel, F., Safi, W., Beust, P., Routoure, J.M., « *Navigation aveugle sur dispositifs mobiles : toucher le Web... pour mieux l'entendre*», 16ème Colloque International sur le Document Électronique, CIDE16, Lille, France, Europa productions, 2013.
- [15] Sanoja, A., Gançarski, S., «*Block-o-Matic: a Web Page Segmentation Tool*», BDA. Nantes, France. 2013. <http://hal.archives-ouvertes.fr/hal-00881693/>
- [16] Deng, C., Shipeng, Y., Ji-Rong, W., Wei-Ying, M., «*VIPS: a Vision-based Page Segmentation Algorithm*», Nov. 1, 2003, Technical Report MSR-TR-2003-79, Microsoft Research. 2003.
<http://research.microsoft.com/pubs/70027/tr-2003-79.pdf>
- [17] Cai, D., He, X., Ma, W-Y., Wen, J-R., Zhang, H., 2004. « *Organizing WWW Images Based on the Analysis Of Page Layout And Web Link Structure*», Microsoft Research Asia, Beijing, China, 2004.
<http://research.microsoft.com/pubs/69080/25.pdf>
- [18] Cao, J., Mao, B., Luo, J., « *A segmentation method for web page analysis using shrinking and dividing*», *International Journal of Parallel, Emergent and Distributed Systems - Network and parallel computing*, Volume 25 Issue 2, April 2010. Pages: 93-104, 2010.
DOI=<http://dx.doi.org/10.1080/17445760802429585>
- [19] Foucault, N., Rosset, S., Adda, G., « *Pré-segmentation de pages web et sélection de documents pertinent en Questions-Réponses*», TALN-RÉCITAL 2013.
- [20] Callan, J.P., 1994. « *Passage- level Evidence in Document Retrieval*», the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Dublin, Pages: 302-310.
Publisher: Springer-Verlag New York, Inc, 1994.
<http://dl.acm.org/citation.cfm?id=188589>
- [21] Kohlschütter, C., Nejd, W., « *A Densitometric Approach to Web Page Segmentation*», USA, CIKM'08, *Proceedings of the 17th ACM conference on Information and knowledge management*, 2008. Pages: 1173-1182.
DOI: <http://dx.doi.org/10.1145/1458082.1458237>
- [22] Chakrabarti, D., Kumar, R., Punera, K., 2008. « *A graph-theoretic approach to webpage segmentation*». *Proceedings of the 17th international conference on World Wide Web, WWW'08*, ACM, USA, 2008. Pages: 377-386. Publisher: ACM New York, NY, USA, 2008.
DOI: <http://dx.doi.org/10.1145/1367497.1367549>
- [23] Liu, X., Lin, H., Tian, Y., 2011. “*Segmenting Webpage with Gomory-Hu Tree Based Clustering*”, *Journal of Software*, Vol 6, No 12, Pages: 2421-2425. 2011.