

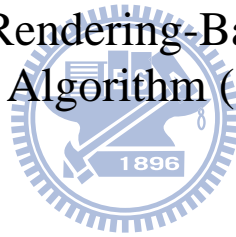
國立交通大學

資訊學院 資訊學程

碩士論文

基於HTML文件佈局之網頁分割演算法

A HTML Rendering-Based Page Segmentation
Algorithm (HRPS)



研究生：余提梵

指導教授：吳毅成 博士

中華民國九十九年六月

基於HTML文件佈局之網頁分割演算法

A HTML Rendering-Based Page Segmentation Algorithm (HRPS)

研 究 生：余提梵

Student: Ti-fan Yu

指導教授：吳毅成 博士

Advisor: Dr. I-Chen Wu

國立交通大學



A Thesis

Submitted to College of Computer Science

National Chiao-Tung University

in partial Fulfillment of the Requirements

for the Degree of Master of Science

in

Computer Science

June 2010

Hsinchu, Taiwan, Republic of China

中華民國九十九年六月

基於 HTML 文件佈局之網頁分割演算法

學生：余提梵

指導教授：吳毅成 博士

國立交通大學 資訊學院 資訊學程碩士班

摘 要

依據統計資料，截至 2010 年來共有 1.13 億個網站存在過，其中有 99.9% 是在近 15 年間成立的，面對這樣龐大又高替換的網頁資料，如何有效地使用是一件很重要的事。對於大量變動態的資料，通常尋求搜尋引擎的協助來正確定位資料；對於已知位址的資料，為了增加使用效率，則會使用資料萃取的技術。而不管是搜尋引擎或資料萃取工具，要對複雜的網頁進行分析，首要就是要對網頁作區塊分類或標記，以濾除噪音 (Noise) 區塊及提取各主題 (Topic) 區域之本文區塊，也就是網頁區塊分割 (Page Segmentation) 。

2003 年微軟團隊發表視覺化網頁分割演算法 (Vision-based page segmentation: VIPS) 後，很多網頁分割研究多參考了視覺化分割技術。但在近幾年來，越來越多網頁的頁面框架設計，採用 DHTML 技術為主時，原始的 VIPS 的方法在使用上，便出現當初設計時沒有顧及的小缺陷，雖然之後的研究，出現很多組合型態的頁面分割演算法來彌補使用上的不足。但因為是採用其它特性的演算法來彌補 VIPS，所以這部份切割區塊也就喪失視覺化分割的特性。

本文提出一個方法，在以視覺化分割為基礎上，帶入網頁文件佈局特性 (HTML Rendering-Based)，以解決視覺化區塊分割在 DHTML 網頁上，可能找不到視覺化分割線 (Separator) 的問題。

A HTML Rendering-Based Page Segmentation Algorithm (HRPS)

Student: Ti-fan Yu

Advisor: Dr. I-Chen Wu

Degree Program of Computer Science National Chiao Tung University

ABSTRACT

According to the statistical datas, Up to 2010, a total of 113 million websites existed , of which 99.9% was established nearly 15 years, the face of such large and high replacement page data , how to effectively use is a very important matter . For the information that we don' t know its location , we usually use search engine to help us to find it out . And for the information that we do know where it is , we use data extraction to increase the efficiency . And whether it is a search engine or information extraction tool , to analyze the complex web , the first steps is to split the Web Page to provide subject area of this location , It' s a important thing that how to use this huge database efficiently .

Since 2003 the team released Microsoft Visual Web segmentation algorithm (Vision-based page segmentation: VIPS) , many papers are mostly used segmentation based on visual segmentation , However, in recent years , more and more web page Layout design , using DHTML technology-based, the original method of VIPS in the use , they are in the original design did not take into account small defects , though after the study , there are many page segmentation algorithm combined patterns to make up for the use of deficiency .

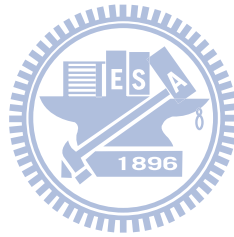
But since they are using other features of the algorithm to make up for VIPS, so this part of the Visual cues is losing the characteristics of visual segmentation , This paper presents a method , in order to split based on visualization , into the HTML document Rendering features , to solve the visual segmentation in DHTML pages , you may not find the visual Separator problems .

誌 謝

首先要感謝我的指導教授，吳毅成博士，由於他的細心指導與協助，本論文才得以順利完成。

另外也要感謝陳隆彬學長、孫德中學長，總在每次的討論之中，給我正確的方向。再來就是實驗室BODE組的各位伙伴，有你們的加油打氣，我才能一路堅持至此，謝謝你們。

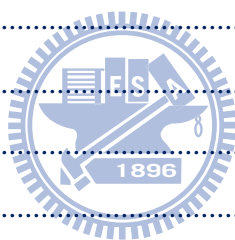
最後要感謝的就是我的爸媽，為了這篇論文，犧牲了很多陪伴他們的時間。他們不曾抱怨、始終要我以學業為重，並在我低潮時給我最大的支持與鼓勵，讓我可以安心的做研究、完成學業。謹以此論文，獻給我最愛的爸媽。



目 錄

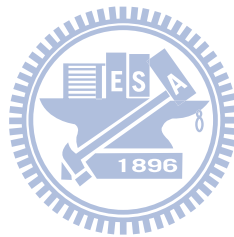
中文提要.....	i
英文提要.....	ii
誌謝.....	iii
目錄.....	iv
表目錄.....	vi
圖目錄.....	vii
符號說明.....	ix
一、 緒論.....	1
1.1 研究背景及動機.....	1
1.2 論文內容概述及大綱.....	2
二、 研究內容與方法.....	3
2.1 網頁區塊分割.....	3
2.2 相關研究.....	3
2.3 視覺化網頁區塊分割.....	5
2.4 VIPS的區塊分割提早停止問題.....	6
2.5 區塊分割提早停止現象的分析.....	7
2.5.1 代表視覺區塊之結點的選擇位置.....	7
2.5.2 隱藏的結點.....	9
2.5.3 重疊結點的位置.....	10
2.5.4 瀏覽器本身的佈局問題.....	11
2.5.5 代表資料的DOM-TREE出現在不只一處.....	12
2.5.6 由小變大的區塊.....	13
2.5.7 使用重疊設計的區塊.....	13
2.6 DHTML頁面設計的效應總結.....	14
三、 理論.....	15
3.1 改良的區塊擷取方法.....	15

3.1.1	DOM-TREE的巡行問題.....	15
3.1.2	視覺區塊的選擇方法.....	15
3.1.3	OFFSET DOM.....	16
3.2	遮蔽區塊調整策略.....	20
3.2.1	DHTML文件布局.....	20
3.2.2	遮蔽區塊的分裂調整.....	20
3.2.3	資訊測量.....	22
3.2.4	遮蔽區塊調整.....	22
3.3	HRPS 程序流程圖.....	24
3.4	分割塊大小調整.....	24
四、	實驗部分.....	26
4.1	系統實作概述.....	26
4.1.1	開發環境.....	26
4.1.2	模組介紹.....	26
4.2	實驗.....	29
4.2.1	範例.....	29
4.2.2	資料來源.....	31
4.2.3	實驗結果.....	34
五、	結論.....	37
5.1	總結.....	37
5.2	相關應用.....	37
5.3	未來工作.....	37
	參考文獻.....	39



表目錄

表 4.1	Data Mining的測試來源表	31
表 4.2	VIPS與HRPS測試結果評估表	36



圖目錄

圖 2.1	Example of FOM.....	4
圖 2.2	VIPS Hierarchical Layout	4
圖 2.3	VIPS演算法示意圖	6
圖 2.4	簡單的頁面分隔線偵測流程圖	6
圖 2.5	VIPS Demo Program.....	6
圖 2.6	Entropy Reduction.....	7
圖 2.7	Empty Separator	7
圖 2.8	Zone of <Table> Tag.....	8
圖 2.9	Zone of Tag.....	8
圖 2.10	Partial DOM Tree of figure 2.8,2.9	9
圖 2.11	隱藏的結點.....	9
圖 2.12	Partial DOM Tree of figure 2.10.....	10
圖 2.13	重疊的結點.....	10
圖 2.14	Partial DOM Tree of figure 2.12.....	11
圖 2.15	Browser Rendering Problem	11
圖 2.16	Repeater Elements.....	12
圖 2.17	Partial DOM Tree of figure 2.15.....	12
圖 2.18	由小變大的區塊.....	13
圖 2.19	使用重疊設計的區塊.....	14
圖 2.20	Partial DOM Tree of 大英博物館	14
圖 3.1	Inline Elements.....	16
圖 3.2	Block Elements	16
圖 3.3	Block Elements as Carriage return.....	16
圖 3.4	HTML/XML Application Program Interface.....	17

圖 3.5	Element Position	18
圖 3.6	Partial Graph of figure 2.15.....	18
圖 3.7	HTML DOM	19
圖 3.8	Offset DOM	19
圖 3.9	Hierarchical layout	20
圖 3.10	(1).No Effect (2).Overlay	21
圖 3.11	Included Block	21
圖 3.12	Intersection Block	21
圖 3.13	Measuring Strategy.....	23
圖 3.14	Example of Segmentation	23
圖 3.15	HRPS Process Flow	24
圖 3.16	Data Regions Adjustment	25
圖 4.1	Command UI.....	26
圖 4.2	WebCrawler.exe.....	27
圖 4.3	WebSegment.exe.....	28
圖 4.4	WebAnalyer.exe.....	28
圖 4.5	WebVerify.exe	29
圖 4.6	Figure 2.18的區塊分割	29
圖 4.7	Figure 2.12的區塊分割	30
圖 4.8	Figure 2.15的區塊分割	30
圖 4.9	VIPS 區塊分割結果	35
圖 4.10	HRPS 區塊分割結果.....	35



符 號 說 明

DoC : Degree of Coherence

DoI : Degree of Information

