

Ekstraksi Data pada Tabel dari Halaman Web Menggunakan Pohon *Document Object Model*

Memem Akbar¹, Cici Patmala², Dini Nurmalasari³

Abstract— Data on the web page can be available in various formats, such as table. With the growing of web pages, the need to extract data from tables is increasing. Results of the extraction can be used for integration with other web tables or stored in a database. This study discusses the extraction of data from a table on a web page using a Document Object Model (DOM) tree. The initial step of this extraction process is to transform the HTML document into a DOM tree. Then, by applying search methods Depth First Search (DFS), part of the data in the table is extracted and stored in a CSV file. An engine has been developed using Visual Basic. The results show that the engine can automatically extract data from the table that has the following characteristics: the number of rows and columns are not limited, able to handle all of the table orientation layout, and able to handle tables that are merged cells.

Intisari— Data pada halaman web dapat tersedia dalam berbagai format, salah satunya adalah tabel. Pada tabel, informasi ditampilkan dalam baris dan kolom. Tabel pada halaman web ditampilkan dengan menggunakan tag HTML. Dengan semakin banyaknya halaman web, kebutuhan untuk mengekstrak data dari tabel tersebut semakin meningkat. Hasil ekstraksi dapat digunakan untuk integrasi dengan tabel web lainnya atau untuk disimpan dalam suatu basis data. Makalah ini membahas ekstraksi data dari tabel pada halaman web dengan menggunakan pohon *Document Object Model* (DOM). Langkah awal proses ekstraksi ini adalah mentransformasi dokumen HTML menjadi pohon DOM. Kemudian, dengan menerapkan metode penelusuran *Depth First Search* (DFS), bagian data pada tabel diekstrak dan disimpan ke dalam file CSV. Sebuah aplikasi telah dikembangkan dengan bahasa Visual Basic. Hasil pengujian terhadap sejumlah halaman web menunjukkan bahwa aplikasi secara otomatis dapat mengekstrak data dari tabel yang memiliki karakteristik: jumlah baris dan kolom yang tidak terbatas, semua orientasi layout tabel, serta dapat menangani tabel yang terdapat *merged cell*.

Kata Kunci— Ekstraksi tabel, DOM, Website, HTML, DFS

I. PENDAHULUAN

Saat ini, halaman web merupakan sumber informasi dengan volume yang sangat besar. Sebuah halaman web dapat terdiri atas teks, gambar, daftar (*list*), dan tabel. Tabel merupakan salah satu cara yang digunakan untuk menampilkan data dalam bentuk baris dan kolom. Tabel menampilkan ringkasan data dari isi halaman web. Ekstraksi data pada tabel menjadi penting untuk dapat melihat keseluruhan isi halaman web.

Selain itu, ekstraksi tabel juga penting sebagai langkah awal dalam proses integrasi data tabel dari beberapa halaman web [1].

Tabel pada halaman web ditampilkan dengan menggunakan *tag-tag* Hypertext Markup Language (HTML). Oleh karena itu, untuk melakukan ekstraksi tabel dari halaman HTML dibutuhkan langkah-langkah yang dapat mempelajari dan menelusuri *tag-tag* penyusun sebuah halaman web, khususnya tabel. Makalah ini membahas tentang ekstraksi tabel pada halaman web dengan menelusuri struktur tag HTML yang membangun sebuah tabel.

Pada makalah ini, tag penyusun sebuah tabel terlebih dahulu dimodelkan dalam pohon *Document Object Model* (DOM). Pada struktur pohon DOM ini, bagian data terletak pada daun paling ujung dari setiap ranting. Struktur data dalam bentuk pohon DOM ini kemudian ditelusuri untuk mendapatkan bagian data pada tabel. Penelusuran pohon yang digunakan adalah *Depth First Search* (DFS). Hasil ekstraksi kemudian disimpan dalam file berbentuk csv.

Pada bagian selanjutnya, diuraikan beberapa penelitian terdahulu yang membahas mengenai ekstraksi data dari tabel pada halaman web. Bagian III membahas perancangan aplikasi untuk menerapkan teknik yang digunakan. Kemudian, bagian IV membahas pengujian aplikasi terhadap beberapa tabel dari sejumlah halaman web.

II. EKSTRAKSI DATA DARI HALAMAN WEB

Penelitian mengenai ekstraksi data dari halaman web telah banyak dilakukan. Sebuah penelitian melakukan ekstraksi data yang ditampilkan dalam bentuk list dan tabel pada sebuah halaman web [2]. Penelitian ini menerapkan *unsupervised learning algorithms* untuk memperoleh struktur dari list dan tabel yang ada di web dengan mengacu pada format halaman dan data yang ada di dalamnya. Dalam penelitian ini, dibutuhkan analisis dari beberapa halaman web sebelum akhirnya dapat dilakukan ekstraksi data menjadi kolom dan baris. Padahal bisa saja pada halaman web hanya ada satu halaman, dan dari satu halaman tersebut terdapat list dan tabel yang butuh diekstrak.

Penelitian yang lain mengembangkan metode untuk deteksi dan ekstraksi tabel dengan melakukan analisis *graphic lines* [3]. Tabel sederhana terdiri atas sel-sel matriks. Semua sel dalam satu baris memiliki tinggi yang sama dan semua sel dalam satu kolom memiliki lebar yang sama. Semua sel tersebut dibatasi oleh *graphic lines*. Penelitian ini dapat menjadi salah satu acuan dalam mengembangkan metode ekstraksi tabel dari file dalam format PDF.

Telah diteliti pula ekstraksi data pada tabel HTML dengan lima jenis layout [4]. Tabel pada halaman web tidak memiliki layout standar seperti tabel pada basis data. Layout tabel pada

^{1,3} Program Studi Teknik Komputer Politeknik Caltex Riau, Jl. Umbansari 1 Rumbai-Pekanbaru 28265, INDONESIA (email: memem@pcr.ac.id, dini@pcr.ac.id)

² Program Studi Teknik Informatika Politeknik Caltex Riau, Jl. Umbansari 1 Rumbai-Pekanbaru 28265, INDONESIA (email: cici12ti@mahasiswa.pcr.ac.id)

halaman *web* sangat bergantung pada kebutuhan dan selera *programmer*. Penelitian ini berkontribusi besar dalam ekstraksi kelima jenis *layout* yang ditentukan.

Pada penelitian lain, ekstraksi merupakan bagian dari proses integrasi data dari beberapa halaman *web* [1], [5]-[8]. Dalam proses ekstraksi, penelitian-penelitian ini menggunakan pendekatan pengenalan *tag* HTML pembentuk sebuah tabel. Karena kebutuhan hasil ekstraksi adalah untuk integrasi data dengan tabel lain, proses ekstraksi dilakukan hingga dapat memisahkan bagian atribut dan bagian data dari tabel. Hasil ekstraksi diubah menjadi pohon semantik [1] dan bentuk XML [5]-[8].

Penelitian lain menerapkan DOM dalam proses ekstraksi konten dari sebuah halaman *web* [9]. Penelitian ini berhasil menghilangkan iklan, gambar, dan *link* yang tidak berhubungan dengan isi *web*.

Telah diterapkan pula algoritme berbasis aturan untuk mengekstrak informasi dari laporan audit pada pernyataan keuangan pemerintahan lokal [10].

Dibandingkan dengan penelitian lain, pada makalah ini dilakukan ekstraksi tabel dari halaman *web* menggunakan metode yang berbeda, yakni menggunakan DOM dalam proses ekstraksi. Pohon DOM merupakan pohon sintaktik untuk melihat struktur dokumen HTML. Penerapan pohon DOM juga memungkinkan ekstraksi untuk kelima jenis orientasi *layout* tabel [4]. Keluaran hasil ekstraksi akan disimpan dalam *file* csv. Bentuk ini dapat digunakan untuk proses berikutnya, untuk disimpan ke basis data maupun diintegrasikan dengan tabel lainnya.

III. HTML, DOM, DAN DFS

A. Hypertext Markup Language (HTML)

HTML adalah sekumpulan simbol-simbol atau *tag-tag* yang dituliskan dalam sebuah *file* yang digunakan untuk menampilkan alamat pada *web browser* [11]. *Tag-tag* HTML selalu diawali dengan `<x>` dan diakhiri dengan `</x>`. Setiap dokumen HTML harus diawali dan ditutup dengan *tag* HTML `<HTML>` dan `</HTML>`. *Tag* HTML memberitahu *browser* bahwa yang ada di dalam kedua *tag* tersebut adalah dokumen HTML.

Tabel digunakan untuk menyajikan data dalam bentuk kolom dan baris. Umumnya setiap kolom menunjukkan data yang sejenis, dan setiap baris menunjukkan kelompok data dalam satu kesatuan. Beberapa *tag* HTML dalam pembuatan tabel beserta fungsinya dapat dilihat pada Tabel I [11].

B. Document Object Model (DOM)

DOM merupakan sebuah ketentuan yang dikembangkan oleh W3C untuk berinteraksi dengan objek-objek yang ada di dalam HTML, XML, maupun XHTML. DOM bersifat *cross-platform* dan *language-independent*. DOM dapat digunakan dengan bahasa pemrograman apapun dan pada sistem operasi manapun. Pengembang bahasa pemrograman atau sistem operasi harus mengimplementasikan antarmuka DOM terlebih dahulu sebelum dapat digunakan pada aplikasi.

TABEL I
DAFTAR REFERENSI TAG HTML

Nama Tag	Keterangan/Kegunaan
<code><table>...</table></code>	Tag untuk membuat tabel
<code><th>...</th></code>	Tag untuk membuat sebuah sel <i>header</i> tabel
<code><tr>...</tr></code>	Tag untuk membuat baris dalam sebuah tabel
<code><td>...</td></code>	Tag untuk membuat sel dalam sebuah tabel
<code><thead>...</thead></code>	Mengelompokkan isi <i>header</i> dalam sebuah tabel
<code><tbody>...</tbody></code>	Mengelompokkan isi tubuh dalam sebuah tabel
<code><tfoot>...</tfoot></code>	Mengelompokkan isi <i>footer</i> dalam sebuah tabel

Sebuah dokumen HTML direpresentasikan oleh DOM dalam bentuk struktur hirarki pohon. Dengan kata lain, cara DOM melihat dokumen HTML seperti struktur hirarki pohon.

C. Depth First Search (DFS)

DFS adalah salah satu algoritme penelusuran struktur graf/pohon berdasarkan kedalaman. Simpul ditelusuri dari *root* kemudian ke salah satu simpul anaknya (misalnya prioritas penelusuran berdasarkan anak pertama [simpul sebelah kiri]), maka penelusuran dilakukan terus melalui simpul anak pertama dari simpul anak pertama level sebelumnya hingga mencapai level terdalam [12].

Setelah sampai di level terdalam, penelusuran akan kembali ke satu level sebelumnya untuk menelusuri simpul anak kedua pada pohon biner [simpul sebelah kanan] lalu kembali ke langkah sebelumnya dengan menelusuri simpul anak pertama lagi sampai level terdalam dan seterusnya. Dalam pencarian menggunakan algoritme DFS, simpul-simpul yang paling dalam pada pohon yang akan dicari paling awal, simpul yang telah dikunjungi disimpan dalam suatu tumpukan (*stack*).

Algoritme penelusuran pohon menggunakan DFS dapat dijelaskan sebagai berikut.

1. Memasukkan simpul ujung (akar) ke dalam tumpukan.
2. Mengambil simpul dari tumpukan teratas, lalu memeriksa apakah simpul merupakan solusi
3. Jika simpul merupakan solusi, pencarian selesai dan hasil dikembalikan.
4. Jika simpul bukan solusi, seluruh simpul yang bertetangga dengan simpul tersebut (simpul anak) dimasukkan ke dalam tumpukan.
5. Jika tumpukan kosong dan setiap simpul sudah diperiksa, pencarian selesai dan hasil solusi tidak ditemukan dikembalikan.
6. Mengulangi pencarian dari langkah kedua.

IV. PERANCANGAN DAN ANALISIS

Bagian ini menjabarkan perancangan yang digunakan untuk membuat aplikasi ekstraksi data tabel dari halaman *web*. Terdapat dua perancangan yang digunakan, yaitu diagram blok dan diagram *use case*.

A. Diagram Blok

Gbr. 1 merupakan diagram blok proses ekstraksi data dari sebuah tabel pada halaman *web*. Yang menjadi *input* pada proses ini adalah alamat *URL* atau kode HTML halaman *web* yang akan diekstraksi. Langkah awal yang dilakukan adalah mengubah dokumen HTML dari halaman *web* ke dalam bentuk teks. Pengolahan data selanjutnya menggunakan pengolahan variabel yang berbentuk *string*.



Gbr. 1 Diagram blok.

Pada sebuah halaman *web* terdapat beberapa bagian yang dapat berupa teks, gambar, *list*, atau tabel. Oleh karena itu, langkah pertama yang dilakukan adalah menentukan bagian yang merupakan tabel dari dokumen HTML. Tabel dideteksi dengan memanfaatkan fungsi pencocokan *string*. Tabel pada dokumen HTML dapat dikenali dengan tag `<table>` pada bagian awal dan `</table>` pada bagian akhir tabel. Jika pada sebuah halaman *web* terdapat dua atau lebih tabel, aplikasi menyimpan tabel tersebut pada indeks kedua dari variabel yang berbentuk *array*.

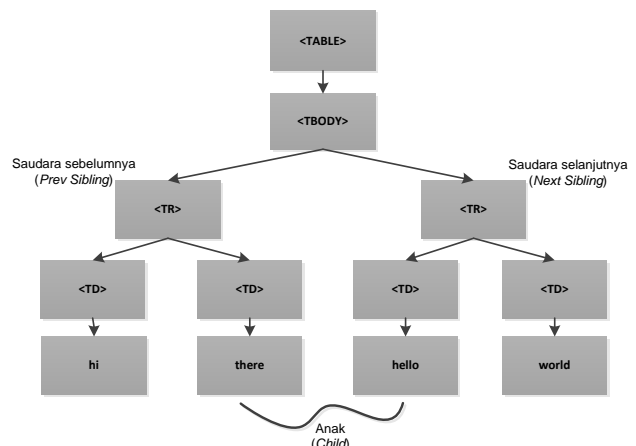
Langkah selanjutnya adalah membuat pohon DOM dari tabel yang dideteksi. Pohon DOM disusun dari tag pembentuk tabel, yaitu tag `<table>`, `<tr>`, `<th>`, dan `<td>`. Keempat tag ini dideteksi untuk menentukan bagian data pada tabel. Bagian data terletak di luar tag `<td>` atau `<th>`. Pada pohon DOM, bagian data ini merupakan daun dengan kedalaman paling rendah dari semua ranting. Sebagai contoh, Gbr. 2 merupakan pohon DOM dari Tabel II. Tampak bahwa data tabel berada pada daun terendah dari setiap ranting.

Langkah ketiga adalah mengekstrak atau mengambil bagian data dari pohon DOM. Untuk menelusuri daun terbawah dari pohon DOM, metode penelusuran yang digunakan adalah DFS. Metode ini menelusuri pohon diawali dari akar kemudian dilanjutkan ke daun sebelah kiri di bawahnya hingga ditemukan daun terbawah dari ranting paling kiri. Kemudian penelusuran dilanjutkan ke ranting di sebelahnya hingga ditemukan daun terbawah dari ranting tersebut. Penelusuran seperti ini dilakukan berulang kali hingga penelusuran mencapai daun terendah pada ranting paling kanan. Algoritme yang digunakan adalah algoritme DFS yang telah dijelaskan pada bagian II.

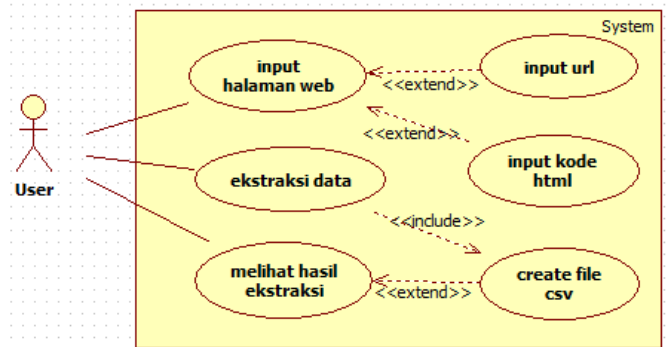
Setelah bagian data pada tabel diperoleh, langkah berikutnya adalah menyimpan data dalam bentuk *file csv*. Pemilihan bentuk *file* penyimpanan ditujukan agar hasil ekstraksi dapat digunakan untuk kebutuhan berikutnya, seperti integrasi dengan data dari tabel lain atau untuk disimpan dalam sebuah basis data.

TABEL II
CONTOH TABEL

hi	There
hello	World



Gbr. 2 Pohon DOM dari Tabel II.



Gbr. 3 Diagram *use case*.

B. Diagram Use Case

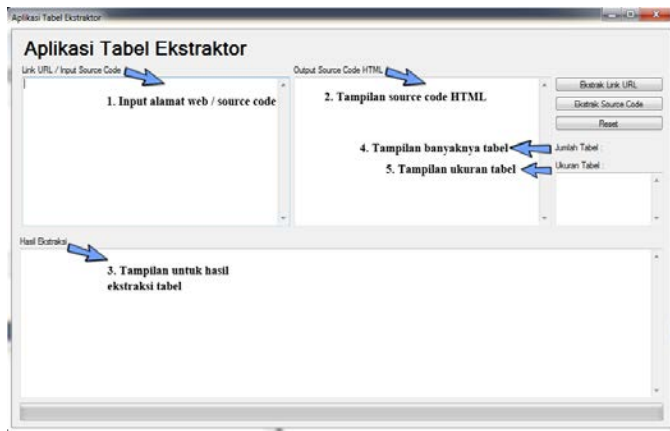
Perancangan diagram *use case* ditujukan untuk menggambarkan fungsionalitas aplikasi yang dibangun. Terdapat seorang aktor secara umum yang dinamakan dengan *user*. Selain itu, terdapat tiga fungsionalitas yang dapat dilakukan oleh *user* terhadap aplikasi. Yang pertama, *user* dapat memasukkan halaman *web* dalam bentuk alamat *URL* atau kode HTML dari halaman *web*. Dengan inisialisasi dari *user*, sistem dapat melakukan ekstraksi data dengan proses yang telah dijelaskan pada bagian sebelumnya. Selain itu, *user* juga dapat melihat data hasil ekstraksi langsung pada aplikasi atau pada *file csv* yang telah disimpan. Diagram *use case* dari aplikasi dapat dilihat pada Gbr. 3.

V. HASIL DAN PENGUJIAN

A. Tabel Ekstraktor

Aplikasi tabel ekstraktor ini dikembangkan dengan menggunakan bahasa pemrograman Visual Basic. Masing-masing fungsionalitas aplikasi yang terdapat pada diagram *use*

case diimplementasikan menjadi beberapa *method*. Ada *method* untuk mengolah *input* hingga diperoleh bagian tabel yang akan diekstrak, ada pula *method* untuk mengubah kode HTML menjadi pohon DOM yang kemudian ditelusuri dengan menggunakan *method* DFS. Kemudian ada juga *method* untuk menampilkan hasil ekstraksi dan mengubahnya menjadi *file* csv.



Gbr. 4 Tampilan aplikasi.

Gbr. 4 merupakan tampilan aplikasi yang dihasilkan. Terdapat satu bagian dari aplikasi untuk memasukkan alamat URL atau kode HTML halaman *web* yang ingin diekstrak. Setelah menekan tombol Ekstrak *Link URL* atau Ekstrak *Source code*, aplikasi akan menampilkan kembali *source code* HTML pada bagian Output *Source code* HTML, kemudian menampilkan hasil ekstraksi pada bagian Hasil Ekstraksi. Selain ditampilkan pada salah satu bagian pada aplikasi, hasil ekstraksi juga secara otomatis tersimpan dalam *file* csv. Aplikasi dapat berjalan jika terhubung dengan koneksi internet.

B. Pengujian

Pengujian dibagi menjadi dua bagian, yaitu pengujian fungsionalitas aplikasi dan pengujian hasil ekstraksi. Kedua pengujian menerapkan prinsip pengujian *black box*.

Pada pengujian fungsionalitas, prosedur pengujian menggunakan *suitability metrics* pada ISO 9126-2. Terdapat tiga metrik yang digunakan, yaitu:

1. *functional adequacy* (FA);
2. *functional implementation completeness* (FICM); dan
3. *functional implementation coverage* (FIC).

Berdasarkan ketiga metrik tersebut diperoleh hasil bahwa semua kebutuhan pengguna telah terpenuhi dan telah lengkap diimplementasikan.

Pengujian kedua, yakni pengujian untuk melihat kebenaran hasil ekstraksi. Pada pengujian ini, aplikasi tabel ekstraktor diujikan pada beberapa halaman *web* dengan beberapa karakteristik. Karakteristik yang diperhatikan pada pengujian ini antara lain eksistensi dan banyak tabel pada halaman *web*, ukuran tabel, orientasi *layout* tabel, serta ada atau tidaknya *merging cell* pada tabel. Tabel III merupakan rekapitulasi data uji yang digunakan.

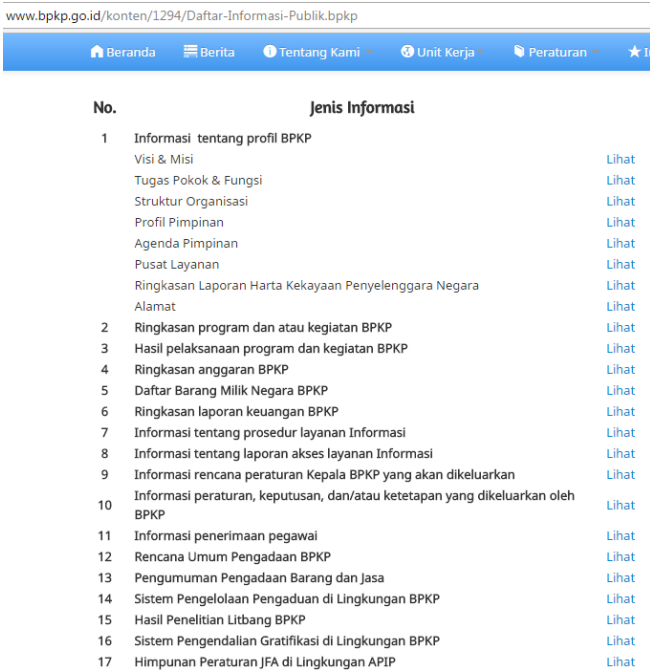
TABEL III
REKAPITULASI DATA PENGUJIAN

No.	Alamat Web	Jumlah Tabel	Ukuran Tabel	Kondisi Tabel
1.	http://www.bpkp.go.id/konten/1294/Daftar-Informasi-Publik.bpkp	1 (satu)	- Baris: 26 - Kolom: 3	- <i>Column wise Table</i> - Tidak ada <i>merging cell</i>
2.	http://www.bengkaliskab.go.id/statis-93-Bidang-Energi-dan-Sumber-Daya-mineral.html	1 (satu)	- Baris: 2 - Kolom: 3	- <i>Column wise Table</i> - Tidak ada <i>merging cell</i>
3.	http://diskominfo.riau.go.id/hal-daftar-ppid-provinsi-riau.html	2 (dua)	Tabel 1 - Baris: 58 - Kolom: 4 Tabel 2 - Baris: 13 - Kolom: 4	- <i>Column wise Table</i> - Terdapat <i>merging cell</i>
4.	https://lpse.lkpp.go.id/eproc4	1 (satu)	- Baris: 18 - Kolom: 4	- <i>Column wise Table</i> - Terdapat <i>merging cell</i>
5.	http://diskominfo.pemkomedan.go.id/statis-6-sumberdayamanusia.html	4 (empat)	Tabel 1 - Baris: 13 - Kolom: 4 Tabel 2 - Baris: 7 - Kolom: 4 Tabel 3 - Baris: 8 - Kolom: 10 Tabel 4 - Baris: 5 - Kolom: 4	- <i>Column wise Table</i> - Tidak ada <i>merging cell</i>
6.	http://www.pemkomedan.go.id/hal-kependudukan.html	4 (empat)	Tabel 1 - Baris: 19 - Kolom: 4 Tabel 2 - Baris: 13 - Kolom: 4 Tabel 3 - Baris: 24 - Kolom: 5 Tabel 4 - Baris: 22 - Kolom: 9	- <i>Column wise Table</i> - Terdapat <i>merging cell</i>
7.	http://diskominfo.riau.go.id/subbidang-3-jarkom.html	1 (satu)	- Baris: 7 - Kolom: 2	- <i>Row wise Table</i> - Terdapat <i>merging cell</i>

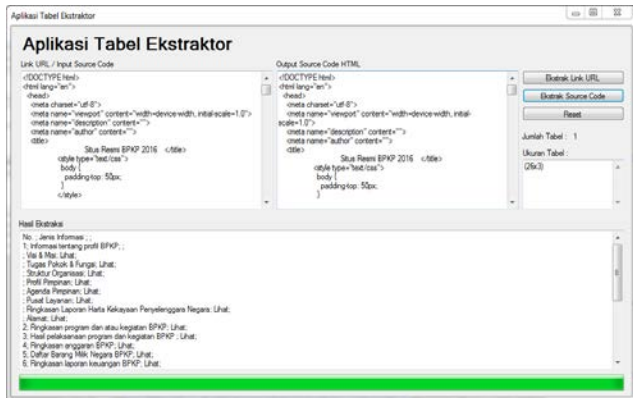
Pada bagian berikut akan dijelaskan beberapa data uji sebagai gambaran hasil pengujian aplikasi tabel ekstraktor.

1) *Data Uji 1*: Gbr. 5 merupakan tampilan halaman *web* data uji nomor 1 yang diakses pada 11 September 2016. Terdapat sebuah tabel yang berukuran 26 x 3 dengan *layout column wise*. Pada tabel tidak terdapat kondisi *merging cell*.

Alamat URL dimasukkan pada aplikasi tabel ekstraktor. Hasil ekstraksi dapat dilihat pada Gbr. 6 dan file csv yang dihasilkan dapat dilihat pada Gbr. 7. Aplikasi berhasil mengekstrak data pada tabel dan menyimpannya dalam file csv.



Gbr. 5 Tampilan halaman web data uji 1.



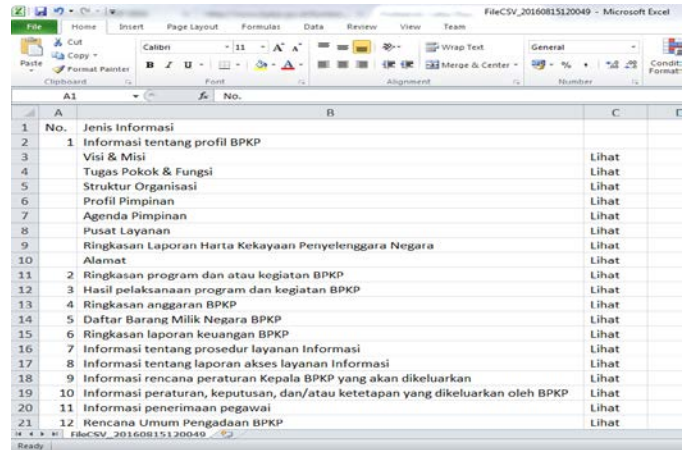
Gbr. 6 Hasil ekstraksi tabel pada data uji 1.

2) *Data Uji 2:* Gbr. 8 merupakan tampilan halaman web pada data uji nomor 3 yang diakses pada tanggal 11 September 2016. Terdapat dua buah tabel pada halaman web dengan ukuran masing-masing 58 x 4 dan 13 x 4. Orientasi layout kedua tabel adalah *column wise*. Pada kedua tabel terdapat *cell* yang merupakan *merging* dari beberapa *cell*.

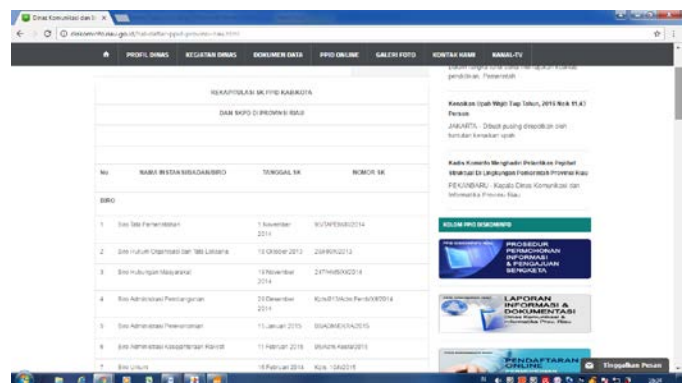
Gbr. 9 merupakan tampilan file csv hasil ekstraksi tabel pada data uji 2. Tampak bahwa aplikasi tabel ekstraktor berhasil mengekstrak data pada kedua tabel dan menyimpannya ke dalam file csv.

3) *Data Uji 3:* Gbr. 10 merupakan tampilan halaman web pada data uji nomor 7 yang diakses pada tanggal 11 September 2016. Pada halaman web terdapat sebuah tabel yang berukuran 7 x 2 dengan orientasi *layout row wise*. Pada

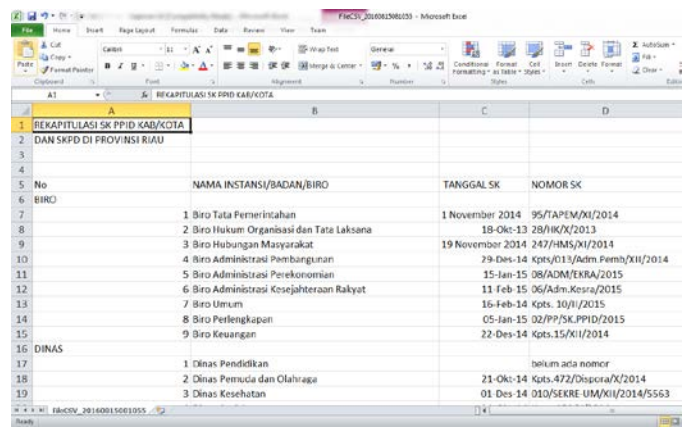
tabel juga terdapat kondisi *cell* yang merupakan *merging* dari beberapa *cell*.



Gbr. 7 Tampilan file csv data uji 1.



Gbr. 8 Tampilan halaman web data uji 2.



Gbr. 9 File csv hasil ekstraksi data uji 2.

Gbr. 11 merupakan tampilan file csv hasil ekstraksi data uji 3. Tampak bahwa aplikasi tabel ekstraktor berhasil mengekstrak data pada tabel dan menyimpannya dalam file csv.

C. Analisis Hasil Pengujian

Berdasarkan hasil pengujian fungsionalitas aplikasi, implementasi fitur pada aplikasi tabel ekstraktor telah sesuai

- [7] Shijun Li, Zhiyong Peng, and Mengchi Liu, "Extraction and integration information in HTML tables," in *Fourth International Conference on Computer and Information Technology (CIT)*, 2004.
- [8] Eko Prasetyo, Lukito Edi Nugroho, and Marcus Nurtiantara Aji, "Perancangan Data Warehouse Sistem Informasi Eksekutif untuk Data Akademik Program Studi," *JNTETI*, vol. 1, no. 3, pp. 13-20, November 2012.
- [9] Suhit Gupta, Gail Kaiser, David Neistadt, and Peter Grimm, "DOM-based content extraction of HTML documents," in *12th International Conference on World Wide Web*, 2003, pp. 207-214.
- [10] Agny Ismaya, "Algoritma Ekstraksi Informasi Berbasis Aturan," *JNTETI*, vol. 3, no. 4, pp. 242-247, November 2014.
- [11] Anhar, *Panduan menguasai PHP & MySQL secara otodidak*. Jakarta, Indonesia: Mediakita, 2010.
- [12] Nathanael T. Black and Wolfgang Ertel, *Introduction to artificial intelligence*. Germany: Springer science; Business Media, 2011.