

## APLIKASI PENGAMBILAN INFORMASI SITUS BERITA SECARA OTOMATIS MENGGUNAKAN *REGULAR EXPRESSION* DENGAN PLATFORM MICROSOFT .NET

Ario Koesnan <sup>1)</sup>, Yudhi Purwananto <sup>2)</sup>, Faisal Johan Atletiko <sup>2)</sup>

<sup>1, 2, 3)</sup> Teknik Informatika, FTIF, ITS, Surabaya, e-mail : [Sherlock\\_h30@yahoo.com](mailto:Sherlock_h30@yahoo.com)

*Abstract, Mosts news website nowadays are rich-content websites where they do not only provide news but also advertisements. Many of the advertisement banners in a news website are large-sized files. Therefore, it would take more time to download a news webpages. The purpose of this research is to help user with limited internet bandwidth to obtain news-only content from a news websites omitting its advertisements. Many news websites have common structures in indexing and directories. Regular expression can be used to adopt the structure of a news websites. By using windows services, the application will run in background to download news-only content. The results of this experiments showed that this application can download the contents of a news websites in background and users can read them offline. File size to be downloaded can also be reduced significantly.*

**Keywords :** *Regular expression, Windows Service, Background process, News Content*

Belakangan ini, masyarakat banyak disuguhkan berbagai macam situs berita mulai dari yang *thin content* dengan kategori yang sedikit seperti *bintang-indonesia.com* sampai situs yang menjadi portal berita yang *rich content* masa kini dan sering diakses dengan bermacam-macam kategorinya semacam *detik.com*.

Berdasarkan pengamatan penulis, semakin lama situs *rich content* semakin memberikan ruang yang banyak untuk penempatan *banner* atau iklan dengan alasan untuk meningkatkan pendapatan situs. Sebagai akibatnya dalam mendapatkan informasi, pengguna harus terbebani oleh *banner* atau iklan tersebut. Sehingga kebanyakan pengguna hanya mengakses situs-situs berita tersebut untuk membaca berita terbaru yang umumnya tersaji melalui kepala berita (*headline*) berita pada halaman utama situs tersebut. Banyak diantara pengguna, terutama dari SOHO (*Small Office Home Office*) dengan *bandwidth* koneksi internet terbatas, jarang mengakses berita utama dari *headline* yang sebagian besar digunakan untuk *men-download* gambar-gambar maupun *banner-banner* iklan dari situs berita yang kebanyakan belum tentu berguna dan tidak berhubungan dengan berita yang dibaca.

Sejalan dengan itu, proses pengambilan informasi secara keseluruhan dari situs berita secara otomatis ini akan memberikan

kemudahan dan keefisienan bagi pengguna dalam memperoleh berita yang lengkap tanpa harus mengakses langsung situs-situs berita tersebut. Proses pengambilan informasi secara otomatis ini juga sangat memudahkan pengarsipan berita dan memahami kronologis kejadian dari berita tersebut.

Tujuan pembuatan aplikasi ini adalah membuat perangkat lunak yang bekerja sebagai proses *background (background process)* yang melakukan pengambilan berita dari situs-situs berita yang telah ditentukan pengguna. Pengambilan dilakukan secara otomatis dan periodik tanpa harus mengakses langsung ke situs berita tersebut dengan melakukan penyaringan (*filtering*) terhadap situs berita sesuai skenario yang telah terdefinisi dalam *profile* yang dikenali melalui *regular expression*.

Permasalahan yang diangkat dalam aplikasi ini adalah :

Bagaimana membuat suatu sistem aplikasi berbasis *web service* untuk *men-download* secara otomatis informasi berita utama dari satu atau lebih situs berita.

Bagaimana merancang dan membuat suatu skenario melalui *regular expression* yang mampu menyaring informasi utama dari situs berita yang memiliki berbagai macam jenis struktur baik *file* maupun datanya. Bagaimana mempelajari dan mengetahui site map situs,

metode pengambilan informasi, serta memberikan saringan (*filter*) yang sesuai pada situs-situs berita yang berbagai macam jenisnya tersebut.

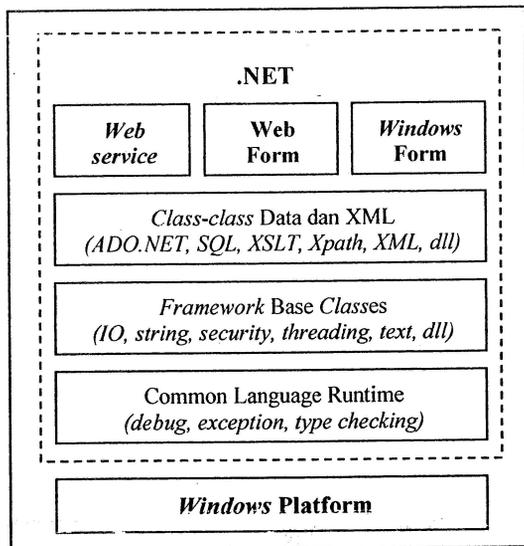
Manfaat dari pembuatan aplikasi ini adalah:

1. Memudahkan pengguna dengan bandwidth internet terbatas untuk mendapatkan berita secara cepat.
2. Pengguna dapat tetap memperoleh informasi berita terbaru tanpa perlu melakukan browsing ke situs berita.

**FRAMEWORK .NET**

*Framework .NET* adalah sebuah *platform computing* baru yang dirancang untuk menyederhanakan pengembangan aplikasi pada lingkungan terdistribusi pada internet [08]. Bahasa yang didukung oleh *Framework .NET* adalah Visual Basic.NET, C++, ASP.NET, Jscript.NET dan C#.

Pada gambar 1 terlihat *framework .NET* berada diatas sistem operasi, yang berupa sistem operasi buatan Microsoft, yaitu *Windows 2000 (Server dan profesional)*, *Windows NT 4.0 (Server dan workstation)*, *Windows Millenium Edition*, *Windows 98* dan *Windows XP Profesional*.



Gambar 1. Framework .NET

**WINDOWS SERVICE**

*Windows service* sebelumnya disebut *NT service*. Fungsi utama dari suatu *windows service* adalah untuk menjalankan aplikasi sebagai *background*. *Windows service* mulai

berjalan jauh sebelum pengguna melakukan *log/masuk* ke sistem (jika *windows service* tersebut diatur untuk mulai/start pada proses *boot-up*). Sebuah *windows service* juga dapat diatur sehingga pengguna harus melakukan *start* secara manual.

*Windows service* memiliki proses sendiri dan oleh karena itu ia dapat berjalan dengan sangat efisien. Secara normal, *windows service* tidak memiliki antarmuka.

**WEB SERVICE**

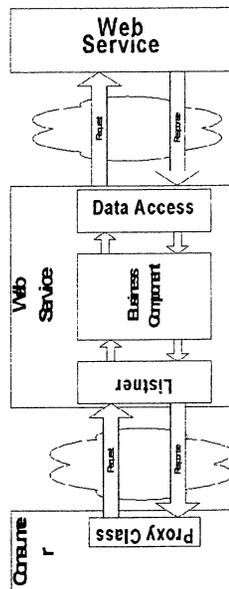
Suatu *Web service* adalah sebuah komponen yang menyediakan layanan untuk pelanggan yang menggunakan protokol internet standar (HTTP, XML) untuk mengakses layanan-layanan ini. *Web service* adalah komponen terprogram yang dapat diakses melalui protokol internet. *Web service* menggunakan XML dan *Simple Object Access Protocol (SOAP)* untuk berkomunikasi dengan pelanggan. XML menyediakan bahasa terstandarisasi untuk pertukaran data dalam format yang diterima secara luas. SOAP adalah protokol berbasis XML yang sederhana yang berjalan melalui HTTP untuk melakukan pertukaran data dalam lingkungan heterogen yang terdistribusi. Gambar 2 menunjukkan arsitektur dari *web service*.

Ketika menggunakan Visual Studio untuk membuat dan menggunakan *web service*, kita tidak perlu mengetahui arsitektur ini. *.NET Framework* mengkonversi semuanya, kita dapat membuat dan menggunakan *web service* tanpa terlebih dahulu harus mempelajari XML dan SOAP.

**REGULAR EXPRESSION**

*Regular expression* menyediakan metode yang kuat, fleksibel dan efisien untuk memproses teks [04]. Notasi pencocokan pola yang luas memungkinkan teks dengan ukuran besar dapat diparsing dengan cepat untuk menemukan pola karakter tertentu; untuk mengekstrak, mengedit, mengganti, atau menghapus *substring* teks; atau menambahkan *string* terekstraksi ke dalam koleksi untuk menghasilkan laporan. Untuk banyak aplikasi yang berhubungan dengan *string* (misalnya pemrosesan HTML, *parsing log file*, dan *parsing HTTP header*), *regular expression* merupakan *tool* yang diperlukan. Dengan desain yang disesuaikan dengan *regular expression* dari Perl 5, *Regular expression*

.NET Framework meliputi fitur pencocokan dari kanan ke kiri (right-to-left matching) dan kompilasi sambil jalan (on-the-fly compilation).



Gambar 2. Arsitektur Web Service

Bahasa *regular expression* didesain dan dioptimasi untuk memanipulasi teks. Bahasanya terdiri dari dua tipe karakter dasar: karakter teks literal (normal) dan metakarakter-metakarakter.

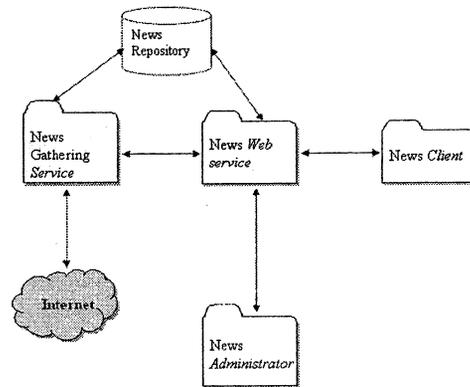
Metakarakter `?` dan `*` adalah metakarakter yang sudah familiar dan digunakan dengan DOS file sistem untuk mewakili suatu karakter apapun atau kumpulan karakter. Perintah DOS COPY \*.DOC A: memerintahkan sistem untuk menyalin file apa saja dengan ekstensi .DOC ke dalam disket di drive A. Metakarakter `*` mewakili nama file apa saja didepan ekstensi .DOC. *Regular expression* juga menyediakan pencarian yang lebih kompleks. Sebagai contoh, *regular expression* `(?<char>w)k<char>`, menggunakan grup bernama dan backreferencing, mencari karakter berdekatan yang berpasangan.

Framework .NET memiliki class-class built-in untuk memproses *regular expression*, yaitu class **Regex**, **Match**, **MatchCollection**, **Group Collection**, **CaptureCollection**, **Group**, dan **Capture**.

**METODOLOGI**  
**Arsitektur Sistem**

*News Administrator* mengatur segala hal yang berhubungan dengan *inserting*, *updating*

dan *deleting* data ke *database*. Aplikasi inilah yang menangani data-data baru yang akan masuk ke sistem, seperti situs baru, alamat URL dan sebagainya. Demikian juga jika ada perubahan maupun penghapusan terhadap *database*. *News Web service* berfungsi sebagai *service* yang mengantarai *server* dan *client*. Fungsi-fungsi yang digunakan *client* umumnya yang berhubungan dengan akses ke *database*, sedangkan *server* menggunakan fungsi-fungsi untuk *parsing* (*filtering*). *News Client* merupakan aplikasi bagi pengguna (*client*) berisi antarmuka untuk membaca dan mencari berita. *News Gathering Service* berfungsi sebagai *server* yang bertugas mengambil berita secara berkala. Aplikasi ini merupakan *windows service* yang berjalan sebagai proses *background* pada saat mengambil berita. Gambar 3 menunjukkan arsitektur sistem



Gambar 3. Arsitektur Sistem

**Pengguna Aplikasi**

Ada tiga macam pengguna aplikasi, yaitu *Administrator User*, *Administrator Regex* dan *Client*.

1. **Administrator User**  
Pengguna ini berfungsi sebagai operator yang bertugas memasukkan data-data ke dalam database, maka *skill* yang diperlukan adalah *entry data*.
2. **Administrator Regex**  
Pengguna ini bertugas menganalisis situs berita dan membuat formula *regex*-nya. *Skill* yang diperlukan adalah mengerti *Regular Expression* dan HTML.
3. **Client**  
Pengguna ini berfungsi sebagai *client* yang memerlukan informasi utama dari situs berita tanpa mengutamakan iklan/*banner*,

gambar, dan menu/link lainnya, hanya berita utamanya saja.

**Perancangan Skenario Regular Expression**

Ada beberapa langkah yang harus dilakukan sebelum merancang skenario ini, yaitu sebagai berikut:

**1. Pengenalan Situs**

Dalam tahap ini dilakukan pengenalan pola situs, mulai dari struktur situs, penamaan direktori, URL, dan lain-lain yang diperlukan bagi perancangan *regex* nantinya. Untuk mengenali suatu situs harus didapatkan *source HTML* dari situs tersebut. Umumnya diambil file yang mengandung indeks berita dari situs yang bersangkutan. Dari *source* ini kemudian diambil pola dengan melihat *tag-tag* khusus yang menentukan posisi *link* indeks berita, waktu berita, judul bahkan isi berita itu sendiri.

**2. Pembuatan regex untuk link, time, dan title**

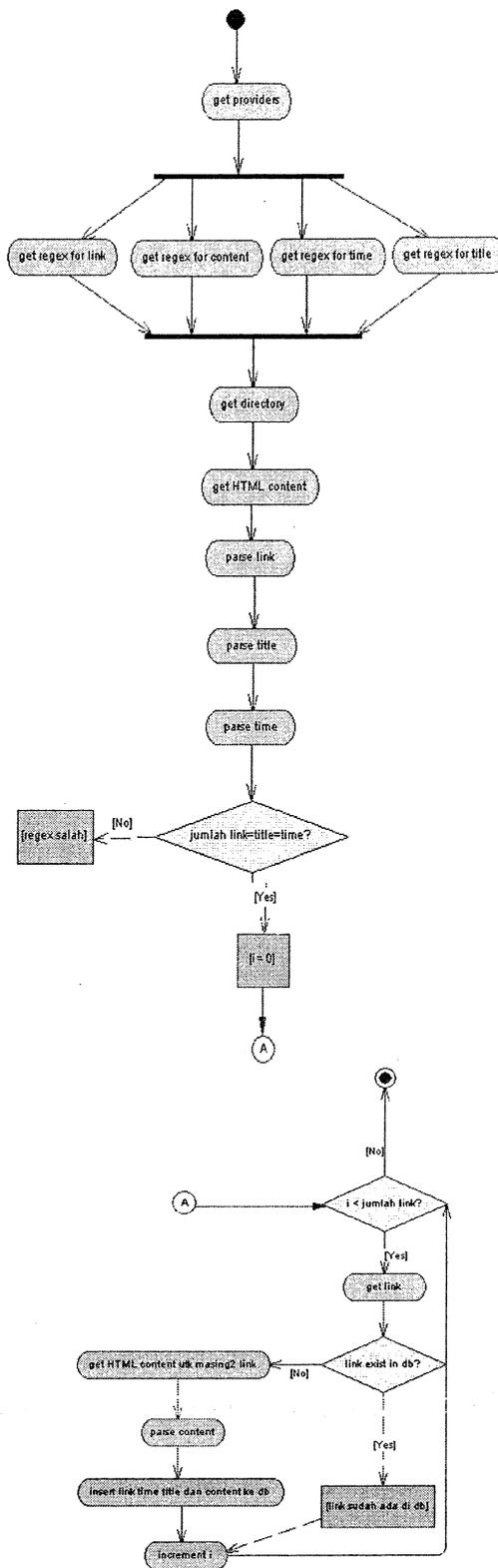
Setelah mengetahui pola situs kita dapat membuat formula *regex* untuk memfilter aspek-aspek yang kita perlukan. Yang paling penting di sini adalah *filtering* terhadap *link* indeks berita, karena *link-link* inilah yang nantinya akan merujuk kepada isi/*content* beritanya. Untuk dapat mengambil *link-link* yang sesuai, harus diketahui dimana posisi mulai dan akhir dari *link-link* yang akan diambil. Selain *link-link* tersebut, dilakukan pula *filtering* untuk mendapatkan waktu berita (*time*) dan judul berita (*title*) juga dengan *regex* tertentu yang dibuat berdasarkan pengenalan pola situs yang telah dilakukan sebelumnya.

**3. Pembuatan regex untuk content**

Setelah mendapatkan *link-link* indeks berita yang akan diambil, dilakukan pengambilan *content/isi* berita juga dengan *regex* hasil pengenalan pola dari *link* yang merujuk ke berita utama tersebut.

**Activity Diagram Pengambilan Berita**

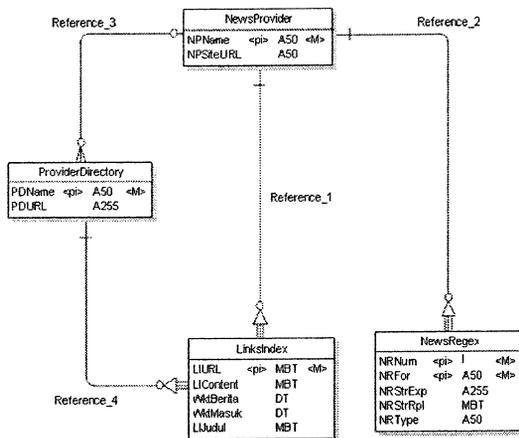
Diagram ini menunjukkan aliran proses utama yang terjadi pada perangkat lunak ini, yaitu proses pengambilan berita. Proses ini dilakukan sebagai proses background (*windows service*) pada *server*. Gambar 4 menunjukkan aliran proses yang terjadi saat *service* melakukan pengambilan berita.



Gambar 5. Activity Diagram Pengambilan Berita

**Perancangan CDM dan PDM**

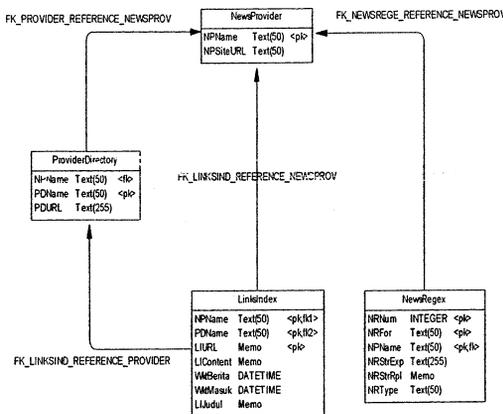
Rancangan *conceptual diagram* (CDM) dari *database* sistem perangkat lunak ini, dapat dilihat pada gambar 6. Sedangkan diagram PDM ditunjukkan oleh gambar 7.



Gambar 6. Diagram CDM

**PEMBAHASAN Implementasi Rancangan Sistem**

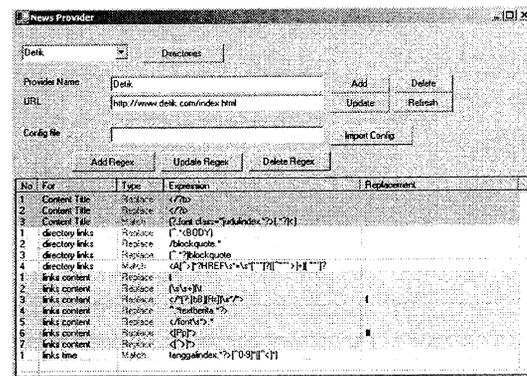
Untuk mengimplementasikan rancangan sistem *client-server* ini, digunakan tools Microsoft Visual Studio .NET dalam membuat *web service*. Bahasa pemrograman yang digunakan adalah VB.NET. *Web service* tidak menghasilkan bentuk antarmuka melainkan hanya API bagi sistem lainnya. Deskripsi API ini dijelaskan dalam bentuk WSDL. *Web service* ini dirancang dengan tujuan sebagai jembatan antara *Client* dengan *Server*. Metode lain yang biasa dipakai adalah *remoting/socket*.



Gambar 7. Diagram PDM

**Aplikasi Administrator**

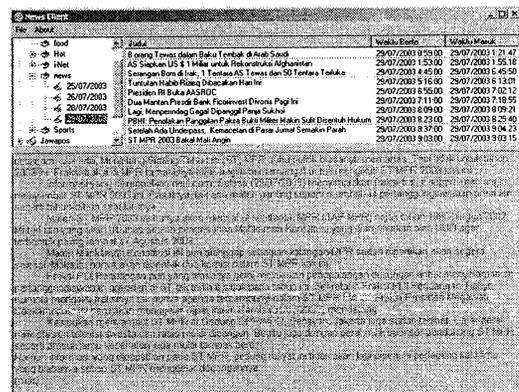
Aplikasi *Administrator* menggunakan *web service* untuk memanggil fungsi-fungsi untuk memanipulasi tabel-tabel dalam *database*, yaitu tabel *NewsProvider*, *ProviderDirectory*, dan *NewsRegex*. Tampilan awal dari aplikasi *Administrator* ditunjukkan pada gambar 8.



Gambar 8. Tampilan Aplikasi Administrator

**Aplikasi Client**

Aplikasi *Client* menggunakan *web service* untuk mengambil isi berita secara *offline* dari *database* untuk dibaca oleh *client/pengguna*. Aplikasi *Client* dibuat menyerupai Microsoft Outlook (*MS.Outlook-like*) untuk memudahkan interaksi dengan *client/pengguna*. Tampilan awal aplikasi ini ditunjukkan pada gambar 9.



Gambar 9. Tampilan Aplikasi Client

**Uji Coba**

Pada bagian ini akan dilakukan uji coba terhadap proses pengambilan berita. Situs-situs berita yang akan digunakan sebagai analisa



PDName	POURL
Ekonomi Bisnis	http://www.republika.co.id/index_benis_online.asp?kal_id=21
Luar Negeri2	http://www.republika.co.id/index_benis_online.asp?kal_id=248
Nasional	http://www.republika.co.id/index_benis_online.asp?kal_id=23
Olahraga2	http://www.republika.co.id/index_benis_online.asp?kal_id=247
Topik Kita	http://www.republika.co.id/index_benis_online.asp?kal_id=156

Gambar 13. Data Direktori Republika

NIRNum	NIRFor	NIRName	NIRSExp	NIRScript	NIRType
1		Republika			2
2		Republika			1
12		Republika			1
22		Republika			1
32		Republika			1
42		Republika			1
52		Republika			1
62		Republika			1
72		Republika			1
82		Republika			1
92		Republika			1
102		Republika			1
112		Republika			1
122		Republika			1
132		Republika			1
142		Republika			1
152		Republika			1
162		Republika			1
172		Republika			1
182		Republika			1
192		Republika			1
202		Republika			1
212		Republika			1
222		Republika			1
232		Republika			1
242		Republika			1
252		Republika			1
262		Republika			1
272		Republika			1
282		Republika			1
292		Republika			1
302		Republika			1
312		Republika			1
322		Republika			1
332		Republika			1
342		Republika			1
352		Republika			1
362		Republika			1
372		Republika			1
382		Republika			1
392		Republika			1
402		Republika			1
412		Republika			1
422		Republika			1
432		Republika			1
442		Republika			1
452		Republika			1
462		Republika			1
472		Republika			1
482		Republika			1
492		Republika			1
502		Republika			1

Gambar 14. Data Regex Republika

Gambar 15 menunjukkan aplikasi client ketika membaca berita.

Judul	Waktu Berita	Waktu Masuk
1. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
2. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
3. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
4. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
5. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
6. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
7. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
8. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
9. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
10. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
11. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
12. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
13. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
14. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
15. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
16. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
17. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
18. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
19. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
20. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
21. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
22. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
23. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
24. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
25. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
26. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
27. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
28. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
29. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
30. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
31. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
32. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
33. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
34. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
35. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
36. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
37. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
38. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
39. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
40. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
41. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
42. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
43. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
44. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
45. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
46. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
47. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
48. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
49. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
50. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00

Gambar 15. Tampilan Aplikasi Client untuk Situs Republika

3) Uji coba Situs [www.kompas.com](http://www.kompas.com)

Uji coba ini dilakukan pada situs kompas dengan direktori *Ekonomi, Gaya Hidup, Kesehatan, Teknologi, Metropolitan, Olah Raga*, dan *Topik Kita*. Skenario regex untuk kompas ini sengaja dibagi menjadi dua, yaitu kompas dengan direktori *Indeks*, dan kompas2 dengan direktori-direktori seperti diuraikan di atas. Hal ini dilakukan karena pola antara indeks berita utama kompas dengan direktori-direktori lainnya berlainan. Gambar 16, 17, dan 18 menunjukkan data direktori dan regex situs kompas.

PDName	POURL
Indeks	http://www.kompas.co.id/utama/index.cfm
Ekonomi	http://www.kompas.co.id/bisnis/bisnis.cfm
Gaya Hidup	http://www.kompas.co.id/gayahidup/index.htm
Kesehatan	http://www.kompas.co.id/kesehatan/index.htm
Metropolitan	http://www.kompas.co.id/metro/index.htm
Olah Raga	http://www.kompas.co.id/olahraga/index.htm
Teknologi	http://www.kompas.co.id/teknologi/index.htm

Gambar 16. Data Direktori Kompas dan Kompas2

NIRNum	NIRFor	NIRName	NIRSExp	NIRScript	NIRType
1		Kompas			2
2		Kompas			1
12		Kompas			1
22		Kompas			1
32		Kompas			1
42		Kompas			1
52		Kompas			1
62		Kompas			1
72		Kompas			1
82		Kompas			1
92		Kompas			1
102		Kompas			1
112		Kompas			1
122		Kompas			1
132		Kompas			1
142		Kompas			1
152		Kompas			1
162		Kompas			1
172		Kompas			1
182		Kompas			1
192		Kompas			1
202		Kompas			1
212		Kompas			1
222		Kompas			1
232		Kompas			1
242		Kompas			1

Gambar 17. Data Regex Kompas

NIRNum	NIRFor	NIRName	NIRSExp	NIRScript	NIRType
1		Kompas2			2
2		Kompas2			1
12		Kompas2			1
22		Kompas2			1
32		Kompas2			1
42		Kompas2			1
52		Kompas2			1
62		Kompas2			1
72		Kompas2			1
82		Kompas2			1
92		Kompas2			1
102		Kompas2			1
112		Kompas2			1
122		Kompas2			1
132		Kompas2			1
142		Kompas2			1
152		Kompas2			1
162		Kompas2			1
172		Kompas2			1
182		Kompas2			1
192		Kompas2			1
202		Kompas2			1
212		Kompas2			1
222		Kompas2			1
232		Kompas2			1
242		Kompas2			1

Gambar 18. Data Regex Kompas2

Gambar 19 dan 20 menunjukkan aplikasi client ketika mengakses informasi berita Kompas.

Judul	Waktu Berita	Waktu Masuk
1. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
2. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
3. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
4. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
5. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
6. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
7. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
8. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
9. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
10. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
11. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
12. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
13. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
14. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
15. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
16. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
17. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
18. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
19. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
20. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
21. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
22. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
23. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
24. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
25. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
26. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
27. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
28. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
29. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
30. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
31. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
32. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
33. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
34. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
35. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
36. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
37. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
38. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
39. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
40. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
41. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
42. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
43. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
44. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
45. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
46. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
47. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
48. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
49. G. Deak	18/07/2003 14:29:00	18/07/2003 14:29:00
50. G. Deak	18/07/2003 14	

#### 4) *Ujicoba perbandingan ukuran file*

Telah dilakukan ujicoba untuk membandingkan ukuran file yang harus didownload jika pengguna langsung melakukan browsing ke sebuah halaman situs berita Detik dengan ukuran file berita (dalam ekstensi .txt) yang didownload oleh aplikasi ini. Untuk kasus browsing langsung, diperoleh data ukuran file yang didownload adalah 94,5 KB dengan *size-on-disk* adalah 184 KB. Sedangkan ukuran file berita yang didownload oleh aplikasi ini adalah 1,49 KB dengan *size-on-disk* adalah 4,00 KB.

#### SIMPULAN

- 1) Untuk mengambil informasi berita utama/main content dari suatu situs berita yang memiliki pola/struktur situs yang tetap atau memiliki web content management, digunakan suatu formula *Regular Expression* yang sesuai dengan pola situs masing-masing.
- 2) Perancangan dan pembuatan skenario *Regular Expression* dapat dilakukan dengan mempelajari struktur/pola suatu situs.
- 3) Formula *Regular Expression* dirumuskan dalam satu file konfigurasi untuk setiap penyedia berita oleh *Administrator Regex* yang kemudian file konfigurasi tersebut akan diimport dan dimasukkan ke database melalui aplikasi *News Administrator* oleh *Administrator User*.
- 4) Satu situs penyedia berita dapat memiliki lebih dari satu pola/struktur situs sehingga memiliki lebih dari satu formula *Regular Expression*.
- 5) Teknologi *Windows Service* digunakan sebagai aplikasi yang berjalan dalam proses *background* untuk melakukan pengambilan berita secara otomatis.
- 6) Teknologi *Web Service* dengan protokol HTTP dapat menjadi penghubung antara *Server* dan *Client*.

Untuk pengembangan :

- 1) Penggunaan *pattern recognition* untuk dapat mengenali pola/struktur situs secara dinamis.
- 2) Penelitian lebih lanjut untuk dapat menghasilkan formula *Regular Expression* secara otomatis.

#### DAFTAR RUJUKAN

- [01] Haryanto, Steve. *Majalah Masterweb MWMag*. Indonesia: Oktober 2001, November 2001, *issue 04, issue 06*
- [02] Jain, Jayesh. *Creating A Windows Service in VB.NET*. [www.devarticles.com](http://www.devarticles.com). Oktober 2002
- [03] JISC TechWatch Report: *Content Management Systems*. [http://www.jisc.ac.uk/techwatch/reports/tsw\\_01-02.pdf](http://www.jisc.ac.uk/techwatch/reports/tsw_01-02.pdf). September 2001
- [04] Microsoft Service Digital Network (MSDN). *.NET Framework Developer's Guide .NET Framework Regular Expressions*, <http://msdn.microsoft.com/library/default.asp?url=/library/en-us/cpguide/html/cpconcomregularexpressions.asp>, 2002
- [05] *Regex Library*; [www.regexlib.com](http://www.regexlib.com)
- [06] Regular Expression di .NET; <http://www.techmedia-online.com/tutorials/t02070101.asp?id=1953>
- [07] Rockford Lhotka. *Implementing a Background Process in Visual Basic .NET*, <http://msdn.microsoft.com/library/default.asp?url=/library/en-us/dnadvnet/html/vbnet09272002.asp>, Magenic Technologies, 2002
- [08] Thai, Thuan, Huang Q. Lam. *.NET Framework Essential*. 1st Edition, O'Reilly, 2001
- [09] *VB.NET Fast Track*. Syngress, 2001
- [10] Situs-situs berita; [www.republika.co.id](http://www.republika.co.id), [www.jawapos.co.id](http://www.jawapos.co.id), [www.kompas.co.id](http://www.kompas.co.id), [www.surya.co.id](http://www.surya.co.id), [www.detik.com](http://www.detik.com), [www.reuters.com](http://www.reuters.com), [www.astaga.com](http://www.astaga.com), [news.bbc.co.uk](http://news.bbc.co.uk) dan situs-situs berita lainnya