

SKRIPSI

*WEB CRAWLING* TERDISTRIBUSI PADA LINGKUNGAN  
HADOOP



Gabriella

NPM: 2014730013

PROGRAM STUDI TEKNIK INFORMATIKA  
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS  
UNIVERSITAS KATOLIK PARAHYANGAN  
2018

**UNDERGRADUATE THESIS**

**DISTRIBUTED WEB CRAWLING ON HADOOP  
ENVIRONMENT**



**Gabriella**

**NPM: 2014730013**

**DEPARTMENT OF INFORMATICS  
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES  
PARAHYANGAN CATHOLIC UNIVERSITY  
2018**

**LEMBAR PENGESAHAN**



**WEB CRAWLING TERDISTRIBUSI PADA LINGKUNGAN  
HADOOP**

**Gabriella**

**NPM: 2014730013**

**Bandung, 31 Mei 2018**

**Menyetujui,**

**Pembimbing**

A handwritten signature in black ink, appearing to be 'Gede Karya', is written over a faint circular stamp.

**Gede Karya, M.T., CISA, IPM**

**Ketua Tim Penguji**

A handwritten signature in black ink, appearing to be 'Vania Natali', is written over a faint circular stamp.

**Vania Natali, M.T.**

**Anggota Tim Penguji**

A handwritten signature in black ink, appearing to be 'Pascal Alfadian', is written over a faint circular stamp.

**Pascal Alfadian, M.Comp.**

**Mengetahui,**

**Ketua Program Studi**

A handwritten signature in black ink, appearing to be 'Mariskha Tri Adithia', is written over a faint circular stamp.

**Mariskha Tri Adithia, P.D.Eng**



## PERNYATAAN

Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa skripsi dengan judul:

### **WEB CRAWLING TERDISTRIBUSI PADA LINGKUNGAN HADOOP**

adalah benar-benar karya saya sendiri, dan saya tidak melakukan penjiplakan atau pengutipan dengan cara-cara yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan.

Atas pernyataan ini, saya siap menanggung segala risiko dan sanksi yang dijatuhkan kepada saya, apabila di kemudian hari ditemukan adanya pelanggaran terhadap etika keilmuan dalam karya saya, atau jika ada tuntutan formal atau non-formal dari pihak lain berkaitan dengan keaslian karya saya ini.

Dinyatakan di Bandung,  
Tanggal 31 Mei 2018



Gabriella  
NPM: 2014730013

## ABSTRAK

*Web crawler* merupakan program yang melakukan *web scanning* dan *data indexing* dengan metode tertentu. Salah satu metode *crawling* yakni *distributed web crawling* yang memanfaatkan banyak agen *crawler* dengan tujuan mempercepat proses *crawling*. Data hasil *crawling* dapat berukuran besar, atau yang sering disebut sebagai *big data*, sehingga diperlukan media penyimpanan yang dapat mengolah *big data*. Untuk mengatasi masalah ini, dapat digunakan Hadoop ecosystem, seperti: HDFS dan HBase. Oleh karena itu, pada penelitian ini dikembangkan aplikasi *distributed web crawler* pada lingkungan Hadoop *Cluster* (HBase). Penelitian juga mencakup eksperimen untuk menjawab pertanyaan: (1) bagaimana pengaruh banyaknya *crawler* terhadap kecepatan *crawling*; dan (2) bagaimana kecepatan pemrosesan data menggunakan Hadoop.

Untuk menyelesaikan permasalahan ini, dilakukan analisis dan perancangan arsitektur dan algoritma *distributed crawling*. Setelah itu, dikembangkan aplikasi *web crawler* yang terdiri atas situs induk, *server crawler*, dan agen *crawler* terdistribusi. Berikutnya, disiapkan Hadoop *Cluster* (HDFS dan HBase) untuk pengujian dan pengumpulan URL *seed* untuk data uji. Langkah terakhir yang dilakukan yakni melaksanakan pengujian dan eksperimen performa *distributed crawler*.

Pengujian dilakukan untuk memastikan perangkat lunak telah berfungsi sebagaimana mestinya. Pengujian ini dilakukan berdasarkan *usecase* dan skenario yang diperoleh dari analisis kebutuhan perangkat lunak. Berdasarkan pengujian fungsional tersebut, didapati perangkat lunak sudah berjalan sebagaimana mestinya. Eksperimen dilakukan untuk melihat performa *crawler*; dilakukan dengan menjalankan aplikasi *web crawler* pada empat komputer yang tergabung dalam *cluster* Hadoop; satu *master* dan tiga *slave*, dimana masing-masing komputer menjalankan lima agen *crawler*. Ada dua eksperimen yang dilakukan, yakni eksperimen *crawling* dan eksperimen *searching*. Eksperimen *crawling* dilakukan untuk melihat bagaimana pengaruh banyaknya *crawler* dan *node* komputer pada *cluster* Hadoop yang digunakan terhadap kecepatan *crawling*. Eksperimen *searching* dilakukan untuk melihat bagaimana pengaruh banyaknya *node* komputer pada *cluster* Hadoop yang digunakan terhadap kecepatan pencarian.

Pada eksperimen *crawling*, didapatkan rata-rata kecepatan *crawling* bertambah sebesar 1.16 kali untuk setiap penambahan satu *node* dan lima agen. Pada eksperimen *searching*, didapatkan rata-rata peningkatan kecepatan pencarian adalah sebesar 1.19 kali untuk setiap penambahan satu *node*. Berdasarkan kedua hasil eksperimen tersebut, dapat disimpulkan bahwa semakin banyak agen dan komputer dalam *cluster* Hadoop yang digunakan, maka semakin cepat pula proses *crawling*. Selain itu, semakin banyak komputer dalam *cluster* Hadoop yang digunakan, semakin cepat pula waktu pencarian.

**Kata-kata kunci:** *web crawler*, *distributed*, Hadoop

## ABSTRACT

Web crawler is a program that performs web scanning and data indexing with certain methods. One of its method is distributed web crawling that utilizes many crawler agents with the aim of speeding up the crawling process. Data from the crawling result can be large, which is often referred to as big data. Therefore, a storage media which can process a big data is required. To solve this problem, Hadoop ecosystem, such as HDFS and HBase, can be used. On this study, a web crawler application was developed on Hadoop Cluster environment (HBase). The study also includes experiments to answer: how the number of crawlers affects crawling speed; and the speed of data processing on Hadoop environment.

To solve these problems, analysis and design of distributed crawler architecture and algorithm are performed. After that, a web crawler application that consists of main site, crawler server, and distributed crawler agent was developed. Next, Hadoop Cluster (HDFS and HBase) for testing was prepared and the URL seed were collected as a test data. The last step is to conduct testing and performance experiment on the distributed crawler.

Testing is done to ensure the software is working properly. This test is based on the usecase and scenarios obtained from the software requirement analysis. Based on the functional testing, the software is found to be working properly. The experiment was conducted to see the performance of the crawler; done by running the web crawler application on four computers which were part of the Hadoop cluster; one master and three slaves, where each computer run five crawler agents. There are two experiments conducted, namely crawling experiment and searching experiment. Crawling experiment were conducted to see how the number of crawlers and computer nodes on the Hadoop cluster used affects crawling speed. The searching experiment was conducted to see how the number of computer nodes on the Hadoop cluster used affects searching speed.

Based on the crawling experiment result, the average crawling speed increased by 1.16 times for each addition of one node and five agents. Based on the searching experiment result, the average searching speed increased 1.19 times for each addition of one node. Based on these two experiments results, it can be concluded that the more agents and computers in the Hadoop cluster are used, the faster the crawling process can be done. In addition, the more computers in the Hadoop cluster are used, the faster the searching process can be done.

**Keywords:** web crawler, distributed, hadoop

*Dipersembahkan untuk orang tua dan adik penulis*

## KATA PENGANTAR

Puji syukur penulis haturkan pada Tuhan yang Maha Esa atas segala berkat dan rahmat-Nya, sehingga skripsi berjudul "*Web Crawling* Terdistribusi Pada Lingkungan Hadoop" dapat penulis selesaikan dengan baik. Skripsi ini dibuat untuk melengkapi persyaratan kelulusan tingkat sarjana (S-1) di Fakultas Teknologi Informasi dan Sains Jurusan Teknik Informatika Universitas Katolik Parahyangan. Selama pembuatan skripsi ini, penulis mendapatkan bantuan dan dukungan dari berbagai pihak. Maka dari itu, pada kesempatan ini penulis hendak mengucapkan terima kasih kepada:

1. Bapak Gede Karya, M.T., CISA, IPM selaku dosen pembimbing yang telah menyediakan waktu untuk membimbing penulis selama penyusunan skripsi.
2. Ibu Vania Natali, M.T. dan Bapak Pascal Alfadian, M.Comp. selaku dosen penguji yang telah memberi kritik dan saran untuk skripsi ini.
3. Orang tua dan adik penulis yang selalu memberi doa dan dukungan selama penyusunan skripsi.
4. Sahabat dan rekan-rekan penulis yang telah banyak membantu dan memberi dukungan selama penyusunan skripsi.

Akhir kata, penulis menyadari bahwa skripsi ini belum sempurna. Namun, penulis berharap skripsi ini dapat bermanfaat bagi para pembaca yang ingin menambah wawasan ataupun melakukan penelitian yang serupa.

Bandung, Mei 2018

Penulis



# DAFTAR ISI

<b>KATA PENGANTAR</b>	<b>xv</b>
<b>DAFTAR ISI</b>	<b>xvii</b>
<b>DAFTAR GAMBAR</b>	<b>xix</b>
<b>DAFTAR TABEL</b>	<b>xxiii</b>
<b>1 PENDAHULUAN</b>	<b>1</b>
1.1 Latar Belakang . . . . .	1
1.2 Rumusan Masalah . . . . .	1
1.3 Tujuan . . . . .	2
1.4 Batasan Masalah . . . . .	2
1.5 Metodologi . . . . .	2
1.6 Sistematika Pembahasan . . . . .	3
<b>2 LANDASAN TEORI</b>	<b>5</b>
2.1 <i>Web Crawler</i> . . . . .	5
2.1.1 <i>Arsitektur Web Crawler</i> . . . . .	5
2.1.2 <i>Algoritma Web Crawler</i> . . . . .	7
2.1.3 <i>File Robots.txt</i> . . . . .	7
2.1.4 <i>Distributed Web Crawling</i> . . . . .	8
2.2 <i>Big Data</i> . . . . .	9
2.3 Hadoop . . . . .	9
2.3.1 HDFS . . . . .	10
2.3.2 MapReduce . . . . .	11
2.3.3 HBase . . . . .	13
2.4 Java 2 Enterprise Edition (J2EE) . . . . .	16
2.4.1 JavaServer Pages (JSP) . . . . .	17
2.4.2 Servlet . . . . .	19
2.5 <i>Library Jsoup</i> . . . . .	20
<b>3 EKSPLORASI DAN ANALISIS</b>	<b>21</b>
3.1 Eksplorasi Hadoop . . . . .	21
3.1.1 Eksplorasi HDFS . . . . .	21
3.1.2 Eksplorasi HBase . . . . .	24
3.2 Analisis Masalah dan Usulan Solusi . . . . .	28
3.2.1 Analisis Masalah . . . . .	28
3.2.2 Usulan Solusi . . . . .	28
3.3 Analisis Kebutuhan Perangkat Lunak . . . . .	34
3.3.1 Deskripsi Analisis Perangkat Lunak . . . . .	34
3.3.2 Diagram <i>Use Case</i> dan Skenario Situs Induk . . . . .	35
3.3.3 Kebutuhan <i>Input</i> Perangkat Lunak . . . . .	38

3.3.4	Kebutuhan <i>Output</i> Perangkat Lunak . . . . .	38
3.3.5	Kebutuhan Atribut Pada Tabel di <i>Database</i> . . . . .	38
3.3.6	Diagram Kelas Sederhana . . . . .	39
<b>4</b>	<b>PERANCANGAN</b>	<b>43</b>
4.1	Perancangan <i>Database</i> Pada HBase . . . . .	43
4.2	Desain <i>User Interface</i> . . . . .	45
4.3	Diagram Kelas Rinci . . . . .	51
4.3.1	Diagram Kelas <i>Crawler</i> . . . . .	51
4.3.2	Diagram Kelas Situs Induk . . . . .	55
4.3.3	UrlChecker . . . . .	56
<b>5</b>	<b>IMPLEMENTASI, PENGUJIAN, DAN EKSPERIMEN</b>	<b>61</b>
5.1	Implementasi . . . . .	61
5.1.1	Lingkungan Implementasi Perangkat Keras . . . . .	61
5.1.2	Lingkungan Implementasi Perangkat Lunak . . . . .	61
5.1.3	Implementasi <i>Database</i> . . . . .	61
5.1.4	Implementasi Situs Induk . . . . .	64
5.1.5	Implementasi <i>Web Crawler</i> . . . . .	71
5.2	Pengujian . . . . .	77
5.3	Eksperimen . . . . .	79
5.3.1	Konfigurasi Lingkungan Eksperimen . . . . .	79
5.3.2	Eksperimen <i>Crawling</i> . . . . .	81
5.3.3	Eksperimen Pencarian <i>Term</i> . . . . .	82
<b>6</b>	<b>KESIMPULAN DAN SARAN</b>	<b>85</b>
6.1	Kesimpulan . . . . .	85
6.2	Saran . . . . .	85
	<b>DAFTAR REFERENSI</b>	<b>87</b>
<b>A</b>	<b>KODE PROGRAM</b>	<b>89</b>
A.1	Kode Program <i>Web Crawler</i> . . . . .	89
A.1.1	<i>Package engine</i> . . . . .	89
A.1.2	<i>Package database</i> . . . . .	94
A.2	Kode Program Situs Induk . . . . .	97
A.2.1	<i>Package engine</i> . . . . .	97
A.2.2	<i>Package database</i> . . . . .	100
A.2.3	<i>Package servlet</i> . . . . .	104
<b>B</b>	<b>HASIL EKSPERIMEN</b>	<b>109</b>
B.1	Hasil Eksperimen Pencarian <i>Term</i> . . . . .	109
B.2	Hasil Perhitungan Banyak URL . . . . .	117
B.3	Penyimpanan Pada Hadoop . . . . .	118
B.4	Foto dan <i>Screenshot</i> Eksperimen <i>Crawling</i> . . . . .	123

## DAFTAR GAMBAR

2.1	Arsitektur <i>Web Crawler</i> . . . . .	6
2.2	Isi <i>File robots.txt</i> . . . . .	8
2.3	<i>High-level Hadoop Architecture</i> . . . . .	10
2.4	Arsitektur HDFS . . . . .	11
2.5	Mekanisme Kerja MapReduce . . . . .	12
2.6	Mekanisme <i>Shuffle and Sort</i> . . . . .	13
2.7	Komponen RegionServer . . . . .	14
2.8	Arsitektur HBase . . . . .	15
2.9	Skema Tabel HBase . . . . .	15
2.10	Pemrosesan JSP . . . . .	18
2.11	JSP <i>Life Cycle</i> . . . . .	19
2.12	Arsitektur Servlet . . . . .	19
3.1	Menghapus <i>Metadata</i> di NameNode . . . . .	21
3.2	Memulai HDFS dan YARN . . . . .	21
3.3	Potongan Informasi NameNode di <i>Command Prompt</i> . . . . .	21
3.4	Potongan Informasi DataNode di <i>Command Prompt</i> . . . . .	22
3.5	Potongan Informasi ResourceManager di <i>Command Prompt</i> . . . . .	22
3.6	Potongan Informasi NodeManager di <i>Command Prompt</i> . . . . .	22
3.7	<i>Overview dfshealth</i> . . . . .	22
3.8	<i>Summary dfshealth</i> . . . . .	23
3.9	Informasi DataNode . . . . .	23
3.10	Memasukkan <i>File</i> ke HDFS dan <i>Me-list File</i> . . . . .	23
3.11	Operasi <i>Wordcount</i> Menggunakan MapReduce di <i>Command Prompt</i> . . . . .	24
3.12	Menghentikan HDFS dan YARN . . . . .	24
3.13	Menjalankan HBase . . . . .	24
3.14	Memulai HBase <i>Shell</i> . . . . .	25
3.15	Melihat Status <i>Cluster</i> , Versi, dan <i>Current User</i> HBase dengan <i>Shell</i> . . . . .	25
3.16	Membuat Tabel Menggunakan HBase <i>Shell</i> . . . . .	25
3.17	Informasi Tabel HBase . . . . .	25
3.18	Informasi RegionServer . . . . .	26
3.19	Memasukkan Data ke Tabel Menggunakan HBase <i>Shell</i> (1) . . . . .	26
3.20	Memasukkan Data ke Tabel Menggunakan HBase <i>Shell</i> (2) . . . . .	26
3.21	Melihat Data Tabel "table1" Menggunakan <i>Command scan</i> di HBase <i>Shell</i> . . . . .	26
3.22	Menghentikan HBase . . . . .	27
3.23	Meng- <i>import</i> HBase <i>Library</i> Pada Program Java . . . . .	27
3.24	Membuat Koneksi ke HBase Menggunakan HBase API . . . . .	27
3.25	Membuat Tabel Menggunakan HBase API . . . . .	28
3.26	Memasukkan Data ke Tabel Menggunakan HBase API . . . . .	28
3.27	Melakukan <i>scan</i> Tabel Menggunakan HBase API . . . . .	28
3.28	<i>Flowchart</i> Proses <i>Distributed Crawling</i> . . . . .	32
3.29	Arsitektur <i>Distributed Web Crawler</i> . . . . .	34

3.30	<i>Use Case</i> Situs Induk	35
3.31	Diagram Kelas Sederhana (Situs Induk)	40
3.32	Diagram Kelas Sederhana ( <i>Web Crawler</i> )	42
4.1	Tampilan <i>Homepage</i>	45
4.2	Tampilan <i>Pop-up</i> Registrasi	46
4.3	Tampilan <i>Pop-up Log in</i>	46
4.4	Tampilan Halaman <i>Admin Menu</i>	47
4.5	Tampilan <i>Pop-up Add URL</i>	47
4.6	Tampilan Halaman <i>Frontier</i>	48
4.7	Tampilan Halaman <i>Web Repository</i>	49
4.8	Tampilan Halaman <i>Search Result</i> Untuk Admin	50
4.9	Tampilan Halaman <i>Search Result</i> Untuk <i>User</i>	50
4.10	Tampilan Halaman <i>Content</i>	51
4.11	Diagram Kelas <i>Crawler (Package engine)</i>	52
4.12	Diagram Kelas <i>Crawler (Package database)</i>	54
4.13	Diagram Kelas Situs Induk ( <i>Package engine</i> )	56
4.14	Diagram Kelas Situs Induk ( <i>Package database</i> )	57
4.15	Diagram Kelas Situs Induk ( <i>Package servlet</i> )	59
5.1	Membuat Tabel <i>Frontier</i> Menggunakan HBase API	62
5.2	Membuat Tabel <i>Web Repository</i> dan <i>Secondary Indexer</i> -nya Menggunakan HBase API	62
5.3	Membuat Tabel Admin Menggunakan HBase API	63
5.4	Memasukkan Data ke Tabel <i>Frontier</i> Menggunakan HBase API	63
5.5	Memasukkan Data ke Tabel <i>Web Repository</i> Menggunakan HBase API	63
5.6	Memasukkan Data ke Tabel Admin Menggunakan HBase API	64
5.7	Tampilan <i>Homepage</i>	64
5.8	Tampilan <i>Pop-up</i> Register	65
5.9	Tampilan <i>Pop-up Log in</i>	65
5.10	Tampilan <i>Admin Menu</i>	66
5.11	Tampilan <i>Pop-up Add URL File</i>	66
5.12	Tampilan Halaman <i>Frontier</i>	67
5.13	Tampilan Halaman <i>Frontier</i> (Kosong)	67
5.14	Tampilan Halaman <i>Repository</i>	68
5.15	Tampilan Halaman <i>Repository</i> (Kosong)	68
5.16	Tampilan Halaman Hasil Pencarian (Admin)	69
5.17	Tampilan Halaman Hasil Pencarian ( <i>User</i> )	69
5.18	Tampilan Halaman Hasil Pencarian Kosong (Admin)	69
5.19	Tampilan Halaman Hasil Pencarian Kosong ( <i>User</i> )	70
5.20	Tampilan Halaman Konten (Admin)	70
5.21	Tampilan Halaman Konten ( <i>User</i> )	71
5.22	Pengecekan Konten Pada Proses Pencarian <i>Term</i>	71
5.23	Memeriksa, Mengambil, dan Menandai URL dari <i>Frontier</i>	72
5.24	Memasukkan URL ke <i>Frontier</i>	73
5.25	Memasukkan URL ke <i>Web Repository</i>	73
5.26	Proses Ekstraksi URL	74
5.27	Menandai URL Untuk Proses <i>Re-crawling</i>	75
5.28	Meng- <i>update</i> Data Hasil <i>Crawling</i> Pada Tabel <i>Web Repository</i> dan <i>Secondary Indexer</i> -nya	75
5.29	Inisialisasi <i>Socket</i> Milik <i>Server</i>	76
5.30	Inisialisasi <i>Socket</i> Milik Agen	76

5.31	Inisialisasi <i>Reader</i> dan <i>Writer</i> . . . . .	76
5.32	Konfigurasi Jaringan Internet . . . . .	80
5.33	Konfigurasi Perangkat Lunak . . . . .	81
5.34	Banyaknya URL Hasil <i>Crawling</i> dengan Penggunaan 1, 2, 3, dan 4 <i>Node</i> . . . . .	81
5.35	Perbandingan Banyak <i>Node</i> Terhadap Kecepatan Pencarian <i>Term</i> . . . . .	82
B.1	Pencarian dengan 1 <i>Node</i> Saat URL bBerjumlah 1105 Buah . . . . .	109
B.2	Pencarian dengan 2 <i>Node</i> Saat URL bBerjumlah 1105 Buah . . . . .	110
B.3	Pencarian dengan 3 <i>Node</i> Saat URL bBerjumlah 1105 Buah . . . . .	110
B.4	Pencarian dengan 4 <i>Node</i> Saat URL bBerjumlah 1105 Buah . . . . .	111
B.5	Pencarian dengan 1 <i>Node</i> Saat URL Berjumlah 1332 Buah . . . . .	111
B.6	Pencarian dengan 2 <i>Node</i> Saat URL Berjumlah 1332 Buah . . . . .	112
B.7	Pencarian dengan 3 <i>Node</i> Saat URL Berjumlah 1332 Buah . . . . .	112
B.8	Pencarian dengan 4 <i>Node</i> Saat URL Berjumlah 1332 Buah . . . . .	113
B.9	Pencarian dengan 1 <i>Node</i> Saat URL Berjumlah 1568 Buah . . . . .	113
B.10	Pencarian dengan 2 <i>Node</i> Saat URL Berjumlah 1568 Buah . . . . .	114
B.11	Pencarian dengan 3 <i>Node</i> Saat URL Berjumlah 1568 Buah . . . . .	114
B.12	Pencarian dengan 4 <i>Node</i> Saat URL Berjumlah 1568 Buah . . . . .	115
B.13	Pencarian dengan 1 <i>Node</i> Saat URL Berjumlah 1721 Buah . . . . .	115
B.14	Pencarian dengan 2 <i>Node</i> Saat URL Berjumlah 1721 Buah . . . . .	116
B.15	Pencarian dengan 3 <i>Node</i> Saat URL Berjumlah 1721 Buah . . . . .	116
B.16	Pencarian dengan 4 <i>Node</i> Saat URL Berjumlah 1721 Buah . . . . .	117
B.17	Perhitungan Banyak URL Untuk Eksperimen dengan 1 <i>Node</i> . . . . .	117
B.18	Perhitungan Banyak URL Untuk Eksperimen dengan 2 <i>Node</i> . . . . .	117
B.19	Perhitungan Banyak URL Untuk Eksperimen dengan 3 <i>Node</i> . . . . .	117
B.20	Perhitungan Banyak URL Untuk Eksperimen dengan 4 <i>Node</i> . . . . .	118
B.21	Melihat Data Penyimpanan Tabel Admin Melalui Hadoop UI di <i>Port</i> 50070 . . . . .	118
B.22	Melihat Data Penyimpanan Tabel <i>Frontier</i> Melalui Hadoop UI di <i>Port</i> 50070 . . . . .	118
B.23	Melihat Data Penyimpanan Tabel <i>Web Repository</i> Melalui Hadoop UI di <i>Port</i> 50070 . . . . .	119
B.24	Melihat Data Penyimpanan Tabel <i>Secondary Indexer</i> untuk <i>web repository</i> Melalui Hadoop UI di <i>Port</i> 50070 . . . . .	119
B.25	<i>File Data columnfamily depth</i> di Tabel <i>Frontier</i> . . . . .	120
B.26	<i>File Data columnfamily crawling</i> di Tabel <i>Frontier</i> . . . . .	120
B.27	<i>File Data columnfamily page</i> di Tabel <i>Web Repository</i> . . . . .	121
B.28	<i>File Data columnfamily agent</i> di Tabel <i>Web Repository</i> . . . . .	121
B.29	<i>File Data columnfamily depth</i> di Tabel <i>Web Repository</i> . . . . .	122
B.30	<i>File Data columnfamily idrvalue</i> di Tabel <i>Secondary Indexer</i> Milik <i>Web Repository</i> . . . . .	122
B.31	Foto Eksperimen <i>Crawling</i> (1) . . . . .	123
B.32	Foto Eksperimen <i>Crawling</i> (2) . . . . .	123
B.33	<i>Screenshot</i> Proses Eksperimen <i>Crawling</i> 1 <i>Node</i> (Komputer <i>master</i> ) . . . . .	124
B.34	<i>Screenshot</i> Proses Eksperimen <i>Crawling</i> 2 <i>Node</i> (Komputer <i>master</i> ) . . . . .	124
B.35	<i>Screenshot</i> Proses Eksperimen <i>Crawling</i> 3 <i>Node</i> (Komputer <i>master</i> ) . . . . .	125
B.36	<i>Screenshot</i> Proses Eksperimen <i>Crawling</i> 4 <i>Node</i> (Komputer <i>master</i> ) . . . . .	125

## DAFTAR TABEL

3.1	Perbandingan Penelusuran URL dengan Metode BFS dan DFS . . . . .	29
3.2	Skenario <i>Log in</i> . . . . .	36
3.3	Skenario Register . . . . .	36
3.4	Skenario Memasukkan URL . . . . .	36
3.5	Skenario Melihat <i>Frontier</i> . . . . .	37
3.6	Skenario Melihat <i>Web Repository</i> . . . . .	37
3.7	Skenario <i>Search Term</i> . . . . .	38
3.8	Atribut Pada Tabel <i>Frontier</i> . . . . .	39
3.9	Atribut Pada Tabel <i>Web Repository</i> . . . . .	39
3.10	Atribut Pada Tabel Admin . . . . .	39
4.1	Tabel <i>Frontier</i> ( <i>RowKey</i> dan <i>ColumnFamily crawling</i> ) . . . . .	43
4.2	Tabel <i>Frontier</i> ( <i>ColumnFamily depth</i> ) . . . . .	44
4.3	Tabel <i>Web Repository</i> ( <i>RowKey</i> , <i>ColumnFamily page</i> , dan <i>ColumnFamily agent</i> ) . . . . .	44
4.4	Tabel <i>Web Repository</i> ( <i>ColumnFamily depth</i> ) . . . . .	44
4.5	Tabel <i>Repository Indexer</i> . . . . .	44
4.6	Tabel Admin . . . . .	45
5.1	Hasil Pengujian Situs Induk . . . . .	78
5.2	<i>Hostname</i> dan <i>IP Address</i> Komputer Pengujian . . . . .	79
5.3	Peningkatan Banyaknya Hasil <i>Crawling</i> . . . . .	82
5.4	Peningkatan Kecepatan Pencarian Untuk 1105 URL . . . . .	83
5.5	Peningkatan Kecepatan Pencarian Untuk 1332 URL . . . . .	83
5.6	Peningkatan Kecepatan Pencarian Untuk 1568 URL . . . . .	84
5.7	Peningkatan Kecepatan Pencarian Untuk 1721 URL . . . . .	84

# BAB 1

## PENDAHULUAN

Pada bab ini dilakukan pembahasan mengenai latar belakang, rumusan masalah, tujuan, batasan masalah, metodologi, dan sistematika pembahasan.

### 1.1 Latar Belakang

Pada jaman sekarang, orang-orang diberi kemudahan dalam mencari informasi dengan adanya *web*. *Web* atau *World Wide Web* merupakan suatu "ruang" informasi yang berisi *resource* dalam berbagai bentuk, misalnya dokumen dan gambar. Tiap *web resource* berada pada suatu tempat yang disebut halaman *web*, dan memiliki alamat atau *Uniform Resource Locator* (URL) sebagai penanda lokasi halaman *web* tersebut. Agar suatu halaman dapat di-*retrieve* ketika melakukan pencarian, *search engine* yang digunakan harus terlebih dahulu mengumpulkan informasi mengenai halaman tersebut. Pengumpulan informasi ini dilakukan dengan menggunakan *web crawler*.

*Web crawler* merupakan program yang melakukan *web scanning* dan *data indexing* dengan metode tertentu. Ada berbagai metode *crawling*, antara lain: *batch and incremental crawling*; *deep crawling*; *distributed crawling*; *focused crawling*; dll. Data yang dihasilkan dari proses *crawling* berupa URL atau halaman-halaman *web*. Pada skripsi ini, dibuat *web crawler* dengan metode *distributed crawling*. *Distributed web crawling* memanfaatkan banyak *crawler* untuk melakukan proses *crawling* dengan tujuan untuk mempercepat proses *crawling*.

Data hasil *crawling* dapat berukuran besar ataupun berjumlah banyak. Jenis data seperti ini sering disebut sebagai *big data*. Oleh karena itu, diperlukan suatu media penyimpanan yang dapat memproses atau mengolah *big data*. Hadoop merupakan salah satu *framework* yang digunakan untuk penyimpanan terdistribusi dan pemrosesan *big data*. Inti dari Hadoop terdiri atas bagian penyimpanan yang disebut *Hadoop Distributed File System* (HDFS) dan bagian pemrosesan yang menggunakan model pemrograman MapReduce.

Aplikasi *distributed web crawler* yang dikembangkan memanfaatkan teknologi HBase, yakni *database* Hadoop yang berjalan diatas HDFS. Oleh karena *web crawler* yang dibuat menggunakan metode *distributed web crawling* dan memanfaatkan Hadoop sebagai media penyimpanan data, maka dilakukan pula eksperimen untuk menjawab pertanyaan: (1) bagaimana pengaruh banyaknya *crawler* terhadap kecepatan *crawling*; dan (2) bagaimana kecepatan pemrosesan data menggunakan Hadoop.

### 1.2 Rumusan Masalah

Salah satu masalah dalam *web crawling* adalah proses yang memerlukan waktu yang cukup lama. Untuk itu, pada skripsi ini ingin diselesaikan masalah peningkatan kecepatan *crawling* dengan agen terdistribusi. Oleh karena itu, permasalahan dapat dirumuskan sebagai berikut:

1. Bagaimana mekanisme pembagian tugas dan koordinasi antar agen sehingga bisa meningkatkan kecepatan *crawling*?

2. Bagaimana mekanisme penyimpanan informasi oleh masing-masing agen pada lingkungan *file system* terdistribusi Hadoop?

### 1.3 Tujuan

Tujuan dari skripsi ini adalah sebagai berikut:

1. Memahami mekanisme pembagian tugas dan koordinasi antar agen sehingga bisa meningkatkan kecepatan *crawling*.
2. Memahami mekanisme penyimpanan data oleh masing-masing agen pada lingkungan *file system* terdistribusi Hadoop.
3. Membangun perangkat lunak *web crawler* dengan metode *distributed web crawling* yang memanfaatkan Hadoop sebagai media penyimpanan URL.
4. Melakukan eksperimen untuk melihat pengaruh banyaknya agen dan *node* komputer pada *cluster* Hadoop terhadap kecepatan *crawling* dan pencarian.

### 1.4 Batasan Masalah

Karena fokus penelitian pada skripsi ini adalah mengenai kecepatan *crawling* dan mekanisme penyimpanan hasil *crawling* pada lingkungan Hadoop, maka batasan-batasan berikut diterapkan pada aplikasi yang dikembangkan:

1. Proses *crawling* tidak memproses *crawler policy*. *Crawler policy* memiliki aturan *politeness* yang memungkinkan suatu *web* tidak bisa di-*crawl* oleh *web crawler*. Jika aturan ini diproses maka keakuratan hasil penelitian kecepatan *crawling* dan pemrosesan data dapat berkurang karena bisa saja ada URL yang seharusnya dapat diproses, tetapi karena *policy* tersebut akhirnya URL tidak jadi diproses.
2. Penelusuran URL hanya berdasarkan HTML. Penelusuran *active link* seperti *javascript* tidak akan dilakukan karena memerlukan teknik penelusuran yang berbeda dan tidak terkait fokus penelitian dari skripsi ini.
3. Hasil yang disimpan berupa teks yang terdapat pada halaman *web*. Penyimpanan *tag* diperlukan jika yang dibuat adalah *search engine*, yang berada diluar fokus penelitian pada skripsi ini.

### 1.5 Metodologi

Berikut ini adalah langkah-langkah yang dilakukan dalam membangun perangkat lunak:

1. Studi pustaka tentang konsep *web crawler*.
2. Studi pustaka dan eksplorasi tentang sistem terdistribusi berbasis Hadoop.
3. Studi teknik *distributed crawling*.
4. Studi teknik penyimpanan pada Hadoop.
5. Menganalisis pemilihan penyimpanan pada Hadoop.
6. Mengembangkan perangkat lunak *web crawler* yang terdiri atas situs induk, *server crawler* dan agen *crawler*.
7. Melaksanakan pengujian dan eksperimen performa pada perangkat lunak.



---

## 1.6 Sistematika Pembahasan

Sistematika pembahasan dalam skripsi ini adalah sebagai berikut:

1. Bab 1 Pendahuluan  
Bab ini membahas tentang latar belakang, rumusan masalah, tujuan, batasan masalah, metodologi, dan sistematika pembahasan.
2. Bab 2 Landasan Teori  
Bab ini membahas tentang teori-teori yang menunjang pengembangan perangkat lunak, antara lain teori *web crawler*, *big data*, Hadoop, Java 2 *Enterprise Edition*, dan *library jsoup*.
3. Bab 3 Eksplorasi dan Analisis  
Bab ini membahas tentang eksplorasi Hadoop, analisis masalah dan usulan solusi, serta analisis kebutuhan perangkat lunak.
4. Bab 4 Perancangan  
Bab ini membahas tentang perancangan *User Interface* dan *database*, serta diagram kelas rinci.
5. Bab 5 Implementasi, Pengujian, dan Eksperimen  
Bab ini membahas implementasi, pengujian, dan eksperimen.
6. Bab 6 Kesimpulan dan Saran  
Bab ini membahas tentang kesimpulan yang didapatkan dan saran yang dapat diberikan berdasarkan hasil analisis, perancangan, implementasi, pengujian, dan eksperimen yang telah dilakukan.