

BAB II

LANDASAN TEORI

2.1. Web Scraping

Menurut (Zhao, 2019), *Web Extraction* atau yang lebih dikenal dengan *web scraping* merupakan sebuah teknik untuk mengekstraksi data dari sebuah website dan menyimpannya ke sistem file atau basis data untuk suatu kepentingan. Biasanya, data pada web yang diekstraksi perlu melalui protokol *Hypertext Transfer Protocol (HTTP)* atau melalui web browser. Kegiatan *web scraping* bahkan sudah dilakukan baik secara manual oleh pengguna atau secara otomatis oleh bot.

2.2. Python

Bahasa pemrograman python merupakan salah satu bahasa pemrograman yang dibuat secara dinamis. Ada banyak librari python yang dibuat untuk keperluan *web scraping*. Beberapa diantaranya adalah *request*, *beautiful soup*, *lxml*, *selenium* dan *scrapy*. Selain untuk mengekstrak data *HTML* yang mentah untuk berbagai keperluan, *request* pada python dapat digunakan untuk mengirim form dan mengakses *API* (Prakash & Rashid, 2017).

Penggunaan bahasa python untuk keperluan komputasi ilmiah telah mendapatkan momentum dalam beberapa tahun terakhir. Menjadi fakta bahwa bahasa python memiliki kompleksitas dan mudah dibaca dengan dilengkapi librari ilmiah yang menjadi pendukung dalam penerapan karakteristiknya (Tejedor et al., 2016).

2.3. Metode Web Scraping

Web Scraping memiliki berbagai macam metode yang telah dikembangkan dalam berbagai penelitian, diantaranya:

a. *XPath*

XPath (XML Path Language) merupakan bahasa kueri untuk memilih bagian-bagian (*nodes*) dari sebuah dokumen *XML*. Selain untuk memilih, *XPath* juga dapat digunakan untuk menghitung nilai seperti: *string*, angka dan *boolean* pada sebuah file *XML* dan *HTML*. *World Wide Web Consortium (W3C)* bahkan telah menetapkan standardisasi dalam penggunaan *XPath*.

b. *CSS Selector*

CSS Selector merupakan sebuah metode untuk menemukan elemen *HTML* di halaman website dan mengekstrak data darinya. *CSS Selector* dideklarasikan sebagai bagian dari sebuah gaya markup yang berlaku untuk mencocokkan dengan tag dan atribut dalam markup tersebut.

c. *HTML DOM*

HTML DOM merupakan objek model standar untuk mendapatkan, mengubah, menambah atau menghapus elemen *HTML*. *DOM* bekerja dengan cara mendefinisikan objek dan properti dari semua elemen *HTML*, dengan metode untuk mengaksesnya. Sebuah web browser tidak mengharuskan untuk menggunakan *DOM* dalam menampilkan dokumen *HTML*. Namun dengan *DOM*, *Javascript* dapat mengakses semua elemen di dalam dokumen *HTML*.

d. *Regular Expression (Regex)*

Regular Expression (Regex) merupakan konstruksi bahasa untuk mencocokkan teks berdasarkan pola tertentu, terutama untuk kasus-kasus yang kompleks. *Regex* juga digunakan untuk mencocokkan pola-pola karakter tertentu dalam suatu kumpulan *string*. *Regex* memiliki dua macam karakter yaitu karakter biasa dan meta karakter.

2.4. Multiprocessing

Multiprocessing merupakan kemampuan suatu sistem untuk mendukung lebih dari satu prosesor pada saat bersamaan. Suatu program yang menggunakan *multiprocessing* akan dipecah menjadi rutinitas yang lebih kecil yang berjalan secara independent. Sistem operasi akan mengalokasikan komputasi ini ke prosesor yang meningkatkan kinerja sistem.

Penggunaan *multiprocessing* dengan menjadikan proses pengolahan menjadi paralel merupakan hal yang perlu untuk mencapai kinerja terbaik. Setiap tugas *multiprocessing* akan berjalan dalam prosesnya sendiri, setiap program yang berjalan di komputer diwakili oleh satu atau lebih proses (Abeykoon et al., 2017; Sherman & Hartog, 2016).

2.5. Penelitian Terkait

Penelitian yang dilakukan (Rizaldi & Putranto, 2017), yaitu membandingkan metode *web scraping* menggunakan *XPath* dan *CSS Selector* dengan parameter perhitungan waktu, jumlah item, pengukuran memori dan pengukuran *bandwidth*. Kemudian penelitian yang dilakukan (Gunawan et al., 2019) masih dalam

perbandingan metode *web scraping* menggunakan metode *Regular Expression*, *HTML DOM* dan *XPath* dengan parameter pengukuran waktu, pengukuran memori dan pengukuran data.

Beberapa penelitian yang telah dilakukan sebelumnya menjadi landasan acuan untuk melakukan pengukuran jumlah objek data, penggunaan *CPU*, penggunaan memori, penggunaan *bandwith* dan *execution time* pada proses *web scraping* menggunakan empat metode, yaitu *XPath*, *CSS Selector*, *HTML DOM* dan *Regular Expression (Regex)* dengan penerapan *multiprocessing*.

Table II-1 Tabel Penelitian Terkait

1	Judul	<i>Parallelizing X-ray Photon Correlation Spectroscopy Software Tools using Python Multiprocessing</i>
	Penulis	Sameera K. Abeykoon, Meifeng Lin dan Kerstin Kleese van Dam (2017)
	Penelitian	Menggunakan modul multiprosesor Python untuk memparalelkan fungsi korelasi waktu dalam <i>scikit-beam</i> .
	Hasil	Memparalelkan perhitungan fungsi korelasi waktu dari perangkat lunak <i>XPCS</i> dalam paket <i>scikit-beam</i> menggunakan modul multiprosesor Python. Modul multiprosesor Python memperkenalkan cara untuk memparalelkan beban kerja.
	Keterkaitan	Konsep pemrograman paralel .
2	Judul	<i>Comparison of Web Scraping Techniques: Regular Expression, HTML DOM and XPath</i>
	Penulis	Rohmat Gunawan, Alam Rahmatulloh, Irfan Darmawan dan Firman Firdaus (2018)
	Penelitian	Membandingkan tiga metode <i>web scraping</i> dengan parameter pengukuran waktu, penggunaan memori dan penggunaan data
	Hasil	Metode <i>HTML DOM</i> lebih unggul pada parameter pengukuran waktu dan penggunaan data, sedangkan metode <i>Regex</i> unggul pada penggunaan memori.
	Keterkaitan	Teknik <i>web scraping</i> dan parameter yang digunakan.

3	Judul	<i>A review of programming languages for web scraping from software repository sites</i>
	Penulis	Mohan Prakash dan Dr. Ekbal Rashid (2017)
	Penelitian	Mereview empat bahasa pemrograman C, Java, PHP dan Python terkait perpustakaan dan metode penggunaannya pada <i>web scraping</i> dan <i>data extraction</i> .
	Hasil	Python merupakan bahasa terbaik yang dapat digunakan untuk <i>web scraping</i> . Hal ini karena adanya modul-modul yang memumpuni seperti <i>Scrapy</i> , <i>Selenium</i> , <i>Spiders</i> dan lain-lain.
	Keterkaitan	Bahasa pemrograman <i>python</i>
4	Judul	Implementasi <i>Web Scraping</i> dan <i>Text Mining</i> untuk Akuisisi dan Kategorisasi Informasi Laman Web Tentang Hidroponik
	Penulis	A Priyanto dan M R Ma'arif (2018)
	Penelitian	Mengumpulkan informasi pada laman web menggunakan <i>web scraping</i> dilanjutkan dengan melakukan pengelompokkan kedalam beberapa kategori menggunakan <i>text mining</i> .
	Hasil	Mengotomasi proses akuisisi informasi khususnya informasi-informasi yang bersumber dari artikel atau tulisan bebas di internet dan mengelompokkannya kedalam beberapa kategori.
	Keterkaitan	Proses <i>web scraping</i>
5	Judul	Perbandingan Metode <i>Web Scraping</i> Menggunakan <i>CSS Selector</i> dan <i>Xpath Selector</i>
	Penulis	Taufiq Rizaldi dan Hermawan Arief (2017)
	Penelitian	Membandingkan dua metode <i>web scraping</i> dengan parameter pengukuran waktu, penggunaan memori, penggunaan data dan jumlah data
	Hasil	Penggunaan metode <i>XPATH</i> untuk <i>web scraping</i> situs berita menghasilkan artikel yang lebih lengkap dibandingkan dengan menggunakan metode <i>CSS Selector</i> . Metode <i>XPATH</i> juga lebih unggul dalam pengukuran waktu.
	Keterkaitan	Teknik <i>web scraping</i> dan parameter pengujian
6	Judul	<i>An Overview On Web Scraping Techniques And Tools</i>
	Penulis	Anand V. Saurkar, Kedar G. Pathare dan Shweta A. Gode (2018)
	Penelitian	Menganalisis teknik <i>web scraping</i> dan <i>tools</i> yang digunakan

	Hasil	Ada banyak teknik <i>web scraping</i> yang bisa digunakan diantaranya <i>Classical copy and paste</i> , <i>Hypertext Transfer Protocol (HTTP) Programming</i> , <i>Hyper Text Markup Language (HTML) Parsing</i> , <i>Document Object Model (DOM) Parsing</i> , <i>Web Scraping Software</i> dan <i>Computer vision web-page analysers</i> dan alat-alat yang digunakan pun beragam diantaranya <i>Mozenda</i> , <i>Visual Web Ripper</i> , <i>Web Content Extractor</i> , <i>Import.io</i> dan <i>Scrapy</i>
	Keterkaitan	Teknik <i>web scraping</i> dan <i>library</i>
7	Judul	<i>DECO: Polishing Python Parallel Programming</i>
	Penulis	Alex Sherman, Peter Den Hartog (2016)
	Penelitian	Mengusulkan penyederhanaan teknik pemrograman paralel tradisional yang meminimalkan interaksi programmer dan tidak memerlukan pengetahuan pemrograman paralel
	Hasil	Penyederhanaan teknik pemrograman bersamaan yang ditargetkan pada pemrogram dengan sedikit pemahaman tentang pemrograman bersamaan
	Keterkaitan	Konsep <i>python parallel programming</i>
8	Judul	<i>Web Scraping and Naïve Bayes Classification for Job Search Engine</i>
	Penulis	C Slamet, R Andrian, D S Maylawati, Suhendar, W Darmalaksana dan M A Ramdhani (2018)
	Penelitian	Penyederhanaan pencarian kerja melalui konstruksi dan kolaborasi teknik pengikisan web dan klasifikasi menggunakan <i>Naïve Bayes</i> di mesin pencari
	Hasil	Menghasilkan aplikasi yang efektif dan efisien bagi pengguna untuk mencari pekerjaan potensial yang sesuai dengan minat mereka.
	Keterkaitan	Konsep <i>web scraping</i>
9	Judul	Analisis <i>Web Scraping</i> Untuk Data Bencana Alam Dengan Menggunakan Teknik <i>Breadth-First Search</i> Terhadap 3 Media Online
	Penulis	Izatul Putri Sonya dan Dr. Prihandoko, Skom (2016)
	Penelitian	Melakukan <i>scraping</i> pada data yang tidak terstruktur di beberapa media online dan mengklasifikasikannya menjadi data yang terstruktur menggunakan teknik <i>Breadth-First</i>
	Hasil	Menghasilkan data yang terstruktur berupa tabel dengan beberapa field yaitu no, hari/tanggal, waktu posting, judul, deskripsi, gambar, dan link halaman artikel.

	Keterkaitan	Konsep <i>web scraping</i>
10	Judul	<i>PyCOMPSs: Parallel computational workflows in Python</i>
	Penulis	Enric Tejedor, Yolanda Becerra, Guillem Alomar, Anna Queralt, Rosa M Badia, Jordi Torres, Toni Cortes dan Jesu's Labarta (2015)
	Penelitian	Melakukan penelitian terhadap permasalahan umum yang dialami pengguna dalam menggunakan pemrograman paralel pada <i>Python</i>
	Hasil	Menyajikan PyCOMPSs, sebuah kerangka kerja yang memfasilitasi pengembangan alur kerja komputasi paralel dalam Python.
	Keterkaitan	Konsep parallel programming
11	Judul	<i>Comparison Of Python Libraries Used For Web Data Extraction</i>
	Penulis	Erdoğan Uzun, Tarik Yerlikaya Dan Oğuz Kirat (2018)
	Penelitian	Membandingkan tiga pustaka ekstraksi terkenal yang berbeda termasuk <i>BeautifulSoup</i> , <i>lxml</i> dan <i>regex</i>
	Hasil	Hasil percobaan menunjukkan bahwa <i>regex</i> mencapai hasil terbaik dengan rata-rata 0,071 ms. Namun, sulit untuk menghasilkan aturan ekstraksi yang benar untuk <i>regex</i> ketika jumlah elemen dalam tidak diketahui, dan <i>lxml</i> menunjukkan hasil terbaik dengan rata-rata 9.074 ms
	Keterkaitan	Penggunaan <i>library web scraping</i> .
12	Judul	<i>Web Scraping</i>
	Penulis	Bo Zhao (2019)
	Penelitian	Menganalisa <i>web scraping</i> secara terperinci
	Hasil	Mendeskripsikan <i>web scraping</i> secara terperinci
	Keterkaitan	Konsep <i>web scraping</i> .

2.6. State of the Art

Beberapa penelitian sudah dilakukan berhubungan dengan penelitian yang sedang dilakukan, peneliti sudah memperluas *state of the art* dari parameter dan teknik yang digunakan pada penelitian sebelumnya, serta penggunaan *paralel programming* yang mampu mengeksekusi beberapa perintah secara bersamaan dengan lebih efisien.

Selain itu, para peneliti sudah menjawab beberapa pertanyaan yang dibutuhkan untuk memulai penelitian ini, dengan mengeksplorasi setiap kelebihan untuk menangkap parameter yang dibutuhkan, mengembangkan teknik-teknik yang sebelumnya telah digunakan agar mendapatkan hasil yang lebih efektif. Beberapa penelitian lainnya berfokus pada penggunaan teknik *web scraping*. Tantangan dan peluang pada penelitian ini yaitu penggunaan *web scraping* pada kasus yang berbeda seperti melakukan *scraping* pada beberapa *URL* bersamaan dan melakukan *scraping* dengan banyak objek data. Hal tersebut dapat dilihat pada matrik penelitian berikut:

Table II-2 Tabel Matriks Penelitian

1	Pengarang	Taufiq Rizaldi, Hermawan Arief
	Tahun	2017
	Judul	Perbandingan Metode Web Scraping Menggunakan CSS Selector dan XPath Selector
	Penjelasan	Penelitian yang dilakukan membandingkan teknik <i>web scraping CSS Selector</i> dan <i>XPath</i> pada <i>weblog</i> http://blog.detik.com dengan membagi tiga kategori yang digunakan untuk <i>scraping</i> yaitu, komunitas, hiburan dan kuliner. Bahasa pemrograman yang digunakan <i>python</i> dengan menggunakan <i>framework scrapy</i> . Parameter yang

		digunakan dalam pengujian yaitu pengukuran waktu, penggunaan memori dan jumlah item.
	Kesimpulan	Untuk penggunaan memori baik <i>CSS Selector</i> dan <i>XPath</i> tidak memiliki perbedaan yang begitu signifikan, sedangkan untuk waktu yang dibutuhkan <i>XPath</i> memiliki waktu pemrosesan yang lebih singkat dari <i>CSS Selector</i> .
	Kekurangan	Objek penelitian yang digunakan tidak diperjelas apakah dari setiap kategori yang dipilih hanya judul artikel saja yang diambil atau isi dari setiap konten yang muncul dari kategori tersebut. Penggunaan <i>framework scrapy</i> dirasa kurang tepat, karena <i>scrapy</i> berfungsi sebagai <i>engine</i> yang sudah mencakup tujuh komponen kompleks <i>Scheduler, Item Pipeline, Downloader, Downloader Middleware, Spiders, Spiders Middleware</i> sehingga akan mengabaikan sedikit teknik yang digunakan.
	Keterkaitan	Penggunaan teknik <i>web scraping</i> dan parameter
2	Pengarang	Rohmat Gunawan, Alam Rahmatulloh, Irfan Darmawan, Firman Firdaus
	Tahun	2018
	Judul	<i>Comparison of Web Scraping Techniques: Regular Expression, HTML DOM and XPath</i>
	Penjelasan	Penelitian yang dilakukan melakukan perbandingan teknik <i>web scraping</i> dengan teknik <i>HTML DOM, Regular Expression, dan XPath</i> pada website http://testing-ground.scraping.pro . Java menjadi bahasa pemrograman yang digunakan pada penelitian ini. Parameter yang digunakan pengukuran waktu, penggunaan data dan penggunaan memori.
	Kesimpulan	Kinerja <i>web scraping</i> metode <i>regular expression</i> paling kecil dalam penggunaan <i>memory</i> dibandingkan metode <i>HTML DOM, dan XPath</i> . Sedangkan <i>HTML DOM</i> membutuhkan waktu paling sedikit dan penggunaan data paling kecil dibanding metode <i>Regex dan XPath</i> .
	Kekurangan	Tidak adanya pengukuran objek yang didapatkan dan pengukuran <i>CPU</i> , tidak dijelaskan mengapa bahasa pemrograman yang digunakan <i>java</i> .
	Keterkaitan	Penggunaan teknik <i>web scraping</i> dan parameter