

NATIONAL TECHNICAL UNIVERSITY OF ATHENS School of Electrical and Computer Engineering Division of Computer Science

Advanced Web Scraping in the Modern Web

Techniques, Prevention, and AI Integration

DIPLOMA THESIS

of

ADONIS M. TSERIOTIS

Supervisor: Vassilios Vescoukis Professor, NTUA

Athens, February 2025



National Technical University of Athens School of Electrical and Computer Engineering Division of Computer Science

Advanced Web Scraping in the Modern Web

Techniques, Prevention, and AI Integration

DIPLOMA THESIS

of

ADONIS M. TSERIOTIS

Supervisor: Vassilios Vescoukis Professor, NTUA

Approved by the examination committee on 7th March 2025.

(Signature)

(Signature)

(Signature)

Vassilios Vescoukis Nikolaos Papaspyrou Zoe Paraskevopoulou Professor, NTUA Professor, NTUA Assistant Professor, NTUA



National Technical University of Athens School of Electrical and Computer Engineering Division of Computer Science

Copyright ©, Adonis M. Tseriotis, 2025 All rights reserved.

The copying, storage and distribution of this diploma thesis, exall or part of it, is prohibited for commercial purposes. Reprinting, storage and distribution for non - profit, educational or of a research nature is allowed, provided that the source is indicated and that this message is retained.

The content of this thesis does not necessarily reflect the views of the Department, the Supervisor, or the committee that approved it.

DISCLAIMER ON ACADEMIC ETHICS AND INTELLECTUAL PROP-ERTY RIGHTS

Being fully aware of the implications of copyright laws, I expressly state that this diploma thesis, as well as the electronic files and source codes developed or modified in the course of this thesis, are solely the product of my personal work and do not infringe any rights of intellectual property, personality and personal data of third parties, do not contain work / contributions of third parties for which the permission of the authors / beneficiaries is required and are not a product of partial or complete plagiarism, while the sources used are limited to the bibliographic references only and meet the rules of scientific citing. The points where I have used ideas, text, files and / or sources of other authors are clearly mentioned in the text with the appropriate citation and the relevant complete reference is included in the bibliographic references section. I fully, individually and personally undertake all legal and administrative consequences that may arise in the event that it is proven, in the course of time, that this thesis or part of it does not belong to me because it is a product of plagiarism.

(Signature)

Adonis M. Tseriotis Electrical & Computer Engineer Graduate, NTUA

7th March 2025

to my family

Περίληψη

Η παρούσα εργασία εξετάζει το web scraping, μια διαδικασία αυτοματοποιημένης εξαγωγής δεδομένων από ιστοσελίδες, εστιάζοντας στις τεχνικές, τις προκλήσεις και τις καινοτόμες λύσεις που καθορίζουν το σύγχρονο περιβάλλον συλλογής δεδομένων. Με την αυξανόμενη πολυπλοκότητα των διαδικτυακών τεχνολογιών και την εφαρμογή μηχανισμών αποτροπής αυτοματοποιημένης πρόσβασης, απαιτούνται προηγμένες στρατηγικές για την αποτελεσματική και ηθικά αποδεκτή συλλογή πληροφορίας.

Η εργασία αναλύει τις θεμελιώδεις μεθόδους scraping, όπως η ανάλυση HTML μέσω HTTP αιτημάτων, η χρήση headless browsers και η αναχαίτιση δικτυακών αιτημάτων, συγκρίνοντας τα πλεονεκτήματα και τα μειονεκτήματά τους. Παράλληλα, εξετάζονται οι τεχνικές προστασίας των ιστοσελίδων, όπως το browser fingerprinting, η ανάλυση μοτίβων κίνησης, οι CAPTCHA προκλήσεις και η απόκρυψη δεδομένων μέσω δυναμικών αποδόσεων περιεχομένου.

Η εισαγωγή της τεχνητής νοημοσύνης (AI) και των μεγάλων γλωσσικών μοντέλων (LLMs) στο web scraping αποτελεί μία από τις πιο καινοτόμες προσεγγίσεις. Η χρήση μηχανικής μάθησης επιτρέπει την ανάπτυξη αυτοπροσαρμοζόμενων εξαγωγέων δεδομένων, την αυτόματη αναγνώριση δομών δεδομένων και την έξυπνη παράκαμψη ανιχνευτικών μηχανισμών.

Ως πραχτική εφαρμογή, η εργασία παρουσιάζει την πλατφόρμα "soniq", ένα ανοιχτού κώδικα, no-code εργαλείο scraping που αξιοποιεί AI για την αυτοματοποίηση και βελτιστοποίηση της συλλογής δεδομένων. Η πλατφόρμα ενσωματώνει LLM-assisted schema inference, προγραμματισμένη εξαγωγή δεδομένων και προηγμένη διαχείριση proxies, επιτρέποντας στους χρήστες να πραγματοποιούν scraping χωρίς εξειδικευμένες τεχνικές γνώσεις.

Η εργασία καταλήγει σε μια συζήτηση για τις μελλοντικές προοπτικές στον χώρο του web scraping, εστιάζοντας στη δημιουργία ευφυών, ανθεκτικών και ηθικά αποδεκτών scraping pipelines, τη διασύνδεση με data warehouses για προηγμένη ανάλυση δεδομένων, και τη χρήση αποκεντρωμένων αρχιτεκτονικών για μεγαλύτερη ανωνυμία και ανθεκτικότητα. Τα ευρήματα αυτής της μελέτης αναδεικνύουν τη σημασία της τεχνολογικής καινοτομίας και της δημιοκρατικοποίησης της πρόσβασης στα δεδομένα, ενισχύοντας τη διαφάνεια και την αποτελεσματικότητα της συλλογής πληροφοριών στον σύγχρονο ψηφιακό κόσμο.

Λέξεις Κλειδιά

Web Scraping, Εξαγωγή Δεδομένων, Αντιμετώπιση Bots, Μηχανική Μάθηση για Εξαγωγή Δεδομένων, Μεγάλα Γλωσσικά Μοντέλα για Εξαγωγή Δεδομένων, Τεχνητή Νοημοσύνη, No-Code Web Scraping, Προγραμματισμένη Εξαγωγή Δεδομένων, Ανίχνευση Bots, Αυτοπροσαρμοζόμενες Τεχνικές Εξαγωγής, Αντίμετρα Ανίχνευσης Scraping.

4

Abstract

The modern web presents both unprecedented opportunities and significant challenges for data extraction, as web technologies evolve to become more dynamic and resistant to automated scraping. Data are the new currency, playing significant role in everything from business intelligence to scientific research, with a particularly vital impact on AI training, where properly structured and high-quality datasets are essential for building accurate and reliable models. This thesis explores the landscape of web scraping, focusing on advanced techniques, countermeasures, and AI integration. It systematically examines the foundational methods for web scraping, detailing strategies for handling both static and dynamic web content while addressing the rising challenges posed by anti-scraping mechanisms such as CAPTCHAs, browser fingerprinting, and traffic pattern analysis.

A major focus of this work is to make web scraping more accessible to non-technical individuals, facilitated through the development of Soniq, an open-source, no-code web scraping platform that integrates LLM-assisted schema extraction, automated proxy management, and real-time adaptability to frontend changes. By leveraging machine learning for intelligent data extraction and generative AI for structured data understanding, this system enables accessible, scalable, and resilient scraping for diverse applications, from market analysis to AI model training.

The findings of this thesis illustrate the transformative impact of LLMs and generative AI in automating, optimizing, and scaling web scraping workflows. Additionally, the study provides insights into ethical considerations and compliance frameworks necessary for responsible data extraction in an era of increasing regulatory scrutiny and cybersecurity challenges. The proposed solutions not only enhance efficiency and accuracy in data extraction but also contribute to the broader effort of making structured web data available to a wider audience beyond large corporations and specialized developers.

Keywords

Web Scraping, AI in Data Extraction, LLM-Assisted Web Scraping, Scraping Prevention, No-Code Web Scraping, Structured Data Extraction, Anti-Scraping Countermeasures, Generative AI.

Acknowledgements

I would like to thank my professor, Vassilios Vescoukis, for his guidance and support throughout this thesis. His feedback and advice helped me a lot.

I also appreciate my supervising Ph.D. candidate, Mr. Christos Hadjichristofi, for always being available to help me and for giving me useful suggestions along the way.

A big thank you to my family and friends for always supporting me and standing by me through this process.

Athens, February 2025

Adonis M. Tseriotis

Table of Contents

Π	ερίλι	ηψη		3
A	bstra	ıct		5
\mathbf{A}	ckno	wledge	ments	7
0	Eхт	τεταμέ	νη Περίληψη στα Ελληνικά	17
	0.1	Εξέλιξ	η και Προκλήσεις του Web Scraping	17
	0.2	Τεχνικ	κές Web Scraping και Σύγκρισή τους	17
	0.3	Μηχαι	ησμοί Προστασίας και Ανίχνευσης Scraping	18
	0.4	Η Εφα	φμογή της Τεχνητής Νοημοσύνης στο Web Scraping	19
	0.5	Η Πλα	τφόρμα "soniq" – No-Code Scraping με ΑΙ	19
		0.5.1	Βασικά Χαρακτηριστικά της Πλατφόρμας	20
Ι	Ba	ckgrou	ind Knowledge	23
1	Mo	dern V	Veb Scraping Architecture	25
	1.1	Introd	uction	25
	1.2	Evolut	ion of Web Technologies	26
	1.3	Advan	ced Web Technologies	27
		1.3.1	Modern Javascript Frameworks	27
		1.3.2	Advanced API Interfaces	36
	1.4	Search	Engine Optimization (SEO) $\ldots \ldots \ldots$	37
2	Tec	hnique	s and Innovations in Web Scraping	39
	2.1	Found	ational Scraping Techniques	39
		2.1.1	Scraping Static Websites	39
		2.1.2	Dynamic JavaScript-Driven Websites	41
		2.1.3	Network Request Interception	42
		2.1.4	Comparing Techniques	44
3	Scr	aping l	Prevention and Countermeasures	47
	3.1	Detect	ion and Mitigation Strategies	47
		3.1.1	Browser Fingerprinting	47
		3.1.2	Traffic Pattern Analysis	48

		3.1.3	САРТСНАя	48
		3.1.4	Honeypotting \ldots	49
		3.1.5	IP Reputation Systems	50
		3.1.6	Behavioral Analytics	50
		3.1.7	$Conclusion . \ . \ . \ . \ . \ . \ . \ . \ . \ .$	50
	3.2	Prever	ntive Development Techniques	52
		3.2.1	Rate Limiting and Throttling	52
		3.2.2	Dynamic Code and CSS Attributes	52
		3.2.3	API Key and Token-Based Authentication	53
		3.2.4	JavaScript Challenges	54
		3.2.5	Cookie-Based Authentication	55
		3.2.6	CDN Security	56
		3.2.7	Data Obfuscation Techniques	56
		3.2.8	Adaptive User Interface Rendering	57
	3.3	Advan	ced Protection Mechanisms	59
		3.3.1	Anti-Scraping SaaS Platforms	59
		3.3.2	Browser Integrity and Verification	59
		3.3.3	Hybrid Defense Architectures	60
4	Inte	ogratio	n of AI in Web Scraping	63
-	4.1	Machi	ne Learning in Scraping	63
		4.1.1	Adaptive Algorithms for Intelligent Data Extraction	63
		4.1.2	Automated Proxy Management	64
		4.1.3	CAPTCHA Solving with AI	65
	4.2	NLP f	or Data Understanding	66
5	Lev	eragin	g Generative AI for Scraping	69
	5.1			
		LLMs	in Scraping Pipelines	69
		LLMs 5.1.1	in Scraping Pipelines	69 70
		LLMs 5.1.1 5.1.2	in Scraping Pipelines	69 70 70
		LLMs 5.1.1 5.1.2 5.1.3	in Scraping Pipelines	69 70 70 71
		LLMs 5.1.1 5.1.2 5.1.3 5.1.4	in Scraping Pipelines	69 70 70 71 71
	5.2	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan	in Scraping Pipelines	 69 70 70 71 71 72
	5.2	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1	in Scraping Pipelines	 69 70 70 71 71 72 72
	5.2	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1 5.2.2	in Scraping Pipelines	 69 70 70 71 71 72 72 72 72
	5.2 5.3	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1 5.2.2 Use C	in Scraping Pipelines	 69 70 70 71 71 72 72 72 73
	5.2 5.3	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1 5.2.2 Use C 5.3.1	in Scraping Pipelines	 69 70 70 71 71 72 72 72 73 73 73
	5.2 5.3	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1 5.2.2 Use C 5.3.1 5.3.2	in Scraping Pipelines	 69 70 70 71 71 72 72 73 73 74
	5.2 5.3	LLMs 5.1.1 5.1.2 5.1.3 5.1.4 Advan 5.2.1 5.2.2 Use C 5.3.1 5.3.2 5.3.3	in Scraping Pipelines	 69 70 70 71 71 72 72 73 73 74 74

77

II Implementation

_			
6	soni	q: No-	code web scraping platform for structured data extraction 79
	6.1	Proble	m definition $\ldots \ldots .$ 79
	6.2	Techno	plogies Used
		6.2.1	Backend
		6.2.2	Frontend
	6.3	Archit	ecture
		6.3.1	Frontend Layer
		6.3.2	Backend Layer
		6.3.3	Deployment Architecture
		6.3.4	System Workflow
	6.4	Use Ca	ase - Scraping Energy News Articles in Seconds
		6.4.1	Creating the Ontology
		6.4.2	Creating a Page
		6.4.3	Inspecting Job Executions
		6.4.4	Scraped Data Exploration

105
. 105
. 105
. 106
. 106
. 107
. 107
. 108

Bibliography	112
List of Abbreviations	113

List of Figures

1.1	Global number of internet users $2005-2024[1]$	28
1.2	Angular Usage Statistics [2]	29
1.3	Component Based Architecture UML Component Diagram	30
1.4	React Usage Statistics [2]	30
1.5	Vue Usage Percentage [2]	32
1.6	SSR Sequence Diagram	33
1.7	ISR Sequence Diagram	33
1.8	SSG Sequence Diagram	34
1.9	Svelte Usage Statistics [2]	35
		~ ~
6.1	soniq Component UML Diagram	80
6.2	soniq API Swagger Docs 1	89
6.3	soniq API Swagger Docs 2	90
6.4	soniq Database UML Class Diagram	91
6.5	soniq Deployment UML Diagram	92
6.6	soniq Activity UML Diagram	93
6.7	soniq Sequence UML Diagram	94
6.8	Ontology Creation	95
6.9	Basic Page Information	95
6.10	Inspect HTML Container	96
6.11	CSS Extraction Schema Generation using LLM \ldots	96
6.12	Modified Generated Schema	97
6.13	Scraping Job Simulation	98
6.14	Scraping Job Scheduling	98
6.15	Page Explorer	99
6.16	energia.gr Pagination Url	99
6.17	Pagination Url Enqueuer	00
6.18	Run Explorer	.00
6.19	Job Execution Details	.01
6.20	Data Explorer	.01

List of Tables

1.1	Angular Usage Percentage	29
1.2	React Usage Percentage	31
1.3	Vue Usage Percentage	32
1.4	Vue Usage Percentage	35
2.1	Comparison of Foundational Scraping Techniques	45
3.1	Comparison of Web Scraping Bot Detection Methods	51
3.2	Comparison of Preventive Development Techniques	58

Chapter 0

Εκτεταμένη Περίληψη στα Ελληνικά

Η εργασία αυτή αποτελεί μια ολοκληρωμένη μελέτη του web scraping, μιας διαδικασίας που επιτρέπει την αυτόματη εξαγωγή δεδομένων από το web, με έμφαση στις τεχνικές, τα εμπόδια και τις καινοτόμες λύσεις που έχουν αναπτυχθεί τα τελευταία χρόνια. Καθώς η πολυπλοκότητα των ιστοσελίδων αυξάνεται και οι μηχανισμοί προστασίας από την αυτοματοποιημένη πρόσβαση βελτιώνονται, η ανάγκη για προηγμένες, προσαρμόσιμες και ανθεκτικές τεχνικές scraping είναι πιο επιτακτική από ποτέ.

Η μελέτη ξεκινά από τις θεμελιώδεις τεχνικές εξαγωγής δεδομένων, αναλύει τους μηχανισμούς ανίχνευσης και αποκλεισμού bots, εξετάζει πώς η τεχνητή νοημοσύνη (AI) μπορεί να μετασχηματίσει το πεδίο, και ολοκληρώνεται με την παρουσίαση του "soniq", μιας no-code, AI-powered πλατφόρμας scraping, σχεδιασμένης για να διευκολύνει την πρόσβαση σε δομημένα δεδομένα χωρίς την ανάγκη τεχνικής εξειδίκευσης.

0.1 Εξέλιξη και Προκλήσεις του Web Scraping

Στις αρχές του διαδικτύου, η δομή των ιστοσελίδων ήταν στατική, βασισμένη σε HTML και CSS, καθιστώντας την εξαγωγή δεδομένων απλή και απρόσκοπτη. Ωστόσο, με την εισαγωγή των μοντέρνων JavaScript frameworks όπως React, Angular και Vue.js, πολλές ιστοσελίδες μετατράπηκαν σε δυναμικές εφαρμογές μονής σελίδας (SPAs), όπου το περιεχόμενο φορτώνεται ασύγχρονα μέσω API κλήσεων και JavaScript rendering.

Αυτή η αλλαγή δημιούργησε σημαντικές δυσκολίες στο scraping, καθώς οι πληροφορίες δεν ήταν πλέον διαθέσιμες στον αρχικό HTML κώδικα. Επιπλέον, η εμφάνιση νέων τεχνολογιών επικοινωνίας μεταξύ πελάτη-διακομιστή, όπως το GraphQL και το gRPC, επέφερε νέα εμπόδια, καθώς οι παραδοσιακές τεχνικές εξαγωγής δεδομένων που βασίζονταν σε REST APIs δεν μπορούσαν να εφαρμοστούν άμεσα. Παράλληλα, ο ρόλος του SEO (Search Engine Optimization) στη δομή και παρουσίαση των δεδομένων επηρέασε σημαντικά τις δυνατότητες συλλογής πληροφοριών, δεδομένου ότι πολλές ιστοσελίδες τροποποιούν τη δημόσια ορατότητα του περιεχομένου τους για να βελτιώσουν την κατάταξή τους στις μηχανές αναζήτησης.

0.2 Τεχνικές Web Scraping και Σύγκρισή τους

Οι μέθοδοι και οι τεχνικές της συστηματικής εξαγωγής δεδομένων απο το διαδίκτυο διαφέρουν με βάση την φύση της κάθε ιστοσελίδας. Συνήθως, η τεχνική είναι ένας αγώνας να λάβεις με χάποιο τρόπο το HTML περιεχόμενο της ιστοσελίδας που περιέχει τα δεδομένα ή να υποχλέψεις την χλήση διχτύου που μέσω αυτής παρέχονται τα δεδομένα στην ιστοσελίδα για προβολή. Στην περίπτωση των στατιχών σελίδων μπορούμε να υλοποιήσουμε το HTML Request Parsing, όπου ο HTML κώδικας μας παρέγεται απο μια απλή κλήση δικτύου τύπου GET και τα δεδομένα εξάγονται από το ακατέργαστο HTML μέσω CSS selectors και XPath queries. Αυτή η μέθοδος είναι γρήγορη και αποτελεσματική για στατικές ιστοσελίδες αλλά αποτυγχάνει σε JavaScript-heavy εφαρμογές. Στις περιπτώσεις όπου η ιστοσελίδα βασίζεται στη Javascript για να φορτώσει το περιεχόμενο της μια απλή κλήση δικτύου δεν αρκεί για να λάβουμε τον HTML χώδιχα. Επομένως, εισάγονται οι Headless Browsers, όπου εργαλεία προορισμένα για application testing όπως Puppeteer και Playwright προσομοιώνουν την ανθρώπινη αλληλεπίδραση με μια ιστοσελίδα, έχοντας την δυνατότητα να τρέξουν την JavaScript και αλληλεπιδρώντας με το UI. Αυτή η τεχνική είναι πιο ευέλικτη, αλλά απαιτεί υψηλή υπολογιστική ισχύ και μπορεί να ανιχνευθεί μέσω browser fingerprinting και άλλων μεθόδων ανίχνευσης αυτοματοποιημένων προγραμμάτων. Ωστόσο, σε ιστοσελίδες οι οποίες βασίζονται εκτεταμένα στις κλήσεις δικτύου για να λάβουν τα δεδομένα που προβάλλουν μπορούμε να παραχάμψουμε την ανάγχη για HTML retrieval χαι parsing υλοποιώντας τη μέθοδο του Network Interception. Συγχεχριμένα, αυτή η μέθοδος εχμεταλλεύεται τις API χλήσεις μιας ιστοσελίδας για την άμεση απόσπαση δεδομένων. Παρόλο που αυτή η τεχνική είναι αποτελεσματική, οι περισσότερες ιστοσελίδες χρησιμοποιούν authentication tokens και encryption για να αποτρέψουν τη μη εξουσιοδοτημένη πρόσβαση στα ΑΡΙ τους. Μετά από συγκριτική ανάλυση, διαπιστώνεται ότι ο συνδυασμός αυτών των τεχνιχών αποτελεί την πιο αποδοτιχή στρατηγική, προσαρμοσμένη στις ανάγκες του εκάστοτε έργου.

0.3 Μηχανισμοί Προστασίας και Ανίχνευσης Scraping

Οι δυσκολίες εξαγωγής δεδομένων απο την εκάστοτε ιστοσελίδα δεν μένουν μόνο στη δυσκολία αναγνώρισης της κατάλληλης μεθόδου και υλοποιώντας την. Οι σύγχρονες ιστοσελίδες εφαρμόζουν εξελιγμένες τεχνικές αποτροπής scraping για να προστατεύσουν τα δεδομένα και τη λειτουργικότητα της σελίδας τους απο κακόβουλα αυτοματοποιημένα bots. Οι τεχνικές αποτροπής scraping που αναλύονται είναι οι εξής:

- Browser Fingerprinting, όπου συλλέγονται δεδομένα για τα χαραχτηριστικά του χρήστη (user agent, installed plugins, time zone) ώστε να εντοπιστούν ανωμαλίες που υποδηλώνουν bot.
- **Traffic Pattern Analysis**, όπου αναλύεται η συμπεριφορά του επισκέπτη (ρυθμός αιτημάτων, χρονικά διαστήματα) για την ταυτοποίηση αυτοματοποιημένης πρόσβασης.
- CAPTCHAs και Honeypots, που λειτουργούν ως δοκιμασίες για τη διαφοροποίηση των ανθρώπινων χρηστών από bots.
- IP Reputation Systems, που αποκλείουν γνωστές IP διευθύνσεις που σχετίζονται με scraping δραστηριότητες.

- Rate Limiting and Throttling, όπου το σύστημα διαχομιστή αναγνωρίζει και απορρίπτει αλλεπάλληλες κλήσεις που έχουν γίνει απο την ίδια διεύθυνση IP σε μικρό χρονικό διάστημα.
- Dynamic Code and CSS Attributes, όπου κατα την υλοποίηση ο προγραμματιστής επιλέγει τη συσκότιση του κώδικα και των CSS attributes κάνοντας δύσκολη την εύρεση ενός σταθερού μοτίβου για εξαγωγή δεδομένων.
- API Key, Token-Based and Cookie-Based Authentication, η συστηματική θωράκιση της ιστοσελίδας χρησιμοποιώντας διαδεδομένες τεχνικές ασφάλειας μπορεί να προβεί σε χαρακτηριστική μεθοδο αποτροπής scraping καθώς τα δεδομένα κάθε χρήστη είναι ιδιωτικά και απροσπέλαστα.
- Javascript Challenges, όπου ο browser του χρήστη δοχιμάζεται σε μια σειρά απο πολύπλοχες διεργασίες έτσι ώστε να επηρρεάσει αρνητιχά τους πόρους του client ετσι ώστε να χαταλάβει αν πρόχειται για αληθινό χρήστη ή ένα αυτοματοποιημένο πρόγραμμα που τρέχει σε διαχομιστές με περιορισμένους πόρους.

Συμπερασματικά, οι πιο αποτελεσματικές μέθοδοι αποτροπής scraping είναι ο συνδυασμός πολλών παραπάνω τεχνικών έτσι ώστε να γίνει μια διεργασία πολύ δύσκολη, η οποία θα αποθαρρύνει κακόβουλο λογισμικό απο το να προσπαθεί να τα προσπελάσει.

0.4 Η Εφαρμογή της Τεχνητής Νοημοσύνης στο Web Scraping

Η ενσωμάτωση της τεχνητής νοημοσύνης στο web scraping αποτελεί έναν από τους πιο καινοτόμους τομείς έρευνας. Μεγάλα γλωσσικά μοντέλα (LLMs) μπορούν να χρησιμοποιηθούν για την αυτόματη αναγνώριση της δομής μιας ιστοσελίδας, την εξαγωγή σημασιολογικά σχετικών πληροφοριών και την κατανόηση του τρόπου με τον οποίο οργανώνονται τα δεδομένα. Η μηχανική μάθηση επιτρέπει τη δημιουργία αυτοπροσαρμοζόμενων συστημάτων scraping που μπορούν να ανιχνεύουν και να προσαρμόζονται σε αλλαγές του περιβάλλοντος μιας ιστοσελίδας. Οι τεχνικές αυτές καθιστούν το scraping πιο ανθεκτικό, καθώς το σύστημα μπορεί να μαθαίνει από τις αλλαγές και να προσαρμόζεται αυτόματα.

0.5 Η Πλατφόρμα "soniq" – No-Code Scraping με AI

Ως πρακτική εφαρμογή, η εργασία παρουσιάζει την πλατφόρμα "soniq", ένα ανοιχτού κώδικα, no-code εργαλείο scraping που επιτρέπει σε χρήστες χωρίς τεχνικές γνώσεις να δημιουργούν, να προγραμματίζουν και να διαχειρίζονται διαδικασίες εξαγωγής δεδομένων από ιστοσελίδες με ευκολία και ακρίβεια. Σε αντίθεση με τα παραδοσιακά εργαλεία scraping, τα οποία απαιτούν εξειδικευμένη γνώση προγραμματισμού και συνεχή τροποποίηση των κώδικα για να προσαρμόζονται στις αλλαγές των ιστοσελίδων, το "soniq" χρησιμοποιεί τεχνητή νοημοσύνη και μηχανισμούς αυτοματοποίησης για την απλούστευση της διαδικασίας.

0.5.1 Βασικά Χαρακτηριστικά της Πλατφόρμας

Η πλατφόρμα ενσωματώνει προηγμένα χαρακτηριστικά, καθιστώντας την ευέλικτη, επεκτάσιμη και φιλική προς τον χρήστη.

- Αυτόματη Αναγνώριση Σχημάτων Δεδομένων (LLM-powered Schema Inference) Η πλατφόρμα αξιοποιεί μεγάλα γλωσσικά μοντέλα (LLMs) για την αναγνώριση και κατανόηση της δομής δεδομένων μιας ιστοσελίδας. Αντί να απαιτείται χειροκίνητος προσδιορισμός των CSS Selectors ή XPath κανόνων, η AI αναλύει το περιεχόμενο και προτείνει τα βέλτιστα πεδία δεδομένων προς εξαγωγή. Αυτό μειώνει δραστικά το χρόνο προετοιμασίας και καθιστά το scraping πιο ανθεκτικό σε αλλαγές UI.
- Σύστημα Χρονοπρογραμματισμού και Αυτοματοποιημένης Εκτέλεσης Το AP-Scheduler επιτρέπει στους χρήστες να προγραμματίζουν scraping εργασίες βάσει προκαθορισμένων χρονικών διαστημάτων. Οι χρήστες μπορούν να ορίσουν ημερήσιες, εβδομαδιαίες ή προσαρμοσμένες ανανεώσεις των δεδομένων που εξάγουν, καθιστώντας την πλατφόρμα ιδανική για επαναλαμβανόμενες εργασίες παρακολούθησης πληροφοριών (π.χ. δυναμικές τιμές προϊόντων, αναλύσεις αγοράς, ενημερωτικά άρθρα).

Αρχιτεκτονική και Τεχνολογίες της Πλατφόρμας Το "soniq" έχει αναπτυχθεί με σύγχρονες τεχνολογίες, διασφαλίζοντας υψηλή απόδοση, ασφάλεια και ευκολία χρήσης.

- Backend: Αναπτύχθηκε με FastAPI, ένα ελαφρύ αλλά ισχυρό framework για ταχύτερη επεξεργασία API αιτημάτων και ενσωμάτωση AI μοντέλων.
- Database: Χρησιμοποιεί MongoDB, επιτρέποντας αποδοτική αποθήκευση και ανάκτηση δομημένων δεδομένων με ευελιξία.
- Task Scheduling: Η υλοποίηση του APScheduler δίνει τη δυνατότητα προγραμματισμένης εκτέλεσης εργασιών scraping, διασφαλίζοντας σταθερή ροή δεδομένων χωρίς την ανάγκη χειροκίνητης παρέμβασης.
- Frontend: Χρησιμοποιεί React και Refine.dev για τη δημιουργία φιλικού προς τον χρήστη interface, μέσω του οποίου οι χρήστες μπορούν να δημιουργούν scraping εργασίες χωρίς την ανάγκη κώδικα.

Σύγκριση με Άλλα Εργαλεία Scraping Η εφαρμογή επιδιώκει να ξεχωρίσει από τα υπάρχοντα scraping εργαλεία καθώς:

- Απλοποιεί τη διαδικασία εξαγωγής δεδομένων, καταργώντας την ανάγκη προγραμματισμού.
- 2. Προσαρμόζεται δυναμικά στις αλλαγές ιστοσελίδων, μέσω LLM-powered schema inference.
- 3. Μειώνει τον κίνδυνο ανίχνευσης και αποκλεισμού, χάρη στην έξυπνη proxy rotation και human-like traffic generation.

 Ενσωματώνει προγραμματισμένες εξαγωγές δεδομένων, καθιστώντας το κατάλληλο για παρακολούθηση δυναμικού περιεχομένου σε μεγάλες κλίμακες.

Χρήσεις και Πιθανές Εφαρμογές Το "soniq" μπορεί να χρησιμοποιηθεί σε διάφορους τομείς όπου η συλλογή δομημένων δεδομένων είναι απαραίτητη:

- 1. Ανάλυση Αγοράς: Παρακολούθηση τιμών προϊόντων, τάσεων αγοράς και συμπεριφοράς καταναλωτών.
- 2. Δημοσιογραφία & Έρευνα: Αυτόματη συλλογή ειδήσεων και άρθρων από πολλαπλές πηγές.
- 3. Εκπαίδευση AI Μοντέλων: Απόκτηση μεγάλων όγκων δεδομένων για εκπαίδευση νευρωνικών δικτύων.
- Διαχείριση Επιστημονικών Δεδομένων: Συλλογή πληροφοριών από ακαδημαϊκές βάσεις δεδομένων.

Συμπεράσματα και Μελλοντικές Βελτιώσεις Το "soniq" αποτελεί ένα καινοτόμο no-code εργαλείο web scraping, το οποίο διευκολύνει τη συλλογή δεδομένων για ερευνητές, επιχειρήσεις και αναλυτές χωρίς να απαιτεί εξειδικευμένες τεχνικές γνώσεις. Η αυτοματοποίηση των διαδικασιών εξαγωγής, η χρήση AI για την αναγνώριση δομών δεδομένων και η διαχείριση proxies καθιστούν την πλατφόρμα ανθεκτική και προσαρμόσιμη στις συνεχείς αλλαγές του διαδικτύου.

Μελλοντικές βελτιώσεις περιλαμβάνουν τη διασύνδεση της πλατφόρμας με data warehouses, την ανάπτυξη μοντέλων reinforcement learning για αυτοβελτίωση των τεχνικών εξαγωγής. Το "soniq" δεν είναι απλώς ένα εργαλείο scraping, αλλά ένα βήμα προς την αύξηση της πρόσβασης στα δεδομένα, καθιστώντας τα διαθέσιμα, οργανωμένα και προσβάσιμα σε όλους.

Part I

Background Knowledge

Chapter 1

Modern Web Scraping Architecture

1.1 Introduction

Web scraping is the practice of programmatically extracting data from web pages, which has become indispensable in the digital age. Its relevance spans industries and disciplines, from data aggregation to provide market data analysis [3] to training data sets for machine learning models [4]. However, with the increasing complexity of modern web technologies, web scraping has turned out to be a more difficult and sophisticated task, requiring deep understanding of the underlying architectural frameworks that power today's web.

Web technologies have changed a lot in the last decade. From static content being delivered through server-rendered HTML to SPAs and real-time data streaming through WebSocket connections, the web has now become dynamic and interactive. These changes bring more richness to user interactions, but also make the extraction of structured data from websites more complicated. Consequently, web scraping, which was once straightforward, has now become an elaborate task that must adapt to these changes in technology.

This evolution has been defined, among other things, by the use of heavy JavaScript frameworks like React, Vue.js, and Angular. Such libraries apply techniques such as Virtual DOM and client-side rendering that improve user experience but, very often, make traditional scraping methods a lot harder because they change page content dynamically after loading. In addition, GraphQL and gRPC have replaced simpler RESTful interfaces in many applications and introduce other challenges while accessing and interpreting data.

Starting below, in this chapter, one will discuss contemporary web technologies that have changed the complexion of web scraping analyzing their key architecture ingredients, discussing every challenge for the extraction of data through them, together with the solutions that could address or defeat these challenges. By understanding them, we actually lay the foundational elements to strongly implement some agile scraping techniques, which have been elaborated on in the following chapters.

1.2 Evolution of Web Technologies

The evolution of web development has been a dynamic interaction between technological advances and changing user needs. From its inception in the early 1990s, the web has undergone a lot of transformations, each era defined by its prevailing methodologies, tools, and user expectations.

- Early Web Development: Text and Hyperlinks The early 1990s marked the birth of the World Wide Web, characterized by static web pages made up mostly of plain text and hyperlinks. It was all about information distribution, with very little interactivity, let alone visual appeal. This was the period when basic HTML dominated the standard for structuring content, enabling users to navigate between pages through hyperlinks. At that time, the web was still an emerging technology, mostly available to academic and research institutions.
- Mid-1990s: The Rise of Client-Side Scripting The mid-1990s witnessed the emergence of JavaScript, which was a game-changing client-side scripting language. Initially used to add interactivity to web pages without communicating with the server, JavaScript enabled developers to add functionality like form validation and simple animations. This was the first step toward more interactive web experiences and set the stage for dynamic, user-centric applications.
- Early 2000s: Dynamic Content with Server-Side Scripting In the early 2000s one can see the introduction of server-side scripting technologies such as PHP, ASP, and JSP. These technologies allowed servers to dynamically create HTML content in response to user actions or database interactions. For the first time, the web was able to become personalized and responsive, as web applications could now provide dynamic content to users in real time. This expanded the role of web applications and became an fundamental part of e-Commerce, social networking, and content management systems.
- Mid-2000s: The Introduction of AJAX The introduction of AJAX in the mid-2000s revolutionized how interactions on the web were performed. AJAX allowed web pages to asynchronously fetch data from servers, enabling pieces of a page to be updated without requiring a full reload. This greatly improved the user experience by providing smoother and faster interactions. The adoption of AJAX braced the creation of more dynamic and responsive web applications, clearing the way for modern web interfaces.
- 2010s: The Emergence of SPAs In the 2010s, Single-Page Applications were made possible with frameworks like AngularJS, React, and Vue.js. SPAs dynamically loaded content within a single page, enhancing performance and user experience by avoiding page reloads and enabling smoother interactions. SPAs leveraged the power of JavaScript and AJAX to create an era of highly interactive and fascinating web applications.

- Mid-to-late 2010s: The Comeback of SSR During the mid to late 2010s, SSR methodologies resurfaced. This was set in motion by frameworks like Next.js. These newer methodologies were implementing the classic server-rendering techniques with modern improvements tackling problems related to search engine optimization and performance of the first load. Notable developments during this time included:
 - Techniques of Server Rendering
 - **SSR** Server-side rendering of HTML, which gave a boost to SEO [1.4] and reduced initial load times.
 - **ISR** It took a hybrid approach by allowing the updating of static pages on demand at runtime and used the best of both worlds between static and dynamic rendering.
 - Hydration
 - This involved serving static HTML from the server and then making it interactive with the help of client-side JavaScript. Hydration was necessary to bridge both static and dynamic web experiences without sacrificing the benefit of faster initial rendering while retaining its interactivity.

This retrospective [5] underscores this transition of the web from static to dynamic, interactive, and intelligent platforms. The roots of basic evolution outlined above hold great importance as a background to present sophisticated web technologies.

1.3 Advanced Web Technologies

Over the past twenty years, the Internet has become a crucial component of daily life, from work-related applications to social media and news. Around the last decade Google, Facebook and other tech leaders, in their efforts to modernize their web pages, developed certain frameworks and protocols that were later released and saw great adoption in the community. These included Angular and React for making Google and Facebook more dynamic, adding interactivity, and overall improving user experience, as well as gRPC and GraphQL which revolutionized the way a client communicates with the server.

1.3.1 Modern Javascript Frameworks

At their essence, frameworks and libraries serve as a method to separate fundamental boilerplate code from the logic of an application. The wheel stops getting reinvented on every web page creation by standardizing web development practices reducing the development time and costs while enhancing cross-browser compatibility and improving performance and security. At their core, the frameworks that will be discussed in this section introduced efficient data binding, component-based architecture, and virtual DOM manipulation.



Figure 1.1. Global number of internet users 2005-2024[1]

Angular

Originally developed by Miko Hevery in 2010 and maintained by Google and the community, AngularJS is a web framework that aims to simplify both the development and testing of web applications by providing a framework for client-side model-view-controller (MVC) architectures, along with components commonly used in web applications [6]. The initial iterations of AngularJS contained numerous flaws, prompting Google developers to completely rewrite the framework and rebrand it as Angular (dropping the JS). Some notable features that might impact web scraping techniques in Angular-based web pages are:

- **Component-Based Architecture** Angular uses a component-based architecture, which allows developers to build encapsulated, reusable user interface elements. Each component encapsulates its own HTML, CSS, and TypeScript, making it easier to manage and test individual pieces of an application.
- **Directives** Angular extends HTML with additional attributes called directives. Directives offer functionality to change the behavior or appearance of DOM elements.
- **SSR** Angular has official support for server-side rendering, which improves an application's load time and performance. Server-side rendering also enhances search engine optimization by making content more accessible to web crawlers.

Numerous websites are using Angular [7], but in the vast Web it still corresponds to only 0.5%.



Figure 1.2. Angular Usage Statistics [2]

	No of Live Websites	Percentage
Top 10k	433	4.33%
Top 100k	$3,\!296$	3.3%
Top 1m	$12,\!679$	1.27%
All Internet	$521,\!384$	0.05%

 Table 1.1. Angular Usage Percentage

React

React, designed by Facebook in 2013, "is basically a web framework that was mainly designed to address the performance issues in web applications. React uses virtual DOM that decides whether the component has to be reloaded or not based on the current state of the component and the changes that have occurred. This prevents the application from re-rendering unnecessarily. Apart from this, React also introduces one-way data flow which helps to control the flow of the data within the application which makes the tracking of the occurred easier and also simplifies the propagation and stability." [8] To emphasize how the React framework could affect web scraping in applications built with it, these features must be addressed.

Virtual DOM The DOM (Document Object Model)[9] is a crucial component in web development, as it divides into modules and executes the code. The standard practice of Javascript Frameworks is to update the DOM at once, which can negatively impact the application's performance. React pioneered the use of a virtual DOM, an exact copy of the real DOM, which gets updated first and then is used to find

the minimal changes that need to be made. While this surely improves performance, it can complicate the process of a web scraper, adding the need to be able to run Javascript in order to get the website's content.

Component-based architecture Similar to Angular, React uses components to structure the UI.



Figure 1.3. Component Based Architecture UML Component Diagram



Figure 1.4. React Usage Statistics [2]

React's adoption is fairly notable in the above statistics as nearly half of the most popular websites are using it.

	No of Live Websites	Percentage
Top 10k	4,246	42.46%
Top 100k	$34,\!952$	34.95%
Top 1m	197,181	19.72%
All Internet	$51,\!355,\!483$	4.5%

 Table 1.2.
 React Usage Percentage

Vue.js

Vue.js [10], created by Evan You in 2014, is a progressive JavaScript framework popular for its simplicity and flexibility. Designed to be incrementally adoptable, Vue.js can function as a lightweight library for enhancing parts of a web page or as a full-fledged framework for building complex Single-Page Applications (SPAs). Its balanced approach between performance and ease of use has contributed to its widespread popularity among developers.

Key features of Vue.js include:

- **Two-Way Data Binding** Vue.js provides a strong two-way data binding mechanism that seamlessly synchronizes the model and the view. This makes it particularly beneficial for developing dynamic interfaces, but for web scraping, it may require handling frequent updates to the DOM efficiently.
- **Component-Based Architecture** Like other modern frameworks, Vue.js emphasizes a component-based structure, allowing developers to wrap and reuse user interface elements. Scrapers can leverage this modularity to identify repetitive patterns in data structures in components.
- **Directives** Vue.js extends HTML with directives such as 'v-bind' and 'v-for', enabling developers to add dynamic behavior to the DOM. These directives, while powerful for interactivity, require advanced handling techniques for web scraping tools to interpret and extract rendered content effectively.
- Virtual DOM Similar to React, Vue.js uses a virtual DOM to optimize rendering performance. This virtual representation can streamline updates but adds complexity for scrapers, as they may need to execute JavaScript to fully load and capture dynamic content.
- **Extensibility** Vue.js is highly extensible with plugins, state management solutions like Vuex, and routing capabilities through Vue Router. These tools enhance development capabilities, but may present additional layers for scrapers to navigate, particularly when dealing with nested routes or complex states.

Vue.js achieves a compelling balance between simplicity and sophistication, making it suitable for a wide range of applications. For web scraping, understanding the dynamic nature of Vue and its modular architecture is crucial to implement effective strategies.



Figure 1.5. Vue Usage Percentage [2]

	No of Live Websites	Percentage
Top 10k	2,361	23.61%
Top 100k	20,058	20.06%
Top 1m	117,128	11.71%
All Internet	4,790,110	0.43%

Table 1.3.Vue Usage Percentage

Next.js

Next.js[11], introduced by Vercel in 2016, represents a significant leap forward in web development by merging modern performance-focused methodologies with traditional server-side rendering (SSR). It was designed to tackle issues related to performance, search engine optimization (SEO[12]), and development efficiency, making it a popular choice for building powerful and scalable web applications. Its feature set provides developers with tools to build static and dynamic web pages with enhanced performance and flexibility.

Key features of Next.js that impact web scraping include:

Server-Side Rendering (SSR) One of the core functionalities of Next.js is SSR, which enables the generation of HTML on the server for each request. This ensures that the content is available to users and search engines immediately, improving initial load times and SEO. For web scraping, this makes data extraction more accessible as the content is pre-rendered and available in the source HTML.


Figure 1.6. SSR Sequence Diagram

Incremental Static Regeneration (ISR) ISR allows developers to update static pages on demand at runtime. This feature combines the speed of static generation with the flexibility of dynamic content updates. For scrapers, ISR introduces challenges in determining when content is updated, as changes may occur asynchronously postinitial rendering.



Figure 1.7. ISR Sequence Diagram

- **Hydration** Next.js employs hydration, a process where the server-rendered HTML becomes interactive on the client side through JavaScript. While this enhances user interactivity, it requires web scrapers to execute JavaScript to capture fully interactive content.
- **API Routes** Next.js simplifies the creation of backend endpoints using its API routes, enabling developers to build serverless functions directly within the application. This segmentation often requires additional steps for web scrapers to identify and interact with these endpoints.
- **Static Site Generation (SSG)** SSG is another feature of Next.js that pre-renders pages at build time, ensuring faster delivery of static content. For web scraping, this typically simplifies the process, as static pages are readily available without needing JavaScript execution.



Figure 1.8. SSG Sequence Diagram

Next.js is a React framework that has become essential for modern web development, offering performance, scalability, and dynamic interactivity. For web scraping, understanding its rendering and update mechanisms is key.

Svelte

Svelte[13], introduced by Rich Harris in 2016, represents a paradigm shift in web development by shifting the workload from the browser to the build process. Unlike traditional JavaScript frameworks such as React and Angular, Svelte compiles components into highly efficient, plain JavaScript code during build time, resulting in faster runtime performance and smaller bundle sizes.

Key features of Svelte include:

- **Compile-Time Optimization** Svelte eliminates the need for a virtual DOM by compiling components into minimal JavaScript code during the build process. This approach reduces runtime overhead and delivers exceptional performance, making it appealing.
- **Reactive Declarations** Svelte's reactivity system allows developers to declare reactive variables directly in the code. This simplifies the state management process, but may require scrapers to handle dynamically updated elements with care.
- **Scoped Styles** With Svelte, styles are scoped to components by default, ensuring better encapsulation. For web scraping, this can introduce challenges in identifying and extracting content tied to dynamically styled elements.
- Minimal Framework Overhead By compiling to plain JavaScript, Svelte avoids the inclusion of a framework runtime, which contributes to its small bundle size. This efficiency can simplify scraping as the application is less dependent on heavy client-side JavaScript.

Svelte has quickly gained popularity for its simplicity and performance. Understanding its compile-time features and reactive nature is essential for effective web scraping strategies.



Figure 1.9. Svelte Usage Statistics [2]

	No of Live Websites	Percentage
Top 10k	168	1.68%
Top 100k	933	0.93%
Top 1m	4,121	0.41%
All Internet	262,802	0.026%

 Table 1.4.
 Vue Usage Percentage

Conclusion

Despite the popularity of modern JavaScript frameworks like React, Angular, and Vue.js, their adoption across the broader web remains surprisingly limited. Statistics indicate that while these tools might dominate the realm of high-profile web applications and tech-driven companies, a significant portion of websites still rely on simpler technologies or legacy systems. The rise of these frameworks also represents a fascinating full-circle in web development. Initially, server-side HTML dominated the web with its straightforward content delivery. Over time, as client-side interactivity and AJAX became more and more important, rendering was moved to the client to improve user experience. However, challenges such as search engine optimization (SEO), initial load times, and the demand for faster performance brought server-side rendering (SSR) back into focus, albeit in a more advanced form. Frameworks like Next.js have revitalized SSR by combining it with modern techniques such as hydration and incremental static regeneration, achieving a balance between performance and interactivity.

1.3.2 Advanced API Interfaces

The evolution of web communication technologies has led to the development of advanced API interfaces, which have revolutionized the way clients and servers interact. These technologies include the foundational REST API as well as more recent innovations such as GraphQL and gRPC, each of which addresses specific limitations of its predecessors while introducing new paradigms for data exchange.

REST API

Representational State Transfer (REST) was introduced by Roy Fielding in 2000 as part of his doctoral dissertation[14]. REST defines a set of principles for designing networked applications, focusing on stateless communication, resource-based structure, and standard HTTP methods such as GET, POST, PUT, and DELETE. Its simplicity and wide adoption have made REST the default standard for APIs over the past two decades.

From a web scraping perspective, REST APIs offer a predictable and structured way to access data. Public and private APIs are often utilized to scrape specific endpoints to extract information. However, reliance on fixed endpoints can sometimes limit its flexibility, particularly in dynamic applications that require tailored queries or efficient batch operations.

GraphQL

GraphQL[15], developed by Facebook in 2012 and open-source in 2015, implemented as a solution to the challenges posed by the rigid structure of REST. Unlike REST, GraphQL allows clients to specify the structure of the response, enabling them to fetch exactly the data they need in a single query. This eliminates the problem of overfetching or underfetching data, which is common with REST APIs.

For web scraping, GraphQL presents both opportunities and challenges. Its query language allows scrapers to tailor requests for precise data extraction. However, its complexity, including nested and deeply interconnected queries, often requires advanced handling. In addition, limiting the rate and enforcing the schema can cause problems to automated scraping tools.

gRPC

gRPC[16], short for Google Remote Procedure Call, was created by Google in 2015 as an open-source framework for high-performance cross-platform communication. Unlike REST and GraphQL, which rely on human-readable formats such as JSON, gRPC uses Protocol Buffers (protobuf) for compact binary serialization. This makes gRPC exceptionally efficient for real-time communication and large-scale applications.

From a web scraping perspective, gRPC introduces significant challenges. The binary format of Protocol Buffers is not inherently human-readable, making it difficult for traditional scraping tools to interpret. In addition, reliance on bidirectional streaming and persistent connections can complicate data extraction efforts. However, for scrapers equipped with the necessary decoding and streaming tools, gRPC efficiency can be harnessed for high-performance data retrieval.

Implications for Web Scraping

These advanced API interfaces reflect the ongoing effort to optimize client-server interactions. Although REST remains widely used due to its simplicity, GraphQL and gRPC provide powerful alternatives for specific use cases. For web scraping, the choice of API interface significantly affects the complexity and efficiency of data extraction. Understanding the nuances of each technology is critical for designing scraping strategies that balance performance, precision, and adaptability.

1.4 Search Engine Optimization (SEO)

Search Engine Optimization, or SEO, is the process of improving a website's visibility on search engine result pages to increase organic traffic. By structuring a website's structure, content, and metadata in harmony with algorithms which search engines like Google use, SEO secures high ranking relevant to queries. In the modern digital ecosystem, SEO is very important, as it directly influences aspects of discoverability and user engagement. [17]

History of SEO

The origins of SEO date back to the late 1990s, which also saw the emergence of search engines such as Yahoo! and AltaVista. Early SEO techniques were unsophisticated, based on keyword stuffing and backlinks. As search engines matured, especially with Google's PageRank algorithm, SEO practices became more complex, shifting their focus to quality content, relevance, and user experience. With the evolution of web technologies from simple static HTML pages to dynamic SPAs and SSR described above, SEO also had to evolve. In modern SEO practices, technical considerations such as page load times, mobile responsiveness, and schema markup have been incorporated to keep up with the complexity of modern web architectures.

Search Engines and Data Crawling

Search engines use automated programs called web crawlers or bots-e.g., Googlebot) to index the web. These crawlers systematically visit websites, analyzing their content, metadata, and structure to determine relevance and ranking. Crawlers mainly use two ways to access data from a website:

HTML Parsing Obtaining information from static or server-side rendered pages.

API Requests Obtaining structured information through publicly available APIs or schemas like JSON-LD.

These processes make sure that search engines are able to provide appropriate and high-quality results to users. For websites, proper indexing is the beginning of visibility; hence, they usually apply SEO methods like sitemaps and meta tags to guide crawlers.

Allowing and Restricting Crawlers: The Role of Robots.txt

Websites regulate crawler permission with a special configuration file called robots.txt. The robots.txt file, located in a website's root directory, instructs crawlers regarding what the website allows or disallows. The most common directives include the following.

Allow Explicitly allows crawlers to access specified directories or files.

Disallow Prevents crawlers from accessing certain areas.

Crawl delay Specifies the time interval between successive requests to reduce server load.

Example of a robots.txt file:

User-agent: * Disallow: /private/ Allow: /public/

Most search engines and other legitimate crawlers honor what is specified in robots.txt. However, it cannot be technically enforced. Most search engines would follow it, but other bots may not, including some scrapers.

Legal and Ethical Considerations

The use of robots.txt introduces legal and ethical dimensions to web crawling. Ignoring robots.txt directives may violate the terms of service (ToS) of a website and, in some jurisdictions, could be considered unauthorized access. In contrast, adhering to robots.txt ensures ethical scraping practices and minimizes the risks of legal repercussions.

Implications for Web Scraping

For web scrapers, SEO-oriented websites can simplify data extraction. Crawlers often make key information accessible and well structured for indexing purposes, reducing the complexity of scraping tasks. However, it is important that scrapers make these allowances responsibly, with attention to directives like robots.txt to not conflict with website administrators or legal frameworks. SEO highlights the mutual reliance between websites and search engines, highlighting that visibility and accessibility are balanced by ethical concerns. Recognizing this dynamic is crucial to develop web scraping strategies that are both technically sound and ethically responsible.

Chapter 2

Techniques and Innovations in Web Scraping

Web scraping is an evolving discipline that combines elements of software engineering, data science, and ethical considerations. As web technologies have become increasingly complex, the methods and tools used for scraping data have also advanced. This chapter provides an in-depth exploration of the foundational techniques, optimized processes, and real-time strategies that support modern web scraping practices.

Having laid the groundwork in Chapter 1 by discussing the evolution of web technologies and the challenges they introduce, this chapter focuses on how data can be efficiently extracted from modern websites. The goal is to bridge the gap between theoretical understanding and practical application, enabling readers to understand the methodologies that make modern scraping effective.

2.1 Foundational Scraping Techniques

Web scraping techniques are the foundation of data extraction from the web, and their evolution parallels the increasing complexity of modern websites. This section examines the basic methods used to navigate static and dynamic web content effectively.

2.1.1 Scraping Static Websites

Static websites are created using only simple HTML and CSS, without complex backend logic or JavaScript frameworks that render content. In static sites, the response to an HTTP GET request returns the whole HTML file of the requested page. This will also, of course, contain the visible content of the webpage, metadata, and structural information. Since the content is pre-rendered and delivered directly to the client, obtaining and parsing the HTML is straightforward. The simplicity of static websites is a direct consequence of their architecture, which separates content from presentation. The contents of a static webpage are precompiled and fixed and stored in fixed files on a web server. This makes them much faster to load and also easier to scrape, since the required data will already be there in the HTML source.

HTML parsing is done with libraries such as Cheerio.js[18] for JavaScript or Beautiful Soup[19] for Python. These libraries enable a developer to navigate and change the Document Object Model (DOM) in order to extract the target data. In general, the DOM is a tree-like representation of an HTML document structure where every node is assigned to an element or text. Among these, Cheerio.js, for instance, is a light yet fast solution inspired by jQuery[20] allows for CSS selectors to effectively target elements within the DOM. CSS selectors, as popularized by jQuery, allow effective targeting of elements based on their attributes, tag name, class, or ID. Examples targeting elements by tag - div, by class - .header or by ID - #main. CSS selectors are central to web scraping. They offer a declarative way to specify what to extract. They are especially effective because they enable intricate queries that merge several criteria, such as selecting all the p elements within a div of a specified class. This adaptability allows the extraction of structured data from web pages that are even moderately intricate.

To illustrate, consider the following HTML content:

```
<!DOCTYPE html>
<html lang="en">
<head>
   <meta charset="UTF-8">
   <meta name="viewport" content="width=device-width, initial-scale=1.0">
   <title>Example Page</title>
</head>
<body>
   <div id="content">
      <h1 class="title">Welcome to the Example Page</h1>
      This is a simple HTML page for scraping.
      Item 1
         Item 2
         Item 3
      </div>
</body>
</html>
```

This HTML contains a simple structure where content is organized into a header, a paragraph, and a list. Each element is marked with attributes such as classes and IDs, which can be used to precisely locate and extract them using CSS selectors.

Using Cheerio.js, the content of this page can be scraped as follows:

```
const cheerio = require('cheerio');
const axios = require('axios');
async function scrapeStaticPage(url) {
   try {
      const { data: html } = await axios.get(url);
      const $ = cheerio.load(html);
```

```
const title = $('h1.title').text();
const description = $('p.description').text();
const items = [];
$('ul.items li.item').each((index, element) => {
    items.push($(element).text());
});
console.log('Title:', title);
console.log('Description:', description);
console.log('Items:', items);
} catch (error) {
    console.error('Error scraping the page:', error.message);
}
```

```
scrapeStaticPage('https://example.com');
```

In this code, the axios[21] library retrieves HTML content via an HTTP GET request, while Cheerio.js parses HTML and applies CSS selectors to extract specific elements. The function \$('h1.title').text() retrieves the text content of the header with the class title, while the \$('ul.items li.item') query iterates through all list items in the unordered list with the class items, collecting their text values.

This approach highlights the theoretical reinforcements of static website scraping: reliance on a predictable and stable document structure, the use of efficient tools for DOM traversal, and the application of precise CSS selectors for targeted data extraction. Together, these elements form the basis for extracting meaningful information from static websites in an organized and efficient manner.

2.1.2 Dynamic JavaScript-Driven Websites

The shift from static to dynamic was a significant evolution in web architecture, driven by the demand for interactivity and richer user experiences. In contrast to static website simplicity, modern dynamic websites generate their content client-side after the first HTML file has loaded. This results in a distinct architecture that requires an alternative method for web scraping, due to the requirement for the scraper environment to execute Javascript.

Dynamic websites rely on CSR: client-side rendering, which means that JavaScript execution in the browser fetches data and builds or updates the DOM. A basic GET request typically yields an HTML file with a tree structure, serving as the primary entry point of each framework without revealing the actual content. The execution of Javascript code is necessary to produce the content. Such dynamic content generation is thus very helpful for interactivity and performance, yet significantly increases the challenge for traditional scraping methods relying on a simple HTTP request to retrieve the HTML file. For example, content may only be loaded after certain user interactions or asynchronous API calls, and thus the scraper would not be able to extract any meaningful data directly from the server response. Once the HTML file has been generated, the process is the same as static web scraping using Cheerio.js or a similar library.

An approach to address these challenges is through the use of headless browsers in the context of web scraping. A headless browser is one that operates without a graphical user interface and, in so doing, can run JavaScript, render dynamic content, and programmatically interact with websites. Although originally designed for use within automated testing of web applications, headless browsers have proved invaluable when it comes to scraping JavaScript-driven websites.

The main purpose of a headless browser is to provide a running environment for the JavaScript code within a page, as would any real user. For this reason, it can load dynamic content, perform updates to the DOM, and generally render the page just as a user would see it. This will be quite important in scraping, as it ensures that scrapers will get the fully loaded and dynamic state of the webpage. This has been made easy with tools such as Puppeteer[22], Playwright[23], and Selenium[24] by giving access to a programmatic method for controlling headless browsers such as Chromium and Firefox.

The described method involves a number of crucial steps. The scraper starts the headless browser and opens the target URL. The browser will execute JavaScript, fetching data and rendering the DOM while the page is loading. The scraper can then interact with the rendered page, either by using DOM traversal techniques or capturing network requests to get structured data from APIs. For example, Puppeteer has the option of waiting until certain DOM elements load so that all dynamic content is present before scraping even starts.

Headless browsing revolutionized the possibility of scraping dynamic websites, but not without challenges. The headless browser consumes much more resources than a traditional parser since it emulates a full browser environment. Additionally, most dynamic websites include some protection against scraping, such as CAPTCHAs, bot detection algorithms, and rate limiting, all of which will be analyzed later in the thesis.

In other words, as websites transitioned from static to dynamic formats, scrapers needed to adapt to the imposed challenges. This gap was closed with the advent of headless browsers, which allowed scrapers to handle content generated by JavaScript. By simulating a user environment, headless browsers enabled developers to interpret data extracted from a contemporary web application, broadening the scope and utility of web scraping for dynamic web technologies.

2.1.3 Network Request Interception

Another advanced technique for handling dynamic websites is based on the interception of network requests, enabling asynchronous communications between the client and the server. Many modern dynamic websites heavily rely on these communications to fetch data and update the UI. Driving factors include AJAX, Fetch APIs, and WebSocket protocols. Although static websites include all the content in the initial HTML response, dynamic websites load the essential data through discrete network requests at runtime. Thus, traditional scraping methods, which rely on DOM parsing, might not be sufficient to fetch the required data. This naturally leads to other ways in which this can be achieved, including intercepting network requests. A technique that is helpful for scrapers to directly access data exchanged between the browser and server.

The interception of network requests relies on the very basic principle of observing the communication channels that take place and helping to exchange data in dynamic web applications. By monitoring these requests, scrapers can determine the endpoints and payloads involved in fetching a required piece of information. It bypasses the need for DOM traversal by capturing raw data in formats like JSON or XML, which is often well-structured and more parse-friendly compared to HTML.

The most important application of network request intercepting in scraping is to emulate the requests a browser would make. That means peering into the headers, parameters, and body of HTTP requests and programmatically re-creating them in order to extract data from a source directly. Browser developer consoles, Mitmproxy[25], and Puppeteer's network interception capabilities are some of the tools that make this process easier. These tools provide developers with the functionality to monitor network activity, filter out relevant requests in order to automate data extraction workflows.

For example, SPAs built using either React or Angular fire API calls on user interaction or page load to dynamically fetch data from a server. These can be observed in the Network tab of the browser's developer tools. From these requests, the scraper could know which endpoint URL, HTTP method, headers, and query parameters are required to fetch the data. Once identified, the scraper can programmatically issue identical requests for that data without having to render a full page or interact with the DOM.

This approach offers some real advantages, especially in terms of efficiency and accuracy. Directly accessing structured data, scrapers avoid processing overhead and potential inconsistencies due to dynamic manipulation of the DOM. Additionally, intercepting network requests reduces resource consumption because it does not need to open and execute JavaScript code. However, the method also means some challenges: Many websites implement protection measures such as authentication tokens, rate limiting, and encryption that should be handled with care to make your data extraction successful.

More specifically, a higher application in network request interception involves the processing of communications done via WebSockets. WebSockets are a special technology that allows for real-time, bidirectional communication between client and server, usually adopted to provide live updates, as in chat systems or financial dashboards. Capturing and decoding these streams demands deeper knowledge of the protocol and tools able to analyze binary data formats, such as Wireshark[?] or specialized WebSocket libraries.

By focusing on the communication layer, the mess in the DOM manipulation and extract structured data can be bypassed. There are lots of pitfalls while developing because of protective mechanisms, but generally, it is a crucial tool in an arsenal of different web scraping strategies to handle such applications driven by JavaScript.

2.1.4 Comparing Techniques

Web scraping techniques differ greatly depending on the structure of the target website, the technology stack it employs, and the level of protection against automated extraction. The three fundamental approaches discussed earlier HTML Request Parsing, Headless Browsers, and Network Interception each have distinct strengths and weaknesses that influence their suitability for different scenarios.

HTML Request Parsing

HTML request parsing is best suited for simple, static websites such as blog pages, news archives, and document repositories where the structure remains stable over time.

Advantages

- Fast and lightweight: Does not require rendering JavaScript, making it efficient for simple pages.
- Low resource consumption: Can run on low-power machines since it does not require a full browser.
- Structured HTML accessibility: Works well for static websites where content is fully available in the initial HTTP response.

Disadvantages

- Fails on JavaScript-heavy sites: The websites that use client-side rendering, where key content is loaded asynchronously via JavaScript, making direct HTML requests insufficient.
- Vulnerable to structural changes: If a website modifies its HTML layout, scrapers relying on hard-coded selectors may break.
- Limited interactivity: Cannot handle authentication flows, infinite scrolling, or other dynamic behaviors.

Headless Browsers

Headless browsers are ideal for web applications with complex JavaScript rendering, such as e-commerce platforms, social media feeds, and interactive dashboards.

Advantages

- Handles JavaScript-heavy sites: Can extract content from Single Page Applications (SPAs) that load dynamically.
- Supports interactivity: Can click buttons, scroll pages, and fill forms, making it effective for sites requiring user interaction.
- Bypasses basic anti-scraping measures: Can simulate real user behavior, reducing detection risks.

Disadvantages

- Resource-intensive: Running a full browser instance consumes significantly more CPU and memory compared to HTML parsing.
- Slower execution: Because the page needs to be rendered and JavaScript executed, scraping speeds are lower.
- Detection risks: Some sites use methods to detect automation tools like Puppeteer, which may require additional evasion techniques.

Network Interception

Network request interception is most effective for applications that rely heavily on APIs, such as stock market data feeds, travel booking platforms, and live sports score updates.

Advantages

- More efficient than full browser automation: Extracts structured data directly without rendering HTML.
- Bypasses UI-based obfuscation: Since it does not rely on DOM elements, it is unaffected by dynamic layout changes.
- Reduces computational overhead: Consumes less resources than a headless browser, since it only captures API responses.

Disadvantages

- Requires knowledge of network protocols: Extracting API calls can be complex and require developer tools, proxies, or packet analyzers.
- APIs may require authentication: Some requests include API keys, tokens, or session cookies, making them harder to replicate.
- WebSocket encryption: Some real-time communication protocols encrypt responses, requiring decryption techniques to access meaningful data.

Technique	Best For	Challenges
HTML Request Parsing	Static websites	Ineffective for dynamic content
Headless Browsers	JavaScript-heavy websites	Resource-intensive
Network Interception	API-driven or dynamic data	Requires protocol expertise

 Table 2.1. Comparison of Foundational Scraping Techniques

Chapter 3

Scraping Prevention and Countermeasures

3.1 Detection and Mitigation Strategies

Web scraping detection and mitigation strategies are at the core of modern web application security to avoid data extraction by unauthorized sources. These methodologies detect and deter attacks by malicious parties through a blend of technological frameworks, behavioral analytics, and artificial intelligence. In this section, we consider key methodologies, their theoretical grounds, and their practical application in order to point out their efficiency and limits.

3.1.1 Browser Fingerprinting

Browser fingerprinting involves capturing and analyzing unique attributes of the browser environment, including the user agent string, the type and version of the browser, as well as the operating system used, any active plugins, the time zone and language of the machine, the screen resolution and various other active settings. These attributes when combined, they form a browser fingerprint that can identify and track individual users. For example, as described by Acar et al. in their FPDetective development, browser fingerprinting can involve analyzing JavaScript and Flash-based attributes to identify fingerprinting behavior across thousands of websites [26]. Laperdrix et al. provide a comprehensive survey on browser fingerprinting, explaining how a combination of device-specific attributes, such as canvas fingerprinting and WebGL parameters, contributes to unique identification [27]. Surprisingly, browsers give a lot of information about the user and, based on research carried out by the Electronic Frontier Foundation [28], 84% of collected fingerprints are globally exclusive, and they found that the next 9% were sets of two. Even though fingerprints are dynamic, new ones can be matched up with old ones with 99.1% correctness, therefore allowing websites to track online behavior in order to serve hyper-personalized advertisements. In other cases, these fingerprints are being used for anti-bot detection as scraping bots, which rely on headless browsers or minimalistic HTTP clients, usually have limited diversity in their fingerprints, making them easier to detect.

However, the latest developments in scraping tools that emulate full browser environments have increased the complexity of this task.

3.1.2 Traffic Pattern Analysis

Among the widely used methods is the analysis of traffic patterns. It involves checking the inbound traffic coming to a website for anomalies that may suggest the presence of some bot-like traffic. For instance, humans tend to produce requests at a somewhat predictable cadence, interspersed by moments of inactivity that reflect human browsing habits. A web scraper, on the other hand, will make a lot of requests in quick succession to scrape off the maximum data. This kind of anomaly is detected by applying various statistical techniques or machine learning models. Machine learning, in particular, excels at distinguishing between normal and anomalous traffic by identifying subtle non-linear patterns in request metadata, such as headers, referrers, and IP addresses.

These characteristics are derived primarily from the frequency of the request, the duration of the response, and the diversity of user agents. The aim is to classify the traffic as human or automated by employing supervised learning algorithms. More recent deep models combine a transformer and convolutional neural network architecture in their designs, aimed at both temporal and spatial feature capture from web traffic streams with greater detection sensitivities toward even very small-scale anomalies indicative of the operation of a scraping bot. [29].

Another powerful strategy is to combine different detection methodologies, such as rule-based systems, statistical anomaly detection, and machine learning classifiers. This layered approach will analyze the traffic indicators of IP diversity, request rates, and browser capabilities to enhance robustness against sophisticated bots. [30].

Unsupervised learning techniques also play a critical role in environments where labeled data is scarce. These methods identify suspicious patterns and behaviors indicative of scraping activities, enabling dynamic anomaly detection across varied traffic datasets [31]. Furthermore, indirect traffic analysis techniques, such as examining packet size distributions and timing, offer a unique perspective in detecting scraping activities, particularly those using obfuscation techniques such as VPNs and proxy servers [32].

Although traffic pattern analysis provides powerful tools to detect web scraping, it is not without challenges. False positives can arise from legitimate high-frequency access, such as API usage or bulk data downloads. Moreover, advanced bots employ sophisticated tactics to evade detection, such as dynamically altering request patterns and mimicking human browsing behaviors. Addressing these challenges requires continued refinement of the detection algorithms and the integration of multiple detection methodologies.

3.1.3 CAPTCHAs

CAPTCHAs (Completely Automated Public Turing tests to tell Computers and Humans Apart) are among the most widely adopted techniques to combat web scraping. CAPTCHAs, introduced by von Ahn et al. [33], serve as a challenge-response mechanism designed to differentiate between human users and automated bots by requiring users to solve problems that are easy for humans but difficult for bots. Common forms of CAPTCHA include identifying distorted text, selecting specific images, or solving simple puzzles. The effectiveness of CAPTCHAs lies in their ability to disrupt the automated workflows of scraping bots. Since bots rely on programmatically sending requests and parsing responses, introducing CAPTCHAs adds an additional layer of complexity that often requires sophisticated AI models to bypass. For example, image-based CAPTCHAs require visual recognition capabilities that many bots lack. Moreover, modern CAPTCHAs, such as reCAPTCHA, employ behavioral analysis by monitoring user interactions like mouse movements and keystrokes to assess whether the user is a bot.

Despite their effectiveness, CAPTCHAs have limitations. Sophisticated bots that use advanced AI and machine learning can now solve many CAPTCHA challenges with high accuracy [34]. Additionally, CAPTCHAs can negatively impact the user experience, particularly when they are overly difficult or appear frequently. This trade-off between security and usability necessitates careful deployment of CAPTCHAs to ensure that they do not deter legitimate users.

In recent years, invisible CAPTCHAs have gained popularity. These CAPTCHAs work silently in the background by analyzing user behavior to identify bots without requiring explicit user interaction. Although this approach minimizes user inconvenience, it raises concerns about privacy and the potential misuse of behavioral data.

In summary, CAPTCHAs remain a valuable tool in web scraping prevention strategies, particularly when combined with other detection methods. Their evolution toward more user-friendly and sophisticated forms continues to make them a crucial component of modern web security.

3.1.4 Honeypotting

Honeypotting is a really smart web security approach aimed at the detection, deflection, or analysis of unauthorized scraping activities. This involves the intentional placement of deceptions, such as links, forms, or hidden content, within a web application in order to detect and evade malicious bot traffic. The principle of honeypotting involves enticing the bots into taking action on the deceptive elements so that indications of malicious behavior are attained and unauthorized data extraction can be barred effectively. Generally, honeypotting means creating traps that would appear very legitimate for bots but should remain undetectable to a human user. This can be performed by honeypots through hidden links and form fields by using the HTML and CSS properties of display: none or visibility: hidden. Since these are not visible or accessible by humans, any action could be counted as that from a bot. If a bot does come into contact with a honeypot, then it can log activity, analyze the data, and block further access.

Honeypots have proven to be very efficient in picking up the presence of scraping bots that get through traditional rate-limiting or CAPTCHAs. The honeypot data allows organizations to further tune their bot detection algorithms and apply countermeasures in focused ways: logging IP addresses, user agents, or request patterns of bots interacting with honeypots to proactively block malicious entities.

One good thing about honeypotting is that it actually does not affect the user experience at all. Unlike CAPTCHAs, which may introduce bad experiences for users, honeypots just silently do their thing in the background, unobtrusive as such. However, advanced bots may sometimes identify honeypots using advanced algorithms, reducing their effectiveness. Moreover, honeypots not configured properly can block legitimate users or leak sensitive data.

Recent research in the field of honeypot technologies has been geared toward integrating machine learning into their detection capabilities. Machine learning models can interpret honeypot interaction data in search of patterns indicative of bot behavior, enhancing the accuracy of bot detection systems. Dynamic honeypots, changing attributes at regular intervals, provide another layer of complexity for a bot seeking to avoid detection.

Honeypotting represents one of the most important tools in web scraping prevention strategies. Complementing traditional methods of detection, such as traffic pattern analysis and CAPTCHAs, honeypots provide a sound and low-impact approach to the identification and mitigation of malicious bot activities. Continuous innovation in the design of honeypots and their integration with machine learning promise to further strengthen security in web applications.

3.1.5 IP Reputation Systems

Another layer of defense is the IP reputation system. These systems maintain whitelists and blacklists of IP addresses and networks with known malicious participation in sending out bots, in order to block or throttle requests from these flagged IP addresses. The advanced ones update their databases dynamically using crowd-sourced intelligence and real-time threat detection. However, this very method is highly challenged by scrapers using rotating proxies or large pools of residential IP addresses.

3.1.6 Behavioral Analytics

Behavioral analytics enhances detection by focusing on how users interact with a website. Legitimate users have complex, non-linear behaviors-like scrolling, hovering over elements, and irregular click patterns, all of which characterize the interaction as organic. Scraping bots programmatically interact with the website and often bypass elements altogether. Advanced behavior analytics leverages machine learning to model typical user interactions and flags deviations indicative of bot activity. These models are trained on massive datasets of real and malicious interactions for the purpose of maximum accuracy. However, a bot can be programmed to simulate human behavior, implementing random delays and scrolls, making it difficult to discover by this method.

3.1.7 Conclusion

Although these techniques are strong in devising ways through which web scraping may be detected and hampered, all of them have certain weaknesses. In cases of legitimate high-frequency access, such as API usage, false positives may result from traffic pattern analysis. Browser fingerprinting is vulnerable to spoofing, in which bots spoof actual browser attributes. CAPTCHAs can be solved with advanced AI bots or using solver APIs. Honeypotting can be avoided by checking that the element is visible before scraping it. In turn, IP reputation systems struggle against the now ordinary use of residential proxies and methods for obscuring one's IP address. Behavioral analytics, promising as it may be, depends on relentless retraining as bot tactics continue to change, with a likely added downside. It might punish users with unconventional browsing behavior. Therefore, successful bot detection implementation requires a multilayered solution that combines different approaches against a multidimensional threat landscape while also applying preventive techniques in the development stage.

Methodology	Advantages	Disadvantages	Bypass Techniques
Browser Finger-	Unique identifi-	Susceptible to	Bots emulating full
printing	cation of users;	spoofing; com-	browser environments;
	effective against	plexity increases	fingerprint spoofing
	simple bots	with advanced	
		bots	
Traffic Pattern	Effective in de-	False positives	Bots mimicking human
Analysis	tecting abnormal	for legitimate	traffic patterns; dis-
	traffic patterns;	high-frequency	tributed requests
	adaptable to vari-	users; complex	
	ous data points	implementation	
CAPTCHAs	Proven effective-	Impacts user	Using CAPTCHA
	ness; disrupts	experience;	solver APIs; employing
	automated work-	advanced	human solvers
	flows	AI can solve	
		CAPTCHAs	
Honeypotting	Minimal user	Can be detected	Bots avoiding hidden
	impact; provides	by sophisticated	elements; inspecting
	valuable bot data	bots; improper	DOM for honeypots
		configuration	
		risks	
IP Reputation	Blocks known	Less effective	Use of proxy pools; dy-
Systems	malicious IPs;	against rotating	namic IP switching
	crowd-sourced	proxies or resi-	
	intelligence im-	dential IPs	
	proves accuracy		
Behavioral Ana-	Detects complex	Requires retrain-	Bots simulating realistic
lytics	bot behavior;	ing; challenges	user interactions
	leverages machine	with sophis-	
	learning	ticated bots	
		mimicking human	
		behavior	

Table 3.1. Comparison of Web Scraping Bot Detection Methods

3.2 **Preventive Development Techniques**

3.2.1 Rate Limiting and Throttling

Rate Limiting and throttling are some of the core web security methods that allow the reduction in the number of requests possible by a client to a server within a certain timeframe. The above measure helps server avoid abuse that includes denial-of-service attacks or data scraping. Rate limiting works by defining the limits on the requests; if these limits are exceeded, further requests will be blocked or refused for some time. To illustrate this, in order to allow a maximum of 100 requests per minute from one IP address, any additional requests can be throttled or rejected. The rate limit is usually implemented on the server side or on the API gateway. A server tracks the timestamp of the moment the request arrived and its origin. Then, the server compares successive requests from that very same origin against the defined threshold within that time window. In such cases, when the limit is exceeded, the servers can delay or refuse processing. Most of the time, this will be combined with an HTTP status code corresponding to an issue, in this case 429 (too many requests). This approach prevents not only overloading on the server, but also shares resources equitably between the users.

Rate limiting works really well in APIs that support programmatic access to sensitive data or high value data. By setting a limit on how often clients can request specific endpoints, APIs can prevent scripted processes from attempting to pull a large volume of data in a very short amount of time. Dynamic rate limiting adjusts those thresholds based on usage patterns or customer authentication levels and provides another layer of security and flexibility. In this regard, for example, authenticated users could have higher thresholds compared to anonymous users; in that respect, their status as trusted would be reflected in that way.

Although rate limiting has proved effective, there are a number of drawbacks. Advanced scraper tools usually implement distributed networks of bots that will send requests from different IP addresses. This is called a "botnet" where the distributed nature of their requests evades rate limiting, given that it seems like requests from different origins for whatever server. Besides that, even legitimate high-frequency users-for instance, data analysts or businesses using APIs for legitimate reasons-get mistakenly rate-limited, and nuanced policies have to be made to create a distinction between benign and malicious traffic.

Rate limiting and throttling remain two highly critical tools when trying to protect against unauthorized access and resource abuses. Although it is very powerful with the aspects of simplicity and adaptability, the weaknesses outlined clearly prove that further security measures should be contemplated in view of the bypass techniques available which have been polished accordingly.

3.2.2 Dynamic Code and CSS Attributes

Another effective approach to impeding web scraping would be frequent changes to the website code or dynamic CSS attributes. This approach includes periodic changes to the structure or presentation of web content; thus, the scraping bots have to be continuously updated, and hence increase the cost and complexity of maintaining them.

Dynamic changes can involve regular renaming of HTML element IDs, classes, or attributes. For example, unique identifiers and class names can be dynamically generated using randomized or hashed strings with each server response. Since most scraping bots rely on static identifiers to locate and extract data, these changes disrupt their ability to locate the desired elements consistently.

The data presentation might be changed similarly using dynamic CSS attributes. For instance, some elements could be moved in the view by using the position or z-index CSS property, for instance, yet nothing actually changed within the real DOM. This allows critical data, for instance, to shift visually without user experience degradation while confusing web scraping bots relying on static DOM parsing.

More advanced would be the combination of server and client-side approaches. On the server side, scripts randomize key aspects of the HTML structure as it is rendered; on the client side, JavaScript manipulates the DOM even more after the page loads. This twolayer approach ensures that even sophisticated bots rendering JavaScript cannot interpret this dynamically changing environment.

It becomes even more powerful when applied together with some content obfuscation techniques, such as encoding or encrypting certain data elements. For instance, critical data fields should appear only after decoding a token passed via JavaScript; this means that a bot should be able to execute and interpret complex scripts.

In turn, frequent changes in code and dynamic attributes need to be weighed against development and maintenance costs. They may bring along some potential problems with debugging and SEO optimization, since the dynamic changes might affect search engine crawlers.

This means that regular code changes and dynamic CSS attributes create one more layer of complexity for a web scraping bot, which cannot be relied upon in its results. An organization that sets up a dynamically changing web reduces the level of effectiveness with which automated tools can scrape information, hence making this tactic a core one in any serious anti-scraping strategy.

3.2.3 API Key and Token-Based Authentication

API key and token-based authentication are considered basic mechanisms to secure programmatic access to web resources. Both methods will use some unique identifier or token that authenticates and authorizes client requests to make sure that only valid users or applications can interact with the API or server.

An API key is a static alphanumeric string assigned to a client after they register or subscribe to an API. The key is a sort of credential that must accompany every request to the API. Returning to our example of a weather service API, let us say that clients must specify the API key as a query parameter or via a request header. Then, upon receiving a request, the server checks the key against its database of issued keys and only processes the request if the key is valid. API keys are relatively easy to implement and use; thus, they are popular for public APIs and low-risk applications.

Token-based authentication, on the other hand, uses dynamically generated tokens, which are, for the most part, short-lived and bound to a session or a specific user. This adds security in that the exposure of keys is minimal. Tokens are often issued after some authentication process, such as login via OAuth. In this case, once authenticated, a token is assigned to the client, and they must specify it within subsequent requests. The server then verifies if the token is valid and has not expired before allowing access.

Both have unique advantages in reducing unauthorized access. API keys are easy to administer and implement; thus, they work well for applications with very minimal security needs. Tokens provide much more flexibility and are more secure, especially when implemented with other more advanced authentication mechanisms like OAuth 2.0. Tokens can also encode additional metadata, such as user roles or permissions, enabling fine-grained access control.

However, these methods are not immune to exploitation. API keys, if compromised, can be used by malicious actors until the key is revoked. Tokens are also susceptible to interception during transmission if secure channels such as HTTPS are not used. Rate limiting, IP whitelisting, and secure storage of keys and tokens are among other measures taken by developers. In addition, refresh tokens and the token expiration policy reduce the time frame during which the intruder might misuse the obtained token.

The bottom line is that API key and token-based authentication is a critical component in today's web application security. The two systems discussed have advantages and disadvantages, but their appropriate execution and incorporation into other security complementing solutions could provide an adequate scenario for safeguarding APIs and web resources against unauthorized access.

3.2.4 JavaScript Challenges

JavaScript challenges are a sophisticated method to block unauthorized access and automated scraping, where many computational obstacles that bots must overcome introduce an access barrier to the Web. They are based on client-side execution, making the way forward very difficult for general scrapers who do not have advanced browser emulation. The website filters off the simple bots and crawlers by embedding code which should run correctly in every successful request.

A common implementation of JavaScript challenges would be generating a dynamic token or changing some structure of a webpage's content according to a JavaScript script executed in the client browser. Services like Cloudflare use challenges based on JavaScript as part of their anti-bot protection mechanisms. Once a request is made, the server sends a script to the client that performs some calculation or data transformation. The result of these operations, normally encoded as a token, is subsequently sent to future requests as evidence of valid client-side execution.

The biggest advantage with JavaScript challenges is that they are able to filter and block bots that cannot run JavaScript. Traditional scraping scripts that directly parse HTML or utilize lightweight HTTP clients will be instantly defeated. Moreover, JavaScript challenges can be dynamic, with the complexity or logic of the script changing to out-smart evolving scraping technologies.

However, with the sophistication of scraping tools, this becomes a big challenge. Advanced bots integrating headless browsers like Puppeteer or Selenium are actually able to emulate the execution of JavaScript, thus bypassing such defenses. If poorly implemented, such challenges to JavaScript may accidentally degrade user experience, particularly for users operating slow devices or with limited browser support. Therefore, it is very critical to balance security and usability when implementing JavaScript challenges effectively.

Overall, JavaScript challenges are a very good addition to web security frameworks, but only in combination with other types of prevention. The fact that they can take advantage of the limitations of basic scraping tools makes them an effective deterrent, though they require careful implementation to remain effective against advanced threats.

3.2.5 Cookie-Based Authentication

Cookie-based authentication is a widely used technique to maintain session integrity and verify user identity in web applications. By storing a unique session identifier in the user's browser as a cookie, the server can validate requests and ensure that they originate from authenticated users. This mechanism is integral to protecting web applications against unauthorized access and provides a seamless user experience by persisting authentication across multiple requests.

The process begins when a user logs into a web application. Upon successful authentication, the server generates a session token, which is stored as a cookie in the client's browser. All subsequent requests to this server contain this cookie, which helps the server validate the session of the user without asking him to log in every time. To improve security, cookies have various attributes like 'HttpOnly', which will not allow client-side scripts to access the cookie, and 'Secure' will send cookies only over HTTPS.

Cookie-based authentication has a number of advantages. It simplifies session management by offloading state persistence to the client, reducing server-side storage requirements. Furthermore, it provides a seamless user experience, by which users remain authenticated across page reloads and navigations. Advanced configurations, such as the 'SameSite' attribute, mitigate risks like cross-site request forgery (CSRF) by restricting the contexts in which cookies are sent.

Meanwhile, despite its broad usage, cookie-based authentication also has vulnerabilities. The leakage of cookies while transmitting over unsecured channels will make it possible for an attacker to hijack the sessions of users. Likewise, wrong handling of storing cookies-for instance, saving sensitive information without encryption-makes users a potential target for a data breach. Bots that have tools for handling cookies will use a weakly implemented authentication mechanism to mimic user behavior.

These risks can be reduced by various protection means that a developer could implement, such as session tokens, periodic rotation of the tokens, and various server-side activity monitors. In addition, combining the results of cookie-based authentication with other means of protection, such as CAPTCHAs or IP whitelisting, enhances it. In other words, cookie-based authentication is still a basis in the security of modern web applications. Its simplicity and effectiveness make it popular, while its vulnerabilities raise the need for thoughtful implementation and integration with complementary security strategies to protect against evolving threats.

3.2.6 CDN Security

Content Delivery Networks, or CDNs, are the backbone of modern web security, distributing web content across a number of geographically dispersed servers. Beyond their core function of reducing latency and improving load times, CDNs boast security features in an attempt to mitigate threats such as web scraping. By acting as intermediaries between the client and the origin server, CDNs provide a layer of abstraction that enables sophisticated request filtering and traffic analysis.

One of the main security features of CDNs is their ability to apply traffic filtering and rate limiting at a scale. CDNs analyze incoming requests in real time for patterns indicative of bot behavior, such as high frequencies of requests or unusual headers. These malicious requests are filtered out before reaching the origin server, reducing the risk of data scraping and server overload. Other CDN providers, such as Cloudflare and Akamai, use machine learning algorithms that dynamically adapt to evolving bot tactics, thus ensuring protection against sophisticated threats.

In addition to filtering traffic, many CDNs also have built-in solutions for bot management. These include challenges such as CAPTCHAs and JavaScript verification that identify legitimate users and automated bots. By centralizing these defenses, CDNs relieve origin servers of the computational overhead associated with detecting and mitigating bots.

However, the effectiveness of CDN-based security depends on proper configuration and integration. Poorly configured rules or overly aggressive filtering can result in false positives, which block real users. Advanced bots that can mimic human behavior can also bypass many of the standard detection mechanisms used by CDNs, requiring constant refinement of security protocols.

In the end, security is one of the most important modern CDN web application defenses. Using distributed infrastructure and advanced analytics, CDNs provide scalable adaptive solutions to combat web scraping and other threats. Their integration with other preventive measures further enhances their efficacy, making them an essential component of a comprehensive security strategy.

3.2.7 Data Obfuscation Techniques

Data obfuscation techniques are designed to obscure the presentation and structure of web content, making it challenging for scraping bots to extract meaningful information. By altering how data are displayed or encoded, these techniques create significant hurdles for automated tools while remaining transparent to legitimate users.

One of the common ways of data obfuscation is to encode critical information. For example, numerical values or textual data can be encoded using base64 encoding or other reversible encoding schemes. The encoded data are then correctly rendered on the client side using JavaScript. Scraping bots that rely on static DOM parsing are often unable to decode this information without executing the accompanying scripts, adding a layer of complexity to their operations.

Dynamic content rendering is another powerful obfuscation technique. With the use of JavaScript, a web application might load or change some information on the page asynchronously, without static content in place. Such data would be shown only to those users whose browser environment was working. Very effective against bypassing JavaScript execution by bots: they cannot get information from dynamically rendered data.

Obfuscation also extends to the structural organization of web content. For example, critical elements can be randomized or scattered across the DOM and their positions corrected visually through CSS. That confuses scraping bots that rely on predictable element structures to extract data.

Although data obfuscation methods significantly interfere with the work of scrapers, it is not without its own weaknesses. Advanced bots with headless browsers and complex parsing algorithms can bypass a lot of these obfuscation techniques. Too much obfuscation can also cause website performance problems and make them difficult to maintain, so balanced implementation is really important.

3.2.8 Adaptive User Interface Rendering

Adaptive user interface (UI) rendering is an innovative approach to web security that dynamically adjusts the presentation of content according to user behavior and context. By tailoring the loading and display of web elements to individual users, this technique complicates scraping attempts, as bots are unable to anticipate or replicate the exact rendering patterns.

An implementation of adaptive UI rendering is lazy loading, where web content is loaded incrementally based on user interaction. For example, images or data tables may load only as the user scrolls through the page. This approach minimizes the amount of content immediately available for scraping, forcing bots to mimic user interactions to access the full dataset.

Another example is content fragmentation, whereby data is fragmented into smaller pieces and served only when certain predefined triggers are activated, for example, mouse movements or specific keystrokes. This ensures that only users exhibiting human-like behavior will get the full content, whereas bots will be left with incomplete or fragmented datasets.

Adaptive UI rendering also uses behavioral analytics to alter content delivery. Analyzing click patterns, cursor movements, and scrolling behaviors, among others, allows a website to detect potential scraping activities in real time and respond accordingly. For example, if an interaction seems abnormal, it could prompt the server to mask or hold sensitive content.

Although very efficient in making scraping workflows fragile, the implementation of adaptive UI rendering needs to be done with care; otherwise, it becomes counterproductive by affecting user experience. Aggressive policies or overreliance on behavioral triggers could easily frustrate legitimate users; a balanced approach would be required.

Technique	Advantages	Disadvantages	Bypass Tech-
			niques
Rate Limiting and Throttling	Effective in con- trolling traffic volume; reduces server overload	Susceptible to distributed at- tacks using botnets; potential false positives for legitimate high-frequency users	Use of distributed bots and rotating IPs to distribute requests
API Key and Token-Based Authentication	Simple to imple- ment; enhances security for pro- grammatic access	Vulnerable to key/token inter- ception; requires secure storage and periodic rotation	Token inter- ception during transmission; reuse of stolen API keys
JavaScript Chal- lenges	Filters out basic bots; adapts dynamically to evolving threats	Advanced bots with headless browsers can bypass; potential user experience degradation	Emulation of JavaScript ex- ecution using headless browsers like Puppeteer
Cookie-Based	Maintains session	Vulnerable to	Session hijacking
Authentication	integrity; seam- less user experi- ence	cookietheft;impropercon-figurationscanexpose sessions	through cookie theft or replay attacks
CDN Security	Scalable and dis- tributed; reduces server load with advanced traffic filtering	Misconfiguration risks; false posi- tives may block legitimate users	Bots mimicking human behavior to evade detec- tion
Data Obfuscation Techniques	Conceals critical information; dis- rupts static scrap- ing tools	Performance degradation; so- phisticated bots can decode obfus- cation	Parsingobfus-cateddatausingadvancedal-gorithmsandheadlessbrowsers
Adaptive User In- terface Rendering	Dynamically ad- justs to user be- havior; fragments data for enhanced security	May affect user experience; com- plex implementa- tion and mainte- nance	Mimicking hu- man interactions to trigger full content rendering

 Table 3.2. Comparison of Preventive Development Techniques

3.3 Advanced Protection Mechanisms

3.3.1 Anti-Scraping SaaS Platforms

The role of Software-as-a-Service platforms, which have become an essential arm in the fight against web scraping, is to offer a suite of anti-scraping solutions with the integration of advanced technologies such as machine learning, behavioral analytics, and cloud computing. They promise ease of integration, scalability, and speed, which positions them well with organizations in dire need of solid web application security.

Anti-scraping SaaS solutions generally work by sitting between the client and the web server, analyzing incoming traffic in real time for malicious requests to block. Cloudflare Bot Management, Akamai Bot Manager, and Datadome are some of those that use advanced algorithms that check request patterns, browser attributes, and user behavior in order to tell a legitimate user from a scraping bot. For example, Cloudflare Bot Management uses machine learning models trained on vast amounts of data to evolve with changing scraping tactics to provide continued protection.

The advantages of SaaS solutions go beyond detection. These platforms can be easily integrated with an organization's infrastructure and provide real-time analytics with insight into traffic patterns and vulnerabilities in the system. In addition, they also come with customized settings that will allow an organization to orient its defenses to use cases. Companies can define very granular rules based on geolocation, user agent, or API usage patterns to allow or block requests.

Of course, such platforms also come with their own challenges. The fact that they are cloud-operated introduces latency that can affect the user experience. Advanced scraping tools that can emulate human behavior may occasionally bypass detection, and thus require constant updates of the algorithms.

3.3.2 Browser Integrity and Verification

Browser integrity and verification techniques involve a set of advanced methods to ensure that client interactions with web applications are well-intentioned and untampered with. This kind of technique validates the expected behavior of a browser, hence distinguishing real users from bots trying to pose as browsers.

One of the most important ways to validate the integrity of the browser is to challenge it with certain cryptographic operations. For example, browsers should calculate some result with the help of their cryptographic key, such as a digital signature. The server then checks the result to confirm who the actual client is. Such mechanisms guarantee that only real browsers-and not light scraping software-will be able to work with the web application.

Another widely used technique is JavaScript-based validation, which uses scripts to evaluate browser attributes and behavior. Examples of JavaScript challenges include testing the execution of certain scripts, the consistency of outputs rendered, and timing metrics. These techniques have been employed by a set of tools, including Cloudflare's browser challenges, to identify deviations that give away the bot. Forcing bots to reproduce subtleties exhibited by real browsers, these challenges pose formidable obstacles to automated tools.

Other emerging solutions include device attestation protocols such as WebAuthn, which are improving browser verification. These protocols leverage hardware-based modules for the attestation of a device's authenticity, hence increasing security against complex bots.

Yet, despite all their power, browser integrity and verification mechanisms have limitations. Advanced bots, equipped with headless browsers and strong emulation capabilities, can take advantage of some of the mechanisms. Poorly implemented challenges also risk disrupting legitimate users, making it essential to balance security with usability.

3.3.3 Hybrid Defense Architectures

Hybrid defense architectures are the epitome of holistic approaches to web security, combining multifaceted anti-scraping techniques into one strong, cohesive framework. Individual methods presented throughout this chapter can easily suffer from various drawbacks and vulnerabilities that advanced bots successfully manipulate. These approaches leverage their complementary strengths, integrated together in an effective dynamic system against a broad landscape of scraping threats.

Generally speaking, a hybrid architecture starts at the very bottom foundational layer with techniques such as rate limiting, API authentication, and cookie-based session management. These measures are generally used to provide a baseline level of security by limiting both the rate and scope of access to web resources. In addition, dynamic capabilities include behavioral analytics and browser integrity checks that detect and respond to anomalous activity in real time. Further complication arises with the addition of JavaScript challenges and adaptive UI rendering, which forces bots to further emulate complex browser interactions and behaviors.

Advanced machine learning-powered bot detection and anti-scraping SaaS platforms develop an additional layer of intelligence with great scalability. These technologies analyze a huge volume of data and self-improve with the constantly evolving tactics of bots, keeping the architecture effective in the face of sophisticated threats. Lastly, honeypotting and dynamic traps form proactive countermeasures based on deception-based strategies, where bots get lured into self-disclosure and therefore enable the operation of targeted mitigation.

What makes this approach effective is the synergy of these methods within a hybrid framework. For example, browser fingerprinting may identify a bot, while rate limiting and traffic analysis contain its activity. Similarly, adaptive UI rendering can disrupt scraping workflows, and machine learning models dynamically adjust defenses based on observed behavior. By layering these techniques, organizations can address both known and emerging threats, maximizing their ability to detect and prevent web scraping.

However, one should not forget the limitations in any anti-scraping strategy. With the continuous development of web scraping tools and techniques, no defense can ensure 100% prevention. Advanced bots with sophisticated emulation capabilities and distributed infrastructures will always find a way to break through even the most comprehensive defenses. The goal of a hybrid defense architecture is not to achieve absolute security, but to make web scraping prohibitively difficult and resource intensive for the average bot operator.

With a hybrid architecture in place, a website can be more consistent in the blocking of non-sophisticated bots, minimizing the overall risk to web applications. At the same time, it can really focus on identifying and mitigating the more advanced threats, ensuring a balance between security and usability. In the dynamics of Web scraping, this multilayered and adaptive approach remains one of the most viable means of protecting online resources by accepting the fact that total protection will never be possible.

Chapter 4

Integration of AI in Web Scraping

4.1 Machine Learning in Scraping

4.1.1 Adaptive Algorithms for Intelligent Data Extraction

The fast-changing nature of web technologies requires the creation of adaptive algorithms that can support efficient data extraction processes against frequent and often unpredictable changes in the structure of web pages. Traditional web scrapers are often based on static configurations fitted to specific page layouts that can quickly become obsolete due to changes in the underlying HTML or JavaScript frameworks. The key development to meet this challenge in web scraping is the design of automatically adaptable web wrappers. These wrappers, through the use of algorithmic techniques, can automatically adapt to structural changes to continue operation without human intervention.

Adaptability in web scraping algorithms is based on two key principles: adaptability to recognize structural patterns and the intelligence to act upon deviations. By embedding machine learning models that identify and classify recurring patterns in HTML and DOM elements, adaptive algorithms can make predictions about future changes to the layout of websites. These models leverage historical data to build probabilistic mappings of likely changes that allow scrapers to adapt their extraction logic preemptively.

One pioneering contribution in this domain is that of Ferrara and Baumgartner on automatically adaptable web wrappers [35]. The proposed framework represents the design for a modular architecture in which scrapers are endowed with a hierarchical understanding of web page structures. In this study, a multilayer representation of the content is adopted, in which the various layers are mined to detect structural consistencies between versions. If there are discrepancies, the system updates the extraction pathways using heuristics and probabilistic methods to maintain accuracy.

Central to this approach is the application of context-aware parsing mechanisms. These mechanisms assess semantic relationships between DOM elements and allow the scraper to focus on content of interest while ignoring irrelevant changes. For example, when a webpage redesign introduces new ornamental elements, the adaptive algorithm can distinguish these from fields containing substantive data, and preserve the integrity of the extracted information. It minimizes a lot of manual reconfiguration and develops and scales the operation of web scraping. Theoretically, its ability to generalize across a range of different Web environments motivates adaptive systems. Their respective algorithms leverage reusable components through which the abstraction of the extraction logics is possible and achieve modularity to a very higher degree to plug into varied architectures seamlessly. The integration of real-time feedback loops allows the scraper to continuously learn from the success or failure of extraction attempts, continuously improving its models. This not only allows the scrapers to be increasingly accurate, but also resilient against increased sophistication in anti-scraping measures.

Adaptive algorithms for intelligent data extraction represent a transformative shift in the practice of web scraping. By embedding adaptability into the core architecture of scrapers, these systems ensure sustained performance in dynamic web ecosystems. The integration of context-aware parsing, probabilistic mapping, and real-time learning mechanisms underscores the theoretical and practical advancements that define modern adaptive scraping technologies.

4.1.2 Automated Proxy Management

In modern web scraping, effective proxy management is crucial, especially in overcoming the restrictions imposed by websites to limit automated access. Proxies are intermediaries that anonymize requests, enabling scrapers to distribute their traffic across multiple IP addresses and avoiding detection and throttling. However, traditional methods of proxy management, relying on static proxy pools and unsophisticated rotation mechanisms, often prove quite insufficient when dealing with the sophisticated anti-scraping mechanisms of modern web platforms. The integration of machine learning and intelligent automation into proxy management systems represents a significant evolution in this domain.

Automated proxy management systems apply data-driven methods for optimal selection, rotation, and usage of proxies. These systems monitor key performance indicators, such as response times, success rates, and error codes, to determine the effectiveness of individual proxies. Then, machine learning models are trained on this data to predict the likelihood of success of future requests routed through particular proxies. Based on these predictions, automated systems dynamically adjust the usage of such proxies in an effort to minimize latency, maximize throughput, and reduce the risk of IP bans.

Anomaly detection is the most important element in intelligent proxy management. Advanced algorithms analyze traffic patterns for irregularities that could indicate possible detection by target websites. For example, sudden spikes in failed requests, or uniform browsing across multiple proxies, may trigger automated actions such as switching to alternative pools of proxies or making the browsing more human-like. This proactive approach not only enhances the efficiency of web scraping operations but also ensures their long-term viability.

Third-party APIs and integrations improve real-time feedback loops to improve the adaptability of proxy management systems. Active feedback mechanisms permit scrapers to iteratively refine their methods based on the feedback from previous interactions of these agents. If one proxy gets flagged off or blocked consistently, the traffic can be assigned to better proxies while retraining the models. Such a process of iteration keeps the model responsive to shifting sands on the Web.

The most accessible kind of bot proxies are data center IPs. Because of their detectable features as non-residential or server-generated interactions, they are blocked by most anti-bot mechanisms. Data center proxies originate either from cloud services or hosting providers; because they lack the diversity and authenticity of residential IPs, they can be detected with the more advanced anti-bot systems. These systems use advanced algorithms to differentiate genuine user activity from automated behavior, considering network attributes such as ASN, geolocation, and browsing patterns. This has consequently made reliance on datacenter proxies alone significantly detrimental in most scraping operations.

Residential proxies have become a necessary solution due to such limitations. Residential proxies reroute traffic via IP addresses assigned to real residential users by ISPs. Unlike datacenter IPs, residential IPs are less suspicious since they mimic actual usage patterns and are more difficult to detect and block. This means that web scrapers can interact with targeted websites under a cover of legitimate user traffic, thus avoiding most of the anti-bot defenses.

Large companies, such as Bright Data or Apify, have pioneered the development and commercialization of residential proxy services. The Bright Data platform connects one to the widest pool of residential IPs anywhere around the globe-unsurpassed in either coverage or flexibility. Often, this might mean adding premium features such as geotargeting abilities, session persistence, and automatic IP rotation, greatly increasing the effectiveness of residential proxies. This technology keeps scrapers ahead in their ability to dynamically respond to changing conditions in Web environments and navigate sophisticated mechanisms of blockage.

The success of residential proxies depends on how well they can be camouflaged within regular web traffic. Anti-bot systems often check behavioral metrics, such as request frequency and navigation patterns, for anomalies. Residential proxies, together with behavioral simulation techniques, enable scrapers to mimic human-like interactions that reduce the possibility of detection. Moreover, diversity and distribution in residential IPs will mean that requests will appear to be coming from a wide variety of different legitimate users, which complicates the task of anti-bot systems in attempts to isolate and block scraping activities.

4.1.3 CAPTCHA Solving with AI

CAPTCHAs are an important set of Web security mechanisms that attempt to distinguish between legitimate human users and automated bots. The challenges presented by these mechanisms range from recognizing distorted characters to performing visual or logical tasks. Despite their utility, significant strides have been made using AI to overcome CAPTCHA systems, putting into question their long-term efficacy. This section will cover the main types of CAPTCHA, ways to solve them with the use of AI, and the availability of third-party APIs providing CAPTCHA-solving capabilities. Text-based CAPTCHAs are among the earliest and most common forms of challenge. They usually ask users to decipher and input alphanumeric characters presented in distorted or obscured formats. AI techniques have been successfully able to solve such challenges, especially those using optical character recognition in combination with supervised machine learning [36]. Deep learning models, trained on a large number of CAPTCHA images, could recognize and interpret text, even distorted, rotated, or with overlapping characters, with high precision. Research shows that such models could achieve very high success rates by undermining the reliability of traditional text-based CAPTCHAs.

In this sense, puzzle-based CAPTCHAs challenge users with tasks that require assembling the pieces to make up a meaningful picture or finishing some logical series; they rely on skills believed to be distinctive to human cognitive capabilities. However, such a solution may be quite easily tackled using a combination of convolutional neural networks and reinforcement learning algorithms. The AI is able to pick up patterns and ways by which it works out the puzzle and can apply them to different formats. Researchers have pointed out the potential of AI-driven approaches to break down such CAPTCHAs, compelling developers to develop newer, more sophisticated designs to keep their strength alive.

Image selection CAPTCHAs require users to scroll through grids of images and to indicate which of them meet a certain criterion, like "select all images containing bicycles." These have also been subjected to attacks from AI, as they depend on human visual recognition and semantic understanding. Deep learning models, which are trained for image classification tasks, do an outstanding job in the analysis and categorization of images, hence finding the right selections [37]. Recent research underscores the capabilities of AI technologies in attacking image-based CAPTCHAs with very high success rates and discloses the shortcomings of their design.

In addition to creating customized AI models, there are also third-party services to solve CAPTCHAs. Such platforms, including 2Captcha and Anti-Captcha, will offer APIs easily integrated into web scraping workflows that automatically solve the CAPTCHAs. These services ensure high accuracy with the help of a mixture of AI-driven algorithms and, where necessary, human operators. In addition, many of these services work with a range of CAPTCHA types, from reCAPTCHA v2 and v3 to Hcaptcha, through their easy-to-implement solutions that make the bypassing process practical for developers who want to skip model development challenges.

These continuously improving AI technologies are a big challenge for the CAPTCHA systems since they keep degrading the separating barriers between human and machine interactions. The review of various types of CAPTCHAs and an overview of the AI-based solutions and third-party services depict the dynamic nature of this field. While AI capabilities continue to grow, innovative and robust CAPTCHA designs will be required to maintain their status as valid security mechanisms.

4.2 NLP for Data Understanding

The integration of Natural Language Processing (NLP) into web scraping methodologies has transformed the ability to extract and process unstructured textual data from the web. The rapid development of digital content on the Internet has created a wide resource base for information, with the major component existing in unstructured forms such as social media posts, customer reviews, news articles, and many others. NLP provides computational tools to process and analyze such information data and convert it into structured and actionable insight. This section discusses some of the key techniques, applications, and challenges in using NLP for the interpretation of unstructured data over the web [38].

Unstructured text data are particularly difficult to analyze due to lack of structure or predefined format. NLP addresses these challenges by employing preprocessing techniques such as tokenization, stemming, and stop-word removal. Tokenization breaks down text into individual words or phrases, forming the basic units for further processing. Stemming reduces words to their root forms, ensuring consistency across variations, while stop word removal eliminates common but uninformative words like "the" or "and." These preprocessing steps prepare the raw text for more advanced analysis, ensuring that only the most relevant information is retained.

Named Entity Recognition (NER) is a crucial NLP technique to understand web data. This process involves the identification and categorization of entities that might be contained in unstructured textual data, such as names, organizations, dates, and locations. For example, in customer reviews scraped from e-commerce platforms, NER can extract references to specific products or brands. Employed with sentiment analysis, another essential application of NLP, one can programmatically assess the polarity of text- whether it is positive, negative, or neutral. Such analysis can enable organizations to gauge public opinion on various products or topics. This is particularly valuable in domains such as social media analysis or overall product satisfaction, where insights into customer attitudes and preferences can inform marketing strategies.

Another significant application of NLP in web scraping is text summarization. Large volumes of textual data, such as news articles or forum discussions, can be distilled into concise summaries using techniques such as extractive and abstractive summarization. Extractive summarization identifies key sentences or phrases from the original text, while abstractive summarization generates new sentences that convey the core ideas.

The combination of web scraping and NLP is not limited to academic or industrial research but extends to practical tools and APIs. Frameworks like SpaCy[39] and NLTK[40] offer pre-built modules for tasks such as tokenization, POS tagging, and entity recognition, simplifying the integration of NLP into web scraping pipelines. Additionally, third-party APIs, such as MonkeyLearn and Aylien, provide ready-to-use sentiment analysis and text classification services, eliminating the need for custom model development. These tools democratize access to advanced NLP capabilities, allowing developers to focus on specific application objectives without digging into the complexities of model training.

As the capabilities of NLP continue to expand, they pave the way for more sophisticated approaches to the understanding of web data. While this chapter has focused on traditional and non-generative NLP methods, the next chapter will explore the transformative potential of generative AI in web scraping, offering new paradigms for data extraction and analysis.
Chapter 5

Leveraging Generative AI for Scraping

The rise of Large Language Models (LLMs) marked an evolutionary era in artificial intelligence. OpenAI's ChatGPT received sweeping attention in 2023 receiving the title of the faster ramp in users in a consumer internet app. These models demonstrated unprecedented capabilities in understanding and generating human-like text, reshaping the landscape of NLP, content generation, and automation. This triggered an everlasting race between AI companies and academics to create the most advanced and affordable model, while industries quickly adopted this technical achievement, easily deploying chatbots, documentation helpers, and other tasks requiring human cognition.

Web scraping has a symbiotic relationship with LLMs [41]. Automated data extraction from online sources using web scraping techniques was crucial to their implementation, as they rely heavily on large amounts of high-quality structured data for training. By extracting diverse datasets, ranging from product reviews to scholarly articles, scraping plays a fundamental role in shaping the contextual richness and performance of LLMs. The accuracy and volume of these data directly influence the model's ability to generate reliable outputs. In addition to training, scraping also enables real-time features such as advanced web search capabilities. LLMs leverage live data to enhance user search experiences by providing detailed source-backed answers based on natural language prompts (e.g. Perplexity, ChatGPT search[42]). Real-time scraping facilitates dynamic browsing interactions used to provide up-to-date information, giving context to the LLM model, thus improving the relevance of responses.

Consequently, generative AI models are great tools for web scraping. Many complex scraping tasks and practices mentioned earlier in this thesis, from NLP to self-modifying scrapers[43] and CAPTCHA solvers, can be transformed using LLMs unlocking new possibilities. This chapter explores how LLMs can be used to automate web scraping, bypass anti-scraping mechanisms, and transform the data attribution layer.

5.1 LLMs in Scraping Pipelines

Introducing LLMs to web scraping has changed the way data are extracted, processed, and understood. They can be integrated in every stage of the pipeline such as crawling, extracting, bypassing anti-bot measures, post processing, making the process smarter cutting down on manual work while boosting accuracy.

5.1.1 Crawling Stage

At the initial stage of web scraping, the crawling phase establishes the foundation for effective data acquisition. Dynamic site mapping, powered by large language models (LLMs), offers a revolutionary approach to navigating complex websites. This involves the generation of detailed visual maps that illustrate how the structure of a site is organized. In addition to mapping visual elements, LLMs also create semantic maps that identify relationships between different sections of a website. These semantic maps go beyond surface-level details to reveal deeper interconnections, such as the hierarchy of pages and the logical flow between them. By predicting page hierarchies, LLMs can infer how various components of a website are related, enabling an optimized crawling strategy. For example, pages that are central to the site's purpose can be prioritized for crawling, while less significant pages are given lower precedence. In addition, LLMs excel at identifying interlink relationships, which are crucial for understanding how users and bots are expected to navigate the site. This predictive capability ensures that the crawling process is not only efficient, but also thorough, capturing critical data while minimizing redundant efforts. These capabilities are particularly useful for websites with nested navigation structures or dynamically loaded content, which pose significant challenges for traditional crawlers.

5.1.2 Data Extraction Stage

The data extraction phase represents the core of web scraping activities, and LLMs have introduced transformative methods to enhance this process. One prominent application involves the use of LLMs to understand free text and structure the results into formats such as JSON. For example, LLMs can process customer reviews written in free-form text and output them in a structured schema, allowing for easier analysis and integration into databases. Tools such as LangChain[44] exemplify this capability, as they are specifically designed to extract semantic meaning from unstructured data. These tools can analyze context, identify key elements, and generate precise structured outputs.

Another powerful approach uses LLMs to identify the necessary CSS or XPath selectors for precise data extraction. This process involves analyzing the DOM (Document Object Model) structure of a webpage to pinpoint the exact elements that contain the desired data. Once identified, the LLM can generate custom scraper code tailored to the specific layout of the page. This method is particularly valuable for websites that frequently update their layouts, as LLMs can dynamically adapt the generated scripts to accommodate changes. For example, if a site's structure is modified, an LLM-powered scraper can adjust its approach by recalibrating the selectors or rewriting sections of the code. This adaptability ensures that data extraction remains functional over time without requiring extensive manual intervention.

Despite these advancements, it is important to recognize the limitations of using LLMs for direct scraping. Employing LLMs to process every request is computationally expensive and introduces significant time overhead. Although the precision and versatility of LLMs make them invaluable during the setup phase, their continuous use for routine scraping tasks is impractical. Instead, a hybrid strategy proves to be the most effective. In this approach, LLMs are employed during the initial configuration phase to create an adaptable scraping framework. Once the setup is completed, traditional methods can take over the repetitive execution of scraping tasks. This combination ensures the benefits of LLMdriven precision while maintaining the cost-effectiveness of conventional techniques.

5.1.3 Antibot Measures Bypassing Stage

One of the most significant challenges in web scraping involves bypassing anti-bot measures, which are designed to identify and block automated access to websites. LLMs offer pioneering solutions to address these barriers. For example, text-based CAPTCHAs, which require users to interpret distorted text, can be solved by LLMs through advanced natural language processing and pattern recognition capabilities [45]. The models analyze the text within the CAPTCHA and generate an accurate response, effectively bypassing the challenge.

For more complex visual CAPTCHAs, LLMs can be integrated with vision-based models to interpret and solve the challenge. These models are capable of analyzing images, identifying patterns, and producing appropriate solutions. Some antibot systems rely on behavioral analytics to detect bots by analyzing user actions, such as mouse movements, scrolling patterns, or click timings. Here, LLMs can simulate human-like behavior, producing realistic interactions that are indistinguishable from those of a genuine user. For example, an LLM can mimic the randomness of human mouse movements or introduce deliberate delays between actions to avoid detection.

By combining these techniques, LLMs provide a versatile toolkit to navigate various anti-bot mechanisms. However, the ethical implications of bypassing such measures should not be overlooked. Although these capabilities expand the potential for data acquisition, they also raise questions about legality and responsible use. Deploying LLMs for antibot bypassing requires careful consideration of both technical feasibility and ethical boundaries.

5.1.4 Post-Processing Data Stage

After data have been successfully extracted, the post-processing phase ensures its usability and relevance. LLMs excel in categorizing and tagging unstructured data, which is critical to gaining actionable insights[46]. For example, social media comments often come in diverse and unorganized forms, making it difficult to identify patterns or trends. LLMs can analyze these comments, categorize them by sentiment, and tag them with relevant metadata. This process not only organizes the data, but also enhances its value by making it easier to interpret.

Another example is the processing of product reviews. LLMs can identify recurring themes, highlight positive and negative feedback, and even summarize the overall sentiment of a dataset. This capability is rooted in the natural language understanding (NLU) features of LLM, which enable them to comprehend context, intent, and sentiment within text. The categorization process can also be extended to include tagging for specific use cases, such as identifying comments related to product quality, shipping issues, or customer service.

The theoretical basis for these capabilities lies in the ability of LLMs to perform contextual analysis at scale. By understanding the relationships between words, phrases, and a broader context, LLMs transform raw data into structured insights. This not only saves time, but also allows organizations to extract deeper meaning from their data. Ultimately, the role of LLMs in the post-processing stage underscores their value as a comprehensive solution for refining and enriching the outputs of web scraping pipelines.

5.2 Advanced Generative AI Capabilities in Web Scraping

5.2.1 Cross-Modal Scraping: Bridging Vision and Text

Cross-modal scraping represents a paradigm shift in the way data are extracted from the Web, combining the strengths of natural language processing (NLP) and vision-based approaches. Unlike traditional methods that rely solely on HTML parsing, cross-modal techniques incorporate visual elements such as screeenshots and multimedia, alongside textual information, to enable more nuanced and comprehensive data extraction. For example, Optical Character Recognition (OCR) can be applied to images containing embedded text, while video frame analysis can extract relevant information from multimedia content. These techniques are particularly beneficial for visually complex web pages, such as those containing dynamic charts, infographics, or non-standard text placements. The integration of vision and NLP expands the range of data that can be scraped, providing richer datasets for downstream applications.

Cross-modal scraping is especially useful in scenarios where structural data is obfuscated or presented in a visually appealing but less machine-readable format. Using visionenabled models, such as those powered by large multimodal models (LMMs), scrapers can navigate and interpret visual cues just as a human user would. This capability not only enhances the depth of data extraction, but also opens new possibilities for analyzing unstructured or semi-structured content that would otherwise be inaccessible.

5.2.2 WebVoyager: A Case Study in Multimodal Web Agents

WebVoyager represents a groundbreaking advancement in the domain of web scraping by utilizing LMMs to create a fully autonomous end-to-end web agent. Designed to interact seamlessly with real-world websites, WebVoyager combines visual and textual inputs to navigate, interpret, and extract data. Unlike earlier approaches that relied on static HTML snapshots or simplified simulations, WebVoyager operates directly on rendered web pages, leveraging the full spectrum of visual and semantic information available in modern web environments [47].

The core functionality of WebVoyager lies in its ability to observe and act iteratively. For example, the agent receives inputs such as screenshots and textual descriptions of interactive web elements. It then formulates a thought process to determine the appropriate action, such as clicking, scrolling, or typing, before executing that action on the live website. This iterative approach mimics human browsing behavior and allows WebVoyager to adapt to real-time changes in web layouts or interactive components.

The utility of WebVoyager extends beyond mere navigation. It excels in handling complex, multi-step tasks that require contextual understanding and decision-making. For example, in a benchmark evaluation, WebVoyager successfully completed tasks such as identifying specific products on e-commerce platforms, locating information from academic websites, and even extracting structured data from multimedia-rich pages. Using both visual analysis and semantic reasoning, WebVoyager achieved a task success rate significantly higher than that of text-only or static methods, demonstrating its effectiveness in real-world scenarios.

One of the distinguishing features of WebVoyager is its innovative evaluation protocol, which uses GPT-4V as an automatic evaluator. This protocol uses a combination of human-like judgment and machine-driven assessment to validate the agent's performance. With a reported 85.3% agreement rate with human evaluations, this method ensures a reliable measure of the agent's capabilities.

Despite its advancements, WebVoyager is not without limitations. The reliance on screenshots as primary input means that text-heavy websites or those with highly complex visual designs can pose challenges. Moreover, while the agent is capable of solving basic web navigation tasks autonomously, further refinements are necessary to handle edge cases such as CAPTCHA solving or interacting with highly customized web components.

In conclusion, WebVoyager exemplifies the potential of multimodal agents in web scraping, showcasing how vision and NLP can be integrated to push the boundaries of what automated systems can achieve. As tools like WebVoyager continue to evolve, they promise to redefine the landscape of data extraction, offering new levels of precision, adaptability, and intelligence to navigate the complexities of the modern web.

5.3 Use Cases and Implications

5.3.1 Training Data Acquisition

One of the most impactful applications of generative AI in web scraping lies in its ability to facilitate the acquisition of training data for machine learning and LLMs[4]. Personalized scraping pipelines for specific domains enable the extraction of relevant highquality data sets essential for the refinement of AI models. For example, generative AI can identify and parse academic articles, extracting metadata, abstracts, and references to construct domain-specific data leads. Similarly, it can collect data sets for conversation agents by extracting interactions from public forums or posts on social networks. The precision offered by LLMs ensures that the collected data is relevant and free of unnecessary noise, significantly improving the quality of downstream applications.

However, this capability also raises ethical concerns. Scraping proprietary or sensitive data without explicit consent can violate privacy laws and intellectual property rights. Many websites for example disallow the ChatGPT User-Agent to protect the data they

expose from using it to train OpenAIs models. Ensuring compliance with regulations such as GDPR and CCPA is critical when designing scraping pipelines. By employing ethical guidelines and transparency in data handling, generative AI can strike a balance between utility and responsibility in training data acquisition.

5.3.2 Domain-Specific Applications

Generative AI-powered scraping extends its utility across diverse sectors, unlocking unique use cases tailored to each domain. In e-commerce, it can track price fluctuations, monitor competitor strategies, and aggregate customer reviews to provide actionable market insights. For healthcare, web scraping can extract and standardize medical research data, enabling healthcare providers to stay up-to-date with the latest treatments and trends. Similarly, in education, generative AI can scrape educational content such as open-access research, course materials, and multimedia resources, creating rich learning repositories.

The financial sector also benefits significantly from these advancements. Scraping financial news, stock market trends, and analyst reports empowers investors and firms with real-time intelligence. LLMs ensure that the extracted data maintain high accuracy, allowing for the automation of time-sensitive analyses. These domain-specific applications highlight how generative AI can transform industries by providing precise and actionable information at scale.

5.3.3 Augmented Search Engines

Search engines are foundational tools for accessing information online, and generative AI offers promising enhancements to their capabilities. By integrating scraped data, search engines can provide richer, more contextual responses to user queries. For example, a search engine powered by generative AI can aggregate data from live news sources, ecommerce platforms, and social media to deliver comprehensive responses tailored to realtime events.

Incorporating LLM-powered real-time insights further refines search engine functionality. Adaptive ranking systems, informed by user trends and dynamic content, ensure that results remain relevant and personalized. This capability not only improves user satisfaction, but also improves the discoverability of content creators, fostering a mutually beneficial ecosystem. The integration of scraped data into augmented search engines demonstrates the potential of generative AI to redefine how information is accessed and utilized.

5.3.4 Implications for Future Development

The growing reliance on generative AI in web scraping brings with it a set of implications that warrant careful consideration. Scalability remains a critical challenge, as the computational requirements for large-scale scraping operations increase alongside the complexity of modern websites. Optimizing these processes to balance performance and cost will be vital for long-term sustainability. Ethical and legal frameworks will also need to evolve to keep pace with these advancements. As generative AI enables deeper and more nuanced data extraction, ensuring compliance with global privacy standards becomes increasingly complex. Collaborative efforts among researchers, policymakers, and industry stakeholders will be essential to navigate these challenges and establish best practices.

Finally, the integration of generative AI into web scraping welcomes a new era of possibilities for data access and use. By addressing its limitations and maximizing its strengths, this technology has the potential to drive innovation in countless fields, shaping a future where information is more accessible, actionable, and impactful than ever before.

Part II

Implementation

Chapter 6

soniq: No-code web scraping platform for structured data extraction

6.1 Problem definition

Web scraping plays a critical role in modern data acquisition, contributing significantly to fields such as search engine indexing, AI model training, and competitive intelligence. However, it presents a technological barrier due to each website's different and potentially complicated structures, and standard web practices that require headless browsers and anti-bot measurements that can be price and resource intensive task. Naive approaches to web scraping consist of different database structures for each data domain to be extracted and different code implementations for each website to be scraped regardless of the framework that will be used to create the scraping pipeline. Therefore, while public information is available to anyone through user interaction, only companies and teams that have the resources to create complex pipelines can automatically acquire and leverage this information. Although there are numerous pay-as-you-go API scraping services, these are predominantly closed-source, limiting users' ability to directly access raw data, define their own extraction schemas, or customize the scraping pipeline. The goal of this thesis is to create an open-source, no-code, supported by LLM technology Software as a Service platform that can give the same advantage to individuals or academic teams, in order to be able to integrate public data to their knowledge systems, train or fine-tune their own LLMs and AI models through a more accessible multi-domain, structured web data extraction without the need to create a scaping pipeline from scratch.

6.2 Technologies Used

The thesis implementation can be broken down to three counterparts: backend, and frontend, which have the ability to inter-communicate and ensure the functionality of the final product.

6.2.1 Backend

The backend counterpart of the application has been built in order to be able to handle multiple users, successfully respond to concurrent requests and efficiently store and query



Figure 6.1. soniq Component UML Diagram

large amount of data. In order to make this possible, open-source, production-grade widely adopted technologies have been utilized:

MongoDB

MongoDB[48] is a NoSQL, document-oriented database designed for high-performance, scalability, and flexibility. Unlike relational databases (SQL), MongoDB stores data in BSON (Binary JSON) format, allowing for schema-less document structures, making it ideal for storing and managing unstructured or semi-structured data. It is widely used in modern web applications, including web scraping pipelines, big data processing, and AI-driven analytics.

Key Features

- 1. Document-Oriented Storage
 - Data is stored in collections of documents, rather than traditional rows and tables.
 - Each document is structured as a JSON-like object (BSON), supporting nested fields and flexible data types.
- 2. Schema-Less Design (Dynamic Schema)
 - Unlike SQL databases, MongoDB does not require a fixed schema.
 - Documents in the same collection can have different structures, making it easy to evolve data models over time.
- 3. High Scalability (Sharding & Replication)
 - Sharding: Distributes data across multiple servers to handle large-scale work-loads.
 - Replication: Ensures high availability and fault tolerance by maintaining multiple copies of data.
- 4. Rich Query Language

- Supports CRUD operations (Create, Read, Update, Delete) with powerful filtering, aggregation, and indexing.
- Queries can be performed using JavaScript-like syntax.
- 5. Indexing for Fast Queries
 - Supports multiple types of indexes, including single-field, compound, and text indexes for efficient query performance.
 - Geospatial indexing enables location-based searches.
- 6. Aggregation Framework
 - Similar to SQL's GROUP BY, the aggregation pipeline allows for complex data transformations, filtering, and computations.
 - Used for data analysis and processing within the database itself.
- 7. Built-in Horizontal Scaling
 - MongoDB is optimized for distributed computing, making it an excellent choice for big data and real-time applications.
- 8. ACID Transactions (Multi-Document)
 - Supports multi-document ACID transactions, ensuring data integrity in complex applications.

Features like schema-less design, rich query language, aggregation framework and horizontal scaling make MongoDB the most responsible decision for this implementation enhancing data versatility and ability to create informative data tables.

FastAPI

FastAPI[49] is a modern, high-performance web framework for building APIs using Python 3.7+. It is designed for speed, scalability, and ease of use, making it an excellent choice for developing web services, including RESTful APIs for data-driven applications like web scraping systems.

FastAPI is built on Starlette (for async web handling) and Pydantic (for data validation and serialization), offering automatic data validation, OpenAPI documentation, and asynchronous support.

Key Features

- 1. High Performance (Asynchronous & Non-Blocking
 - Built on ASGI (Asynchronous Server Gateway Interface), allowing concurrent request handling.
 - Supports async/await, making it faster than Flask for I/O-bound operations.

- 2. Automatic OpenAPI & Swagger Documentation
 - Enables self-documenting APIs without additional effort.
 - FastAPI automatically generates OpenAPI documentation (/docs endpoint using Swagger UI).
 - Provides Redoc UI for alternative API exploration (/redoc).
- 3. Data Validation with Pydantic
 - Uses Pydantic for strict data validation and serialization.
 - Ensures API requests contain correctly formatted data.
 - Allows definition of strict types for request/response models.
- 4. Dependency Injection System
 - Allows clean separation of concerns with dependency injection.
 - Useful for integrating authentication, database connections, and middleware efficiently.
- 5. Built-in Security (OAuth2, JWT)
 - Supports OAuth2 and JWT authentication out of the box.
 - Provides automatic handling of API security, reducing development effort.
- 6. WebSocket & GraphQL Support
 - Enables real-time communication through WebSockets.
 - Supports GraphQL APIs alongside traditional REST endpoints.
- 7. Easy Integration with Databases
 - Works well with MongoDB (Motor), PostgreSQL (SQLAlchemy), Redis, and other databases.
 - Supports async database operations for high scalability.

Therefore, FastAPI is a great addition to the implementation's stack checking every potential need there is in order to design a robust and well documented API infrastructure.

crawl4ai

Crawl4AI[50] is a next-generation AI-powered web scraping framework that leverages large language models (LLMs) for intelligent data extraction. Unlike traditional web scrapers that rely on static CSS selectors and XPath queries, Crawl4AI uses machine learning and NLP techniques to dynamically detect, structure, and extract web content, making it more resilient to UI changes and less prone to anti-scraping mechanisms.

It is designed to work with structured and unstructured web data, integrating autocorrecting pipelines, proxy management, and API-based scraping for a scalable and adaptive scraping experience.

- 1. AI-Powered Smart Extraction
 - Uses LLMs and NLP models to understand webpage layouts and extract structured data automatically.
 - Can infer relationships between elements (e.g., extracting product details even if HTML structures vary across different sites).
 - Supports semantic extraction, allowing it to capture contextual information.
- 2. Hybrid Extraction Methods
 - AI-Assisted Extraction: Automatically detects and extracts key data fields without predefined rules.
 - Manual Extraction: Allows users to define CSS selectors and XPath queries for precision.
 - LLM-Assisted Schema Generation: Generates CSS selectors for users based on textual descriptions.
- 3. Multi-Threaded & Asynchronous Execution
 - Supports parallel scraping for high efficiency.
 - Uses asynchronous I/O operations to handle multiple pages simultaneously, reducing execution time.
- 4. Open Source
 - Code is available to inspect, build, contribute.
 - Ability to extend and add extra features applying to each use case.
- 5. Proxy and CAPTCHA Handling
 - Integrates with residential proxies, rotating IPs, and CAPTCHA solvers.
 - Uses adaptive proxy selection to minimize detection.
 - Supports third-party CAPTCHA solving services like 2Captcha and Anti-Captcha.

Crawl4AI functionality is very crucial for this implementation as it provides an abstraction layer for tasks that would otherwise need to be implemented from scratch such as undetectable headless browsers, uniform css extraction, and LLM prompt engineering. While it provides these features, is highly flexible allowing this application to effortlessly use and extend it.

APScheduler

APScheduler[51] (Advanced Python Scheduler) is a lightweight, flexible, and efficient job scheduling library for Python. It allows developers to schedule tasks (jobs) to run at specific intervals, on a fixed date, or in response to events. APScheduler is commonly used for task automation, cron-like scheduling, and background job execution, making it an ideal choice for scheduling web scraping jobs in the web scraping system.

- 1. Multiple Job Scheduling Options
 - One-time jobs \rightarrow Execute a task at a specific datetime.
 - Interval-based jobs \rightarrow Run tasks every X seconds/minutes/hours.
 - Cron-like jobs \rightarrow Schedule jobs using cron expressions for precise timing.
- 2. Persistent Job Storage
 - Supports SQLite, PostgreSQL, MySQL, and MongoDB for storing scheduled jobs.
 - Ensures that jobs persist across application restarts.
- 3. Asynchronous & Threaded Execution
 - Supports multi-threaded execution (default) and async job execution.
 - Ensures that scheduled tasks do not block the main application thread.
- 4. Error Handling & Job Monitoring
 - Logs job failures, execution times, and retries.
 - Provides job state tracking, including paused, running, and completed jobs.
- 5. Flexible Job Triggers
 - Uses built-in triggers for different scheduling needs (date, interval, cron).
 - Allows custom event-based triggers.

6.2.2 Frontend

The frontend of the application has been built on two pillars of modern software design: user experience and extensibility. Following production-grade practices, the UI counterpart of this application attempts to be user friendly, intuitive while being attentive to performance and productivity. The technologies utilized to achieve this are the following:

React

React[52] is a JavaScript library for building user interfaces, developed and maintained by Meta (formerly Facebook). It is widely used for single-page applications (SPAs) and component-based UI development. React is particularly well-suited for dynamic and interactive web applications, making it an ideal choice for the frontend of the web scraping system.

React's virtual DOM (Document Object Model) and declarative programming paradigm ensure high performance and efficient UI updates, making it perfect for rendering and managing user-defined scraping jobs.

- 1. Component-Based Architecture
 - React applications are built using reusable components, allowing modular and maintainable code.
 - Each component manages its own state and can be reused across the application.
- 2. Virtual DOM for Efficient Rendering
 - React uses a virtual DOM to efficiently update only the necessary UI elements when the state changes.
 - Reduces the number of actual DOM manipulations, improving performance and responsiveness.
- 3. Declarative UI
 - React describes the UI state declaratively, making code easier to read and debug.
 - Instead of manually updating the UI, React updates the view when data changes.
- 4. React Hooks for State and Effects
 - Introduces functional components with React Hooks (useState, useEffect, etc.).
 - Simplifies state management and side effects handling without using class components.
- 5. Context API for Global State Management
 - Built-in alternative to Redux, allowing global state management without prop drilling.
- 6. React Router for Navigation
 - Enables client-side routing, making navigation fast and seamless without full page reloads.
- 7. Integration with APIs
 - Works efficiently with RESTful APIs (FastAPI backend), GraphQL, and Web-Sockets.

Refine.dev

Refine.dev[53] is an open-source React framework for building data-intensive applications. It is designed to simplify CRUD (Create, Read, Update, Delete) operations, making it a powerful tool for managing structured data applications like admin dashboards, internal tools, and scraping monitoring systems.

Refine.dev is an ideal choice for the web scraping system's frontend, as it allows quick development of dashboards for job management, scheduling, and data monitoring.

- 1. Rapid CRUD API Integration
 - Automatically connects to REST or GraphQL APIs.
 - Provides hooks (useList, useCreate, useUpdate) to fetch and manipulate data.
- 2. Auto-Generated Admin UI
 - Generates CRUD interfaces automatically from API schemas.
 - Provides table views, forms, and filters with minimal configuration.
- 3. Authentication & Authorization
 - Supports JWT authentication, OAuth, and third-party login providers.
 - Provides role-based access control (RBAC).
- 4. Powerful Data Fetching with React Query
 - Uses React Query to fetch and cache API data efficiently.
 - Provides real-time updates and optimistic UI.
- 5. React Router & Navigation
 - Seamlessly integrates with React Router for SPA navigation.
 - Provides built-in support for nested routes and protected pages.
- 6. Custom Hooks for Business Logic
 - Developers can override default behavior using custom hooks.

Refine is built on top of React, React Query, and other modern libraries, providing outof-the-box support for many application features. Adding opinion to common repetitive layers that are crucial for each application make the code more readable, efficient file structure and many other features that prevents the developer from reinventing the wheel while making the fronted of their application scalable

ShadCN UI

ShadCN UI[54] is a modern, accessible, and customizable component library for React applications. Unlike traditional UI libraries (e.g., Ant Design, Material UI), ShadCN provides unopinionated, developer-friendly components built on Tailwind CSS and Radix UI.

ShadCN UI is ideal for building beautiful, lightweight, and highly customizable UIs, making it perfect for the frontend of the web scraping system.

- 1. Tailwind CSS-Based Styling
 - Uses Tailwind utility classes for styling instead of complex stylesheets.
 - Allows full customization without the need for additional CSS overrides.
- 2. Headless & Accessible Components
 - Uses Radix UI for keyboard navigation and screen reader support.
 - Components are unstyled by default, giving developers full control over UI styling.
- 3. Server-Side Rendering (SSR) Friendly
 - Fully compatible with Next.js and React Server Components.
 - Supports progressive enhancement for fast page loads.
- 4. Dark Mode Support
 - Auto-detects system themes and allows manual toggling between light and dark mode.
- 5. Lightweight & Fast
 - Unlike Ant Design or Material UI, ShadCN does not ship unnecessary CSS or JS.
 - Works out of the box with Vite, Next.js, and React.
- 6. Component Reusability
 - Components are imported per use case, reducing bundle size and improving performance.

6.3 Architecture

The architecture is built to accommodate a broad range of users, from non-technical individuals requiring a no-code solution to advanced users who may customize and extend its functionality. The system consists of three primary layers: the frontend, which enables users to interact with the system through a well-structured user interface; the backend, which is responsible for API handling, job scheduling, ensures the persistence of extracted data and scraper coordination; and the proxy layer, which creates an isolated network between the other layers while exposing it to the rest of the world. Each of these layers operates within its own dedicated Docker container [55], ensuring a clear separation of concerns and allowing for tailored environmental configurations based on their specific requirements. Despite this isolation, the proxy layer acts as an intermediary, seamlessly managing communication between containers within their private network while also interfacing with the external environment. Once each service is individually configured, Docker Compose [56] orchestrates the multi-container deployment, synchronizing the lifecycle of all components to maintain system coherence and operational efficiency.

This architecture follows a modular monolithic approach, where different components interact within a unified backend system without being entirely decoupled as separate services. The decision to follow this design was motivated by the need to balance simplicity, maintainability, and efficiency resulting in a streamlined and performant architecture, wellsuited for handling structured data extraction tasks.

The following sections provide a detailed breakdown of each architectural component, its role within the system, and the interactions between different layers.

6.3.1 Frontend Layer

The frontend of the system is designed to provide a modern, responsive, and interactive user interface that allows users to create, manage, and monitor their scraping jobs. It is developed using React, leveraging Refine.dev for simplified CRUD operations and ShadCN UI for a clean and customizable design. This choice of technologies ensures a balance between usability, flexibility, and efficiency.

The primary responsibilities of the frontend include allowing users to define ontologies, specifying the structure of the data they wish to extract, as well as configuring scraping jobs by providing target URLs and selecting extraction methodologies. The frontend also facilitates data visualization, enabling users to inspect extracted data through interactive dashboards. Authentication and authorization mechanisms are implemented to ensure secure access control, allowing only registered users to define and manage scraping operations.

The frontend interacts with the backend through a RESTful API, exchanging structured JSON data. Data related to scraping jobs, ontologies, and extracted content is retrieved dynamically from the backend and displayed in the user interface. Secure communication is facilitated through Traefik, which acts as a reverse proxy, ensuring proper request routing and enforcing authentication policies.

6.3.2 Backend Layer

The backend is implemented using FastAPI, a high-performance web framework optimized for asynchronous execution. It is responsible for securely handling user requests, managing scraping jobs, scheduling extractions, and processing data interactions with the database. The backend exposes a REST API, which the frontend consumes to enable users to create and manage their scraping workflows.

The core functionality of the backend revolves around ontology management, scraping job scheduling, and data retrieval. Users can define custom ontologies that structure the data extraction process, specifying the type of information they need to extract. These ontologies are then used to define page schemas, which determine the specific extraction rules for each web page. The backend ensures that these configurations are properly stored and retrieved when needed.

soniq ^{0.1.0} ^{0AS 3.1}

	Authorize 🔒
login	^
POST /api/v1/login/magic/{email} Login With Magic Link	~
POST /api/v1/login/claim Validate Magic Link	
POST /api/v1/login/oauth Login With Oauth2	v
POST /api/v1/login/refresh Refresh Token	
POST /api/v1/login/revoke Revoke Token	≜ ∨
POST /api/v1/login/recover/{email} Recover Password	~
POST /api/v1/login/reset Reset Password	
users	^
GET /api/vl/users/ Read User	≙ ∨
PUT /api/vl/users/ Update User	≙ ∨
POST /api/v1/users/ Create User Profile	~
GET /api/vl/users/all Read All Users	≙ ∨
POST /api/v1/users/new-totp Request New Totp	≜ ∨
POST /api/vl/users/toggle-state Toggle State	≜ ∨
POST /api/v1/users/create Create User	≜ ∨
GET /api/vl/users/tester Test Endpoint	~

Figure 6.2. soniq API Swagger Docs 1

To handle scheduled scraping jobs, the system utilizes APScheduler, which periodically triggers web scraping tasks based on user-defined schedules. Unlike architectures that rely on distributed task queues, the scheduler directly invokes the Crawl4AI scraper, ensuring that scraping jobs are executed efficiently without unnecessary complexity. Once a job is initiated, the scraper fetches the target web pages, applies the specified extraction schema, and stores the extracted data in the MongoDB database.

Data retrieval is another key function of the backend. Users can request previously extracted data through dedicated API endpoints. The system supports advanced filtering and querying mechanisms, allowing users to extract insights from large datasets efficiently. Logs of all executed jobs are maintained, enabling users to track scraping performance and debug any potential issues.

Web Scraping Engine

The web scraping functionality is powered by Crawl4AI, an AI-enhanced scraping engine that enables both rule-based and AI-driven extraction methods. Traditional web scrapers rely on static CSS selectors to extract specific elements from web pages, making them fragile in the face of website layout changes. Crawl4AI, however, incorporates machine learning techniques to identify structured data patterns dynamically, improving its ability to adapt to changes in webpage structures.

Scraping can be performed in two modes: an AI-assisted mode, where the system automatically detects and extracts structured content, and a manual mode, where users

ргоху	^
POST /api/v1/proxy/{path} Proxy Post Request	ê v
GET /api/v1/proxy/{path} Proxy Get Request	≜ ∨
ontology	^
POST /api/v1/ontology/ Create Ontology	â V
PUT /api/vl/ontology_id} Update Ontology	â V
GET /api/vl/ontology_id} Get Ontology	~
POST /api/v1/ontology/List List Ontologies	≜ ∨
page	^
POST /api/v1/page/ Create Page	≜ ∨
PUT /api/vl/page_id} Update Page	≜ ∨
GET /api/vl/page/{page_id} Get Page	~
POST /api/v1/page/list List Pages	â V
job	^
GET /api/vl/job_id} Get.ucb	~
DELETE /api/v1/job/{job_id} Dukeo Job	~
POST /api/vl/job/list ListJobs	~
run	^
GET /api/v1/run/{run_id} Get Run	~
POST /api/vl/run/list List Puns	≜ ∨
data	^
GET /api/vl/data/{data_id} GetData	~
POST /api/vl/data/list Ust Data	â V

Figure 6.3. soniq API Swagger Docs 2

define CSS selectors to specify exact elements to be extracted. This flexibility allows the system to cater to both technical and non-technical users.

Once a scraping task is triggered, Crawl4AI fetches the requested web pages and processes them according to the defined ontology and extraction schema. The extracted data is validated, structured, and stored in MongoDB. In cases where an error occurs during extraction, the system logs the failure and provides debugging information to the user.

Data Storage and Processing

All structured and unstructured data collected by the scraper is stored in MongoDB, a NoSQL document-oriented database optimized for handling large-scale web data. The database schema is designed to support the flexible nature of the data being extracted while maintaining data consistency and fast retrieval times. The primary collections within the MongoDB instance include:

Ontologies, which define the structure of extracted data.

Pages, which store webpage configurations and their corresponding extraction schemas.

Jobs , which schedule the scraping tasks.

JobRuns, which tracks the execution of scraping tasks.

ExtractedData , which contains the final structured output from web scraping operations.



Figure 6.4. soniq Database UML Class Diagram

6.3.3 Deployment Architecture

The system is designed to be containerized, ensuring that all components can be easily deployed and managed. The deployment is structured around Docker, where each component runs as an isolated container. Traefik is used as a reverse proxy to route requests securely between the frontend and backend while also handling SSL termination and load balancing. The key components of the deployment include:

Frontend (React + Refine.dev + ShadCN UI), deployed as a Docker container.

Backend (FastAPI + APScheduler + Crawl4AI), handling API requests and job execution.

MongoDB, which persists all scraping configurations and extracted data.

Traefik, which manages routing and API exposure.

The deployment architecture ensures scalability and maintainability, and easy deployment allowing components to be updated independently without affecting the entire system.



Figure 6.5. soniq Deployment UML Diagram

6.3.4 System Workflow

The system follows a well-defined workflow for executing scraping tasks. First, the user defines an ontology and configures the extraction schema for a target web page. This configuration is stored in MongoDB, ensuring that the system retains the extraction rules for future jobs. Next, the user schedules a scraping job, which is handled by APScheduler and triggered at the specified time. When execution begins, the Crawl4AI scraper retrieves the target page, extracts data according to the ontology, and stores the results in MongoDB. The user can then access the extracted data through the frontend interface.

This workflow enables an end-to-end automated scraping pipeline, from configuration and scheduling to data retrieval and visualization. The combination of AI-assisted extraction, manual selector-based scraping, and task scheduling ensures that users can extract data efficiently without requiring deep technical knowledge.



Figure 6.6. soniq Activity UML Diagram



Figure 6.7. soniq Sequence UML Diagram

6.4 Use Case - Scraping Energy News Articles in Seconds

In this section, one will analyze the usage described above applied to the energy news domain. Extracting news from websites is a very common practice that provides researchers with powerful data points that can be utilized to train AI models, gain insights about current advancements, or even maintain a system up-to-date by repeating the extraction process every day.

6.4.1 Creating the Ontology

The first step of scraping a semantic domain in "soniq" is to create the desired ontology. In the context of the implementation of this thesis, an ontology is the basic entity that describes and structures the extracted data.



Figure 6.8. Ontology Creation

Once an ontology is created, the user can continue defining the websites they desire to extract data from.

6.4.2 Creating a Page

A page is a secondary entity that refers to an existing ontology. During its creation, the user needs to define the URL from which "soniq" will source the data, the extraction schema to be followed, and the frequency in which the job will run.

Step 1: Basic Information



Figure 6.9. Basic Page Information

Step 2: Schema Generation

In the next step of the workflow, the user can optionally find the css selector of the HTML element that contains the data they want to extract. This process is only useful to limit the tokens sent to the LLM and therefore the price of the request, but can also be skipped in cases where the user is not so technical.



Figure 6.10. Inspect HTML Container

Once this process is finished and the user has found a potential css selector, they can move on to querying the LLM to generate the actual css selectors that can successfully retrieve the data structure defined in the ontology.

soniq								
	Page > Create >							
	2. CSS Schema Generation							
e_	.post-wrap Generate							
	Egr energia.gr							
	<pre>> Crawling https://www.energia.gr/eidiseis/1/list < Successfully crawled in 7.49 seconds < Cleaned htal in 0.43 seconds</pre>							
	> Generating css extraction schema							
	V Generated CSS extraction schema in 3.45 Seconds Base Selector							
	article.post							
	CSS Schema							
	title	h3 a	• •					
	date	.post-meta .time	⊕ 🛍					
	summary	p.excerpt-entry	⊕ 🗇					
	content	div.content	⊕ ₫ ^					
	title	h3 a	⊕ til					
	date	.post-meta .time	⊕ [†]					
	summary	p.excerpt-entry	⊕ m̂					
	Add Field							
Test								
	Back		Next					

Figure 6.11. CSS Extraction Schema Generation using LLM

One can notice that the LLM has hallucinated and added fields that the user did not specify in the ontology (content has nested fields). For that reason, the generated schema is completely modifiable in order to add, edit, or remove fields that have not been generated correctly.

soniq			
<>	Page > Create >		
\otimes	2. CSS Schema Generation		
€	.post-wrap Generate		
	E₊rr energia.gr		
	 > Crawling https://www.energia.gr/eidiseis/1/list > Successfully crawled in 7.49 seconds > Cleaned html in 0.03 seconds > Generating css extraction schema > Generated css extraction schema in 3.45 seconds 		
	Base Selector		
	article.post		
	CSS Schema		
	title	h3 a	ن ال
	date	.post-meta .time	⊙ ti
	summary	p.excerpt-entry	 ⊕
	Add Field		
		Test	
	Back		Next

Figure 6.12. Modified Generated Schema

In order to verify the LLM results, the user can simulate the scraping job and retrieve the data using the generated strategy by clicking the "Test" button. By viewing the extracted data, the user can verify if there is a field without content and can cross-validate using the real website that every data point they needed to scrape has been successfully extracted.

soniq		Test	
	Title	Date	Summary
\$ ()) ())	Για Ενέργεια και Στρατιωτική Συνεργασία Συνομίλησαν Τηλεφωνικά οι ΥΠΕΞ ΗΠΑ και Κυπριακής Δημοκρατίας	Τετ, 5 Μαρτίου 2025 - 19:51	Την εκτίμηση του για τον σημαντικό ρόλο που επιτελεί η Κυπριακή Δημοκρατία στην Ανατολική Μεσόγειο και για τη βοήθειά που παρέχει στις επιχείρησης εκκένωσης Αμερικανών πολιτών από ζώνες περιφερειακών συγκρούσεων εξέφρασε ο υπουργός Εξωτερικών των ΗΠΑ Μάρκο Ρούμπιο κατά τη διάρκεια τηλεφωνικής συνομιλίας που είχε με τον Κύπριο ομόλογο του Κωνσταντίνο Κόμπο
°o	Cenergy Holdings: Αύξηση Λειτουργικής Κερδοφορίας Κατά 27% το 2024 – Μεγαλύτερο Κατά 75% το Προτεινόμενο Μέρισμα για Εφέτος	Τετ, 5 Μαρτίου 2025 - 19:44	Το 2024, η Cenergy Holdings συνέχισε να αξιοποιεί τη διαρκώς αυξανόμενη ζήτηση στον ενεργειακό τομέα. Η ζήτηση για προϊόντα καλωδίων παρέμεινε ικανοποιητική συμβάλλοντας στη στήριξη των τιμών, ενώ τα ενεργειακά έργα εκτελέστηκαν ομαλά, αποτελώντας τον κύριο πυλώνα της αύξησης της κερδοφορίας του τομέα καλωδίων. Ο τομέας σωλήνων χάλυβα πέτυχε καλύτερες επιδόσεις από το 2023, λόγω βελτιωμένων περιθωρίων κέρδους, που προέκυψαν από το μείγμα εκτελεσμένων έργων
	ΗΠΑ: Ο Νέος Εμπορικός Πόλεμος του Τραμπ και τα Αντίμετρα των Άλλων Χωρών	Τετ, 5 Μαρτίου 2025 - 19:23	Ο δεύτερος εμπορικός πόλεμος του Ντόναλντ Τραμπ είναι πλέον πραγματικότητα και οι πρώτες επιπτώσεις έχουν ήδη καταγραφεί. Με τους εμπορικούς εταίρους των ΗΠΑ να γνωρίζουν πλέον τη νοοτροπία του Αμερικανού Προέδρου, οι απαντήσεις δεν άργησαν να έρθουν, εξασφαλίζοντας πως αυτός ο πόλεμος θα έχει παράπλευρες απώλειες. Με τη διεθνή αβεβαιότητα να αυξάνεται, οι αγορές λαμβάνουν ως δεδομένο ότι η κατάσταση θα γίνει πολύ χειρότερη πριν βελτιωθεί
	Prodexpo North 2025: Τα Κτίρια της Θεσσαλονίκης Αλλάζουν Μορφή - Η Σημασία των Πράσινων Κτιρίων	Τετ, 5 Μαρτίου 2025 - 19:16	«Στόχος μας ήταν να δημιουργήσουμε ένα κτίριο που θα έχει ευελιξία στην χρήση, θα λειτουργεί περισσότερο ως χώρος προορισμού και όχι μόνος ως χώρος εργασίας κι εκτιμώ πως το καταφέραμε», τόνισε η Evie Leon, Διευθύνουσα Σύμβουλος, Deutsche Telekom Cloud Services, στη συζήτηση που πραγματοποιήθηκε στο πλαίσιο του 8ου συνεδρίου της Prodexpo North, στη Θεσσαλονίκη
	Την Πώληση του Πλοίου «Κρήτη ΙΙ» για Ανακύκλωση, Έναντι 3,6 Εκατ. Δολ. Ανακοίνωσε η Attica Group	Τετ, 5 Μαρτίου 2025 - 18:43	Την πώληση του πλοίου Κρήτη ΙΙ, έναντι 3,6 εκατ. δολ., με σκοπό την ανακύκλωσή του, ανακοινώνει η Attica AE Συμμετοχών
	Νετανιάχου: «Είμαστε Αποφασισμένοι να Νικήσουμε»	Τετ, 5 Μαρτίου 2025 - 18:20	Ο πρωθυπουργός του Ισραήλ Μπενιαμίν Νετανιάχου δήλωσε σήμερα στον νέο αρχηγό του ισραηλινού στρατού ότι είναι αποφασισμένος να οδηγήσει την χώρα στην νίκη σχεδόν 17 μήνες μετά την έναρξη ενός πολυμέτωπου πολέμου που ξεκίνησε η επίθεση της Χαμάς κατά του Ισραήλ στις 7 Οκτωβρίου 2023
	Οι Καλλιέργειες Τροφίμων σε Κίνδυνο Από την Αύξηση της Θερμοκρασίας	Τετ, 5 Μαρτίου 2025 - 17:29	Καθώς οι θερμοκρασίες συνεχίζουν να αυξάνουν, το ένα τρίτο της παγκόσμιας παραγωγής τροφίμων θα μπορούσε να τεθεί σε κίνδυνο. Αυτό διαπιστώνει μελέτη που δημοσιεύθηκε στο περιοδικό «Nature Food»
	Η Επέτειος του Γιάννου Κοντόπουλου: Το Πλήρωμα του Χρόνου και το +80% στην 3ετία	Τετ, 5 Μαρτίου 2025 - 17:17	Τρία χρόνια συμπληρώνει σε λίγες ημέρες στο «τιμόνι» του Χρηματιστηρίου Αθηνών ο Γιάννος Κοντόπουλος, ο οποίος βλέπει τον μεγάλο στόχο της αναβάθμισης στις Αναπτυγμένες Αγορές να πλησιάζει ολοένα και περισσότερο, σηματοδοτώντας την πλήρη επιστροφή στην κανονικότητα. Κι αυτό ενόσω η εγχώρια κεφαλαιαγορά πραγματοποιεί σταθερά βήματα για να «ξεχάσει» τελείως την εποχή των διαδοχικών κρίσεων και της πολυετούς οικονομικής περιπέτειας
	Θεσσαλονίκη: Σε Εξέλιξη Τροποποιήσεις της Μελέτης για το Ενιαίο Παραλιακό Μέτωπο Εν Αναμονή του ΠΔ	Τετ, 5 Μαρτίου 2025 - 16:53	Τροποποιήσεις της μελέτης για την ανάπλαση του παραλιακού μετώπου της Θεσσαλονίκης βρίσκονται σε εξέλιξη, όπως ανέφερε σήμερα ο αντιπεριφερειάρχης Υποδομών και Δικτύων της Περιφέρειας Κεντρικής Μακεδονίας Πάρις Μπίλλιας, μιλώντας στο 8ο συνέδριο Prodexpo North, με αντικείμενο την αξιοποίηση της ακίνητης περιουσίας στη βόρεια Ελλάδα
	ΒΕΑ: Φορολογία και Κόστος Ενέργειας τα	Τετ, 5 Μαρτίου	Η υψηλή φορολογία, με βάση το τεκμαρτό εισόδημα, το υψηλό μη μισθολογικό κόστος και οι συσσωρευμένες οφειλές του παρελθόντος, εξακολουθούν να αποτελούν «τροχαπέδη για την ανάπτυξη των βιοτεχνικών

Figure 6.13. Scraping Job Simulation

Step 3: Job Scheduling

It is now time to schedule the extraction frequency. This is the last step of the page creation workflow and once it is completed, the scraping job is scheduled using the aforementioned technologies.

soniq						
<>	Page > Create >					
⊗	3. Schedule th	ne scraping job				
8						
2 0	Minute	Hour	Day of Month	Month	Day of Week	
	0	~ 10	✓ Every day	~ Every month	✓ Every day	
	At 10:00 Alv Next run time	1 e: 06/03/2025, 10:00:00				
	Back					Next

Figure 6.14. Scraping Job Scheduling

The user can verify that the page has been created and explore other pages in the exploration data table.

Figure 6.15. Page Explorer

Step 4: Handle pagination

Pagination is the most common practice of limiting data availability on all types of websites that expose valuable information. For this purpose, "soniq" manages to resolve this issue by providing the user the ability to define the pagination variables. It is visible from the browser's url input that the current website the user intends to scrape adds a query parameter named **p** in order to handle pagination.



Figure 6.16. energia.gr Pagination Url

Once the user has identified all the required variables to handle pagination in the website of their choice, they can edit the page and configure it to enqueue more urls for extraction. Chapter 6. soniq: No-code web scraping platform for structured data extraction



Figure 6.17. Pagination Url Enqueuer

6.4.3 Inspecting Job Executions

In the last section, the user created and scheduled a page. The "soniq" platform has the responsibility to enqueue this scraping job when the user has defined. An additional feature is the job execution exploration and detailed review that the platform provides to the user to monitor and debug their defined jobs.



Figure 6.18. Run Explorer

The detailed review of the job execution provides valuable information to the user, such as its status, the number of items successfully scraped, the date it took place, its duration, and for advanced users the actual service logs in order to identify potential errors.

ob Details w details and res	S sults of your web scraping job			Succes
Job Overvie	ew			Copy ID
Basic information	about the scraping job			
Job ID:	67c8bfcce39f357c1e0a54a5	Date:	🛱 Mar 5, 2025, 09:20:12 PN	4
Status:	success	Duration:	I7.90 seconds	
Items Scraped:	52	Schedule:	0 10 * * *	
Target Infor Details about the I	mation target website and scraping configuration s://www.energia.gr/eidiseis/1/list			
Target Infor Details about the t URL: https Ontology: 1 Schema Confi	rmation target website and scraping configuration :://www.energia.gr/eidiseis/t/list Energy news – News articles about energy iguration			
Target Infor Details about the I URL: https Ontology: I Schema Confi Execution L Detailed logs of th	mation target website and scraping configuration :://www.energia.gr/eidiseis/1/list Energy news - News articles about energy iguration LOGS Le scraping process			
Target Infor Details about the I URL: https Ontology: Schema Confi Execution L Detailed logs of th View Execution	rmation target website and scraping configuration :://www.energia.gr/eidiseis/1/list Energy news - News articles about energy iguration cogs es scraping process			

Figure 6.19. Job Execution Details

6.4.4 Scraped Data Exploration

Finally, the most important feature of this platform is the ability to view, explore, and export scraped data to analyze and manage the data.

soniq						
0	note >					
ē	Energy news					
8						Refresh Create
	www.smergia.gr	«Εξοικονομώ 2025»: Με Ξεχειριστά Κίνητρα Εντάσσονται οι Πληγέ	Ter, 5 Maptiou 2025 - 16:36	Με νέα απόφαση του Υπουργού Περιβάλλοντος και Ενέργειας, κ.	0	
	www.energia.gr	Το Σχέδιο άράσης της ΕΕ για το Μέλλον της Αυτοκινητοβιομηχανία	Tet, 5 Maptiou 2025 - 16:30	Το σχέδιο δράσης της ΕΕ για το μέλλον της αυτοκινητοβιομηχανία	34	
	www.eoergia.gr	Στις Κάλπες η Γροιλανδία: Ανεβαίνει Ξαφνικά το Κίνημα Ανεξαρτησ	Tet, 5 Maptiou 2025 - 16:04	Οι ψηφοφόροι ζυγίζουν το πο σημαντικό θέμα των εκλογών τους:		
	www.energia.gr	Η ΠΟΕΠΗΕΦ Ζητά Παρέμβαση της ΡΑΑΕΥ για Αλλαγές στην Υπουρ	Ter, 5 Maptiou 2025 - 15:33	Με σημερινή επιστολή της προς την Ρυθμιστική Ακεξάρτητη Αρχ	n	
	www.energia.gr	Ρωσική Επίθεση σε Ενεργειακή Μονάδα στην Οδησσό	Ter, 5 Maptiou 2025 - 15:19	Ρωσικές δυνάμεις επιτέθηκαν σε ενεργειακή μονάδα στην περιφέ	þ	
	www.energia.gr	Τράπεζα Πειραιώς: Με Στρατηγική Συνεργασία με την ΜΑΝΤΙLITY,	Ter, 5 Mapricu 2025 - 15:15	Η Τράπεξα Παρακώς, στο πλαίσιο της Εταιρικής της Υπευθυνότητ	10	
	www.snergia.gr	Eurelectric: Η Ευελιξία των Η/Ο Εξοικονομεί Χρήματα για τους Πελι	Ter, 5 Mapricu 2025 - 15:15	Ένας ιδιοκτήτης ηλεκτρικού οχήματος στην Ευρώπη θα μπορούσ	8	
	www.energia.gr	Πολυμελός Πρωτοδικείο Κοζάνης: Απογορεύει την Κατασκευή Φ/Β	Tet, 5 Mapticu 2025 - 15:00	Με την υπ' αριθμ. 12/2025 απόφαση του Πολυμελούς Πρωτοδικοί	0	
	www.energia.gr	Επίτιμη Διδάκτορος του Τμήματος Φυσικής του ΑΠΘ Αναγοροίεται	Ter, 5 Mapricu 2025 - 14:55	Επίτιμη διδάκτορας του Τμήματος Φυσικής της Σχολής Θετικών Ε		
	www.energia.gr	Το Κρεμλίνο Χαιρετίζει την Δήλωση Ζελένσκι για την Προθυμία του	Ter, 5 Maptiou 2025 - 14:48	Η Ρωσία χαιρετίζει τη δήλωση του προέδρου Βολοντίμερ Ζελάνσκι		
	www.energia.gr	13η Solaire Expo Maroc: Αυξημένο Ενδιαφέρον των Ελληνικών Επη	Ter, 5 Mapticu 2025 - 14:43	Από τις 25 έως τις 27 Φεβρουαρίου 2025 πραγματοποιήθηκε στη		
	www.coorpia.gr	Γιατί Είναι Τόσο Σημαντικός οι Σπάνιες Γαίες της Οικρανίας (και Ο	Ter, 5 Mapricu 2025 - 14:05	Την περασμένη Παρασκευή, άλος ο πλανήτης side live τον «εξευτ	z	
	www.energia.gr	Γερμανία: Βουτιά 76% για τις Πωλήσεις των Tesla τον Φιέβρουάριο	Ter, 5 Mapricu 2025 - 13:51	Η αμερικανική αυτοκινητοβιομηχανία Tesla συνέχισε να χάνει έδα	ø	
	www.energia.gr	ΕΥ: Οι CEOs Παγκοσμίως Παραμένουν Αισιόδοξοι, «Βλέποντας» Νι	Tet, 5 Maptiou 2025 - 13:46	Η εμπιστοσίνη των CEOs στην ανάπτυξη κερδίζει άδοφος παγκοσ	n	
	www.energia.gr	Wood για ΓΕΚ ΤΕΡΝΑ: Αναβαθμίζει την Τιμή-Στόχο στα €27 με 489	Ter, 5 Maptiou 2025 - 13-32	Στα 27 ευρώ ανά μετοχή (από 18 ευρώ) αυξάνει την τιμή-στόχο η		
	www.energia.gr	Ο Fabio Quaranta είναι ο νέος Δευθύνων Σύμβουλος της Zeniθ	Tex, 5 Maptiou 2025 - 13:20	Η Zenið, η θυγατρική εταιρεία της Plenitude στην Ελλάδα με ηγει		
	www.energia.gr	ΕΕΑΑ: Δυναμικό «παρών» στην 7η Δικθνή Verde.Tec - Αροσίωση σ	Ter, 5 Mapricu 2025 - 13:05	Η Ελληνική Εταιρεία Αξιοποίησης Ανακύκλωσης (ΕΕΑΑ) συμμετιές	a a a a a a a a a a a a a a a a a a a	
	www.energia.gr	Βιώσιμη Ανάπτυξη στην Ελλάδα: Το Νέο Δεδομένο & Τάσεις στην Ε	Tet, 5 Mapticu 2025 - 12:49	Το Κέντρο Αειφορίας (CSE) παρουσιάζει την 12η ετήσια έρευνά τ		
	www.energia.gr	Η INTRACOM DEFENSE και η ΕΤΙΜΑD Ενώνουν τις Δυνάμεις τους	Tet, 5 Maptiou 2025 - 12:25	Η Intracom Defense (IDE) υπέγραψε Μνημόνιο Συνεργασίας με τη		
	www.energia.gr	Νέα Οδός - Κεντρική Οδός: Χρισή Διάκριση στα ΑΙ & DATA Award:	Ter, 5 Mapriou 2025 - 12:01	Μία ακόμα Χρυσή διάκριση μετρούν η Νέα Οδός και η Κεντρική Ο		
			I← < → → Page 2 of 3 Rows per pag	e 20 v Total Rows 42		

Figure 6.20. Data Explorer



Epilogue


6.5 Overview

The journey through advanced web scraping in the modern web has underscored the growing complexities and innovations shaping the field. This thesis has explored both offensive (scraping methodologies) and defensive (anti-scraping techniques) perspectives, while also proposing an alternative vision: the democratization of web scraping through a no-code AI-enhanced platform.

As the arms race between scrapers and anti-bot defenses continues, this work highlights the potential of AI-driven adaptability in overcoming technical barriers while maintaining ethical and legal compliance. The introduction of LLM-assisted schema generation, adaptive pipeline correction, and automated API-based extraction provides a path toward a more accessible and structured approach to data extraction. The implementation of soniq demonstrates how open-source, AI-enhanced web scraping can level the playing field, allowing smaller teams, researchers, and independent developers to access structured web data without extensive engineering expertise.

Looking forward, the interplay between web security, AI, and data accessibility will only grow more intricate. The future of web scraping lies in intelligent, ethical, and highly adaptive solutions that balance efficiency with compliance. Further research should focus on enhancing AI-driven automation, improving adversarial robustness, and refining realtime scraping frameworks.

Ultimately, this work serves as a bridge between technological advancement and accessibility, reaffirming that structured web data should not be a privilege of large-scale enterprises but an open and democratized resource for all.

6.6 Future Work

While this thesis has laid a strong foundation for AI-driven, no-code web scraping democratization, several avenues for future research and development remain open. The evolution of web technologies, security measures, and AI capabilities presents both challenges and opportunities for advancing the field of intelligent data extraction. Below are key areas for potential future work:

6.6.1 Enhancing AI-Driven Adaptability in Scraping

One of the most critical areas for future development is increasing the adaptability of web scraping systems to handle dynamic and ever-changing web environments. Websites frequently update their DOM structures, anti-bot measures, and data formats, often rendering traditional scrapers obsolete. While this thesis has introduced LLM-assisted schema detection and self-correcting pipelines, further work is needed to develop autonomous scraping agents that can learn and evolve over time.

Future advancements should focus on self-learning scrapers that leverage reinforcement learning and continual adaptation techniques. These agents should monitor UI changes over time, detect patterns of modifications, and automatically adjust their extraction strategies without human intervention. Additionally, the integration of multimodal AI, combining computer vision for layout interpretation with NLP for text structuring, can significantly enhance the robustness of web scrapers. This hybrid approach would allow scrapers to navigate web interfaces more naturally, much like a human user would, ensuring accurate data extraction even in the face of sophisticated obfuscation techniques.

Beyond adaptability, another key area for future work is optimizing AI-driven scraping for speed and efficiency. Current AI-assisted extraction models introduce higher computational overhead, making real-time scraping impractical for large-scale applications. Research into lightweight AI models, edge processing techniques, and optimized inference pipelines will be necessary to enable high-speed, cost-effective AI-driven scraping at scale.

6.6.2 Expanding No-Code Customization & User Experience

Democratizing web scraping requires further development in no-code and low-code solutions, making data extraction more accessible to researchers, business analysts, and non-technical users. While this thesis introduces a no-code scraping platform, future work should focus on enhancing user experience, automation, and customization to create a seamless and intuitive data extraction process.

One of the most impactful directions for expansion is the development of a visual ontology builder, where users can define structured data schemas through an interactive, drag-and-drop interface. This would eliminate the need for manual ontology scripting, allowing users to visually map relationships between extracted data fields and generate structured outputs without writing a single line of code. Such a feature would significantly lower the entry barrier for users who need structured data but lack technical expertise in web scraping.

Additionally, integrating AI-driven guidance into the no-code interface can help users automatically generate optimized scraping configurations based on natural language descriptions. Users should be able to describe their data needs conversationally, and the system would generate, test, and validate the necessary extraction workflows. Furthermore, expanding the multi-domain capabilities of the platform—such as enabling seamless scraping across e-commerce, finance, academia, and social media—will ensure that different industries can easily adopt structured web data extraction without requiring extensive customization.

Finally, integrating real-time monitoring dashboards that provide insights into scraping efficiency, data quality, and system performance will enhance usability. These dashboards should allow users to identify errors, track historical trends, and receive automated suggestions for improving scraping efficiency. Future research should explore how real-time analytics and visualization tools can empower non-technical users to make data-driven decisions while keeping scraping operations efficient and transparent.

6.6.3 Exploring Decentralized & Federated Scraping Approaches

Web scraping currently relies heavily on centralized architectures, where individual scrapers operate in isolation or within controlled server environments. However, decentralized and federated approaches offer exciting possibilities for scalability, anonymity, and data resilience while reducing the risk of centralized failures and legal vulnerabilities.

One promising avenue for future work is the development of peer-to-peer (P2P) scraping networks, where users contribute scraping resources in a distributed fashion, effectively creating a crowdsourced data extraction ecosystem. In such a system, scrapers can collaborate, sharing learned extraction patterns while distributing crawling workloads across a decentralized network. This would not only improve scraping efficiency but also make it harder for websites to detect and block scraping operations, as requests would originate from diverse, geographically dispersed nodes rather than a single centralized entity.

Additionally, federated learning techniques could be applied to web scraping, enabling scrapers to collaboratively improve extraction models without sharing raw data. This method would allow different scraping agents to learn from each other's experiences, refining AI-assisted extraction techniques in a privacy-preserving manner. Such an approach could lead to a more adaptive, intelligent, and community-driven scraping ecosystem.

6.6.4 Ethical and Regulatory Considerations

The growing legal and ethical concerns surrounding web scraping necessitate further research into compliance frameworks and responsible data extraction methodologies. As governments and organizations impose stricter data protection regulations such as GDPR, CCPA, and AI transparency laws, future work should focus on developing AI-driven compliance tools that ensure scraping remains within legal and ethical boundaries.

One direction for future research is the development of privacy-preserving scraping techniques that leverage differential privacy and anonymization methods to protect both website owners and users. By ensuring that extracted data is sanitized, de-identified, and aggregated in a compliant manner, organizations can minimize ethical risks while maintaining access to valuable public data.

Additionally, AI-generated dataset attribution mechanisms should be explored, ensuring that any data scraped and used in AI training or analytics is properly cited and referenced. The introduction of automated transparency reports, which log data provenance and usage, could help address concerns about AI model biases and dataset integrity.

Beyond compliance, future research should also investigate the ethical implications of AI-generated scraping decisions, particularly in cases where models autonomously determine which data to extract, retain, or discard. Establishing ethical guidelines and AI governance policies for web scraping will be crucial to ensuring fair, responsible, and transparent data extraction practices.

6.6.5 Extracted Data Ingestion in Data Warehouses for Better Processing and Analytics

As web scraping becomes more sophisticated, integrating extracted data seamlessly into modern data warehouses is a crucial next step. Current scraping solutions often focus on extraction but lack optimized pathways for long-term storage, querying, and large-scale analysis. Future research should explore how scraped data can be efficiently structured, transformed, and ingested into cloud-based data lakes and analytics platforms.

A key direction is the development of automated ETL (Extract, Transform, Load) pipelines that allow scraped data to be cleaned, normalized, and integrated into structured warehouse environments like BigQuery, Snowflake, or Apache Druid. This would enable organizations to leverage scraped data in real-time analytics, business intelligence, and AI model training, significantly enhancing its value beyond raw collection.

Further work should focus on real-time ingestion pipelines, allowing scraped data to continuously update dashboards, predictive models, and decision-making systems. By creating a direct bridge between web scraping and enterprise analytics, organizations can harness the full potential of structured web data for strategic insights.

6.6.6 Conclusion

Future work in web scraping must focus on enhancing adaptability, accessibility, decentralization, ethical compliance, and seamless data integration. As AI-driven automation continues to evolve, these directions will shape the next generation of intelligent, responsible, and scalable web scraping solutions.

Bibliography

- Statista. Global number of internet users 2005-2024. https://www.statista.com/ statistics/273018/number-of-internet-users-worldwide/. Access Date: 05-02-2025.
- [2] Built With. https://builtwith.com. Access Date: 20-12-2024.
- [3] Kasereka Henrys. Importance of Web Scraping in E-Commerce and E-Marketing. Bugema University, January 19, 2021.
- [4] SCM De S Sirisuriya. Importance of Web Scraping as a Data Source for Machine Learning Algorithms - Review. 2023 IEEE 17th International Conference on Industrial and Information Systems (ICIIS), σελίδες 134–139, 2023.
- [5] Vijay Panwar. Web Evolution to Revolution: Navigating the Future of Web Application Development. International Journal of Computer Trends and Technology, 72:34–40, 2024.
- [6] Radu Bucea-Manea-Tonis. Angular JS The Newest Technology in Creating Web Applications. Annals of Spiru Haret University Economic Series, 16:103, 2016.
- [7] madewithangular Lalith Polepeddi. Made With Angular. https://www. madewithangular.com/sites.
- [8] Arshad Javeed. Performance Optimization Techniques for ReactJS. 2019 IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), σελίδες 1–5, 2019.
- [9] Joe Marini. Document Object Model. McGraw-Hill, Inc., USA, 1η έχδοση, 2002.
- [10] Vue JS. https://vuejs.org/. Access Date: 22-11-2024.
- [11] Next JS. https://nextjs.org/. Access Date: 08-12-2024.
- [12] Nursel Yalçın xaı Utku Köse. What is search engine optimization: SEO? Procedia
 Social and Behavioral Sciences, 9:487–493, 2010. World Conference on Learning, Teaching and Administration Papers.
- [13] Svelte. https://svelte.dev/. Access Date: 27-01-2025.
- [14] Roy Thomas Fielding. Architectural Styles and the Design of Network-based Software Architectures. University of California, Irvine, 2000. Access Date: 08-11-2024.

- [15] GraphQL. https://graphql.org/. Access Date: 31-01-2025.
- [16] gRPC. https://grpc.io/. Access Date: 16-02-2025.
- [17] Harold Davis. Search Engine Optimization. O'Reilly Media, Inc., May 2006.
- [18] Cheerio.js. https://github.com/cheeriojs/cheerio. Access Date: 21-11-2024.
- [19] BeautifulSoup. https://www.crummy.com/software/BeautifulSoup/bs4/doc/. Access Date: 13-12-2024.
- [20] *jQuery*. https://jquery.com/. Access Date: 06-01-2025.
- [21] Axios. https://axios-http.com. Access Date: 09-02-2025.
- [22] Puppeteer. https://github.com/puppeteer/puppeteer. Access Date: 04-12-2024.
- [23] Playwright. https://github.com/microsoft/playwright. Access Date: 18-01-2025.
- [24] Selenium. https://github.com/seleniumhq/selenium. Access Date: 07-02-2025.
- [25] MitmProxy. https://github.com/mitmproxy/mitmproxy. Access Date: 25-12-2024.
- [26] Gunes Acar, Marc Juarez, Nick Nikiforakis, Claudia Diaz, Seda Gürses, Frank Piessens και Bart Preneel. FPDetective: dusting the web for fingerprinters. Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security, CCS '13, σελίδα 1129–1140, New York, NY, USA, 2013. Association for Computing Machinery.
- [27] Pierre Laperdrix, Nataliia Bielova, Benoit Baudry xa Gildas Avoine. Browser Fingerprinting: A Survey. ACM Trans. Web, 14(2), 2020.
- [28] Peter Eckersley. How Unique Is Your Web Browser? Privacy Enhancing TechnologiesMikhail J. Atallah και Nicholas J. Hopper, επιμελητές, σελίδες 1–18, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
- [29] Yundi He, Runhua Shi xai Boyan Wang. WT-CFormer: High-Performance Web Traffic Anomaly Detection Based on Spatiotemporal Analysis, 2024.
- [30] Pedro Marques, Zayani Dabbabi, Miruna Mihaela Mironescu, Olivier Thonnard, Frances Buontempo, Ilir Gashi και Alysson Bessani. Using diverse detectors for detecting malicious web scraping activity. 2018 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W), σελίδες 67– 68. IEEE, 2018.
- [31] Viktor Petrov. Network Traffic Analysis: Studying Anomaly Detection Approaches for Network Traffic Analysis to Identify Suspicious Patterns and Behaviors Indicative of Cyber Threats. Cybersecurity and Network Defense Research, 3(1):13–24, 2023.

- [32] Xun Gong, Nikita Borisov, Negar Kiyavash και Nabil Schear. Website detection using remote traffic analysis. Privacy Enhancing Technologies: 12th International Symposium, PETS 2012, Vigo, Spain, July 11-13, 2012. Proceedings 12, σελίδες 58–78. Springer, 2012.
- [33] Luisvon Ahn, Manuel Blum, Nicholas J. Hopper και John Langford. CAPTCHA: Using Hard AI Problems for Security. Advances in Cryptology — EUROCRYPT 2003Eli Biham, επιμελητής, σελίδες 294–311, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg.
- [34] Meriem Guerar, Luca Verderame, Mauro Migliardi, Francesco Palmieri xat Alessio Merlo. Gotta CAPTCHA 'Em All: A Survey of 20 Years of the Human-or-computer Dilemma. ACM Comput. Surv., 54(9), 2021.
- [35] Emilio Ferrara xa Robert Baumgartner. Design of Automatically Adaptable Web Wrappers. CoRR, abs/1103.1254, 2011.
- [36] Ye Wang xai Mi Lu. An optimized system to solve text-based CAPTCHA, 2018.
- [37] Huthaifa Mohammed Kanoosh, Ammar Farooq Abbas, Noora Nazar Kamal, Zainab Mejeed Khadim, Duaa A. Majeed και Sameer Algburi. Image-Based CAPTCHA Recognition Using Deep Learning Models. Proceedings of the Cognitive Models and Artificial Intelligence Conference, AICCONF '24, σελίδα 273–278, New York, NY, USA, 2024. Association for Computing Machinery.
- [38] Vijayaragavan Pichiyan, S Muthulingam, Sathar G, Sunanda Nalajala, Akhil Ch xau Manmath Nath Das. Web Scraping using Natural Language Processing: Exploiting Unstructured Text for Data Extraction and Analysis. Proceedia Computer Science, 230:193–202, 2023. 3rd International Conference on Evolutionary Computing and Mobile Sustainable Networks (ICECMSN 2023).
- [39] spaCy. https://github.com/explosion/spaCy. Access Date: 30-11-2024.
- [40] *nltk.* https://github.com/nltk/nltk. Access Date: 20-01-2025.
- [41] Aman Ahluwalia xai Suhrud Wani. Leveraging Large Language Models for Web Scraping, 2024.
- [42] ChatGPT Search. https://openai.com/index/introducing-chatgpt-search/. Access Date: 03-02-2025.
- [43] Samuel Zuehlke, Joel Nitu, Simone Sandler, Oliver Krauss και Andreas Stöckl. Self-Repairing Data Scraping for Websites. 2024 4th International Conference on Electrical, Computer, Communications and Mechatronics Engineering (ICECCME), σελίδες 1-4, 2024.
- [44] LangChain. https://github.com/langchain-ai/langchain. Access Date: 15-02-2025.

- [45] Gelei Deng, Haoran Ou, Yi Liu, Jie Zhang, Tianwei Zhang xat Yang Liu. Oedipus: LLM-enchanced Reasoning CAPTCHA Solver, 2024.
- [46] Xuanhe Zhou, Xinyang Zhao xai Guoliang Li. LLM-Enhanced Data Management, 2024.
- [47] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan xai Dong Yu. WebVoyager: Building an End-to-End Web Agent with Large Multimodal Models. Annual Meeting of the Association for Computational Linguistics, 2024.
- [48] MongoDB. https://github.com/mongodb/mongo. Access Date: 15-11-2024.
- [49] FastAPI. https://github.com/fastapi/fastapi. Access Date: 07-12-2024.
- [50] UncleCode. Crawl4AI: Open-source LLM Friendly Web Crawler Scraper. https: //github.com/unclecode/crawl4ai, 2024.
- [51] APScheduler. https://github.com/agronholm/apscheduler. Access Date: 05-01-2025.
- [52] React. https://github.com/facebook/react. Access Date: 12-12-2024.
- [53] Refine.dev. https://github.com/refinedev/refine. Access Date: 19-11-2024.
- [54] ShadCN UI. https://github.com/shadcn-ui/ui. Access Date: 23-12-2024.
- [55] Docker. https://docs.docker.com/. Access Date: 10-12-2024.
- [56] Docker Compose. https://github.com/docker/compose. Access Date: 26-12-2024.

List of Abbreviations

SEO	Search Engine Optimization
API	Application Programming Interface
UI	User Interface
AI	Artificial Intelligence
LLM	Large Language Model
NLP	Natural Language Processing
NLU	Natural Language Understanding
NER	Natural Entity Recognition
SSR	Server-Side Rendering
SPA	Single-Page Application
CSR	Client-Side Rendering
P2P	Peer-to-Peer
ETL	Extract, Transform, Load
OCR	Optical Character Recognition
CDN	Content Delivery Network
BPF	Band Pass Filter
CAPTCHA	Completely Automated Public Turing test to tell Computers and Humans Apart
JWT	JSON Web Token
REST	Representational State Transfer
SQL	Structured Query Language
JSON	JavaScript Object Notation
RDF	Resource Description Framework
gRPC	gRPC Remote Procedure Call
SSG	Static Site Generation
ISR	Incremental Static Regeneration
DOM	Document Object Model
PWA	Progressive Web Application
ACL	Access Control List
TLS	Transport Layer Security
IoT	Internet of Things
CCPA	California Consumer Privacy Act
GDPR	General Data Protection Regulation
ACL	Access Control List
RBAC	Role-Based Access Control
SSO	Single Sign-On

VPN	Virtual Private Network
MITM	Man-in-the-Middle
SQL	Structured Query Language