

Google Scholar zoekresultaten voor wetenschappelijke projecten: linked data & natural language processing

Dhr. Bart De Paepe

Scriptie voorgedragen tot het bekomen van de graad van
Professionele bachelor in de toegepaste informatica

Promotor: Dhr. Jan Claes

Co-promotor: Dhr. Milan Lamote

Academiejaar: 2024–2025

Eerste examenperiode

Departement IT en Digitale Innovatie .

**HO
GENT**

Woord vooraf

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Samenvatting

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Inhoudsopgave

Lijst van figuren	viii
Lijst van tabellen	ix
Lijst van codefragmenten	x
1 Inleiding	1
1.1 Probleemstelling	1
1.2 Onderzoeksvraag	2
1.3 Onderzoeksdoelstelling	2
1.4 Opzet van deze bachelorproef	3
2 Stand van zaken	4
3 Methodologie	7
4 Google Scholar alert	10
4.1 Instellen van de zoekopdracht & aanmaken van het alert	10
5 Web scraping	12
5.1 Inleiding	12
5.2 Web scraping met gebruik van een LLM	14
5.2.1 Web scraping met OpenAI	14
5.2.2 Web scraping met Mirascope en Anthropic	17
5.2.3 Web scraping met een lokaal model	19
5.3 Web scraping via het parsen van de DOM	22
5.4 Meest geschikte Web scraping	24
6 Natural Language Processing	25
6.1 Inleiding	25
6.2 Relevantiescore	25
6.3 Tekstverwerking	27
7 Linked data	29
7.1 Inleiding	29
7.2 De DOI opzoeken in de URL	29
7.3 De DOI opzoeken in Crossref	30
7.4 De DOI opzoeken op de webpagina van de publicatie	30

8	Semantic search	35
8.1	Inleiding	35
8.2	Chroma	35
8.3	MongoDB	38
9	Conclusie	42
10	Google Scholar zoekopdracht	44
10.1	Zoeken in Google Scholar ("Google Scholar Guide", 2025).	44
10.1.1	Basis zoeken	44
10.1.2	Geavanceerd zoeken	46
10.2	E-mail alerts	47
10.3	Opstellen van een zoekopdracht	49
11	Parsen van Google Scholar zoekresultaten	51
11.1	Doelstelling	51
11.1.1	E-mail lezen	51
11.1.2	Body van de e-mail parsen	51
11.1.3	DOI opzoeken in de link van het geparseerde resultaat	52
11.2	Analyse	53
11.2.1	Use Cases	53
11.2.2	Domeinmodel	53
11.3	Implementatie	57
11.3.1	E-mail lezen	58
11.3.2	Body van de e-mail parsen	59
11.3.3	DOI opzoeken in de link van het geparseerde resultaat	60
11.4	Resultaat.	63
11.4.1	E-mail lezen	63
11.4.2	Body van de e-mail parsen	64
11.4.3	DOI opzoeken in de link van het geparseerde resultaat	66
12	Crossref	68
12.1	Doelstelling	68
12.2	Analyse	68
12.3	Implementatie	68
12.4	Resultaat.	68
13	Natural language processing	72
13.1	Doelstelling	72
13.2	Analyse	74
13.3	Implementatie	74
13.4	Resultaat.	76

A	Onderzoeksvoorstel	78
A.1	Inleiding	78
A.2	Literatuurstudie	80
A.3	Methodologie	81
A.4	Verwacht resultaat, conclusie	83
	Bibliografie	84

Lijst van figuren

5.1	Google Scholar SERP.	13
5.2	OpenAI dashboard.	16
5.3	HTML structuur van de GS alert.	22
6.1	Bag of Words.	26
6.2	term frequency-inverse document frequency.	26
8.1	Semantic search.	35
8.2	Chroma resultaat.	38
8.3	Mongodb resultaat.	39
10.1	Google Scholar basis zoeken.	45
10.2	Google Scholar zoekresultaten.	45
10.3	Google Scholar geavanceerd zoeken.	47
10.4	Google Scholar melding maken.	48
10.5	Google Scholar email alert.	50
11.1	E-mail lezen.	52
11.2	E-mail body parsen.	52
11.3	DOI opzoeken in de link van het geparseerde resultaat.	52
11.4	Use case diagram.	53
11.5	Domeinmodel.	53
11.6	Technology Stack.	57
11.7	Python release cycle.	58
11.8	HTML structuur van de GS alert.	60
11.9	Databank JSON view.	63
11.10	Databank table view.	63
11.11	Databank structuur e-mail.	64
11.12	Databank structuur body.	64
11.13	Databank structuur zoekresultaat.	65
11.14	Databank structuur link.	66
12.1	Databank structuur Crossref.	70
13.1	Google Scholar SERP.	73
13.2	OpenAI dashboard.	76
A.1	Chronologische oplistings van de uit te voeren stappen	79

Lijst van tabellen

11.1	GS stuurt e-mails naar google-scholar@marineinfo.org. Deze e-mails moeten automatisch verwerkt worden.	54
11.2	De body van een e-mail wordt ontleed en de informatie wordt gestructureerd opgeslaan.	55
11.3	Zoek een DOI in een tekstueel document.	56
12.1	Crossref.	69
13.1	Overzicht zoekresultaten met OpenAi en Anthropic.	77

Lijst van codefragmenten

5.1	Prompt htmlfragment	14
5.2	Cleaning codefragment	15
5.3	Prompt codefragment	15
5.4	Parse opties codefragment	16
5.5	Field codefragment	17
5.6	Schema codefragment	18
5.7	Model codefragment	18
5.8	Extract codefragment	19
5.9	Ollama	19
5.10	Scrapegraph configuratie codefragment	20
5.11	Prompt codefragment	21
5.12	Beautiful Soup select codefragment	23
5.13	Beautiful Soup navigatie codefragment	23
6.1	term frequency-inverse document frequency	26
6.2	relevantiescore	27
6.3	Coreference resolver	28
7.1	Originele URL van de publicatie	29
7.2	Opzoeken van DOI in de URL van de publicatie	31
7.3	Query op basis van de bibliografie	31
7.4	Crossref commons codefragment	32
7.5	Pymupdf codefragment	33
7.6	Selenium codefragment	34
8.1	Chroma codefragment	36
8.2	Embeddings codefragment	37
8.3	Query codefragment	37
8.4	MongoDB model codefragment	38
8.5	Embeddings codefragment	39
8.6	Index codefragment	40
8.7	Query codefragment	41
11.1	Voorbeeld header	61
11.2	Embedded url	61
11.3	Location replace codefragment	61
11.4	Pymupdf codefragment	62
12.1	Crossref commons codefragment	69
13.1	Prompt codefragment	74

13.2 Parse opties codefragment	75
13.3 Anthropic overloaded codefragment	76

1

Inleiding

Het Vlaams Instituut voor de Zee (VLIZ) ("Vlaams Instituut voor de Zee", [2024](#)) is een pionier in zeekennis. Dit wetenschappelijk instituut gelegen in Oostende heeft onder andere een mandaat om een complete en geactualiseerde catalogus bij te houden van alle wetenschappelijke publicaties in de mariene sector. Al meer dan 20 jaar bouwt het Integrated Marine Information System (IMIS) aan deze catalogus die intussen beschikt over meerdere collecties van marien wetenschappelijk referentiemateriaal.

Naast wetenschappelijke literatuur zitten er ook collecties van mariene wetenschappelijke projecten (vb. het World Register of Marine Species (WoRMS) ("World Register of Marine Species", [2024](#))) in het systeem. Het is belangrijk voor het VLIZ om te weten door hoeveel publicaties er naar een project verwezen wordt. Deze maatstaf geeft een indicatie van het draagvlak van elk project binnen de wetenschappelijke gemeenschap en is één van de belangrijkste criteria tijdens projectevaluaties. Daarom is het cruciaal om continu nieuwe publicaties waarin verwezen wordt naar die projecten op te zoeken. Daarvoor wordt op heden Google Scholar gebruikt die bekend staat als de meest uitgebreide en geactualiseerde index. Op die manier blijft IMIS up-to-date en zijn de projectreferenties steeds geactualiseerd.

1.1. Probleemstelling

Momenteel verloopt dit proces binnen het VLIZ grotendeels handmatig. De zoekresultaten, volgens een bepaalde zoekfilter per project, afkomstig van Google Scholar worden manueel gefilterd en de Digital Object Identifier (DOI) van de geselecteerde artikels wordt opgezocht. Vervolgens worden de basisgegevens van elk artikel zoals titel, auteurs, datum en uitgever opgevraagd op basis van de DOI in Crossref ("Crossref", [2024](#)). Met deze informatie wordt manueel beslist om het artikel al dan niet toe te voegen aan een collectie binnen IMIS.

Dit is een tijdrovend proces. Daarom is er vraag naar automatisatie die de zoekresultaten verwerkt en gestructureerd opslaat. Het systeem moet ook kunnen aangeven of de publicatie reeds in IMIS zit zodat er geen duplicaten opgeslaan worden. De beslissing om een artikel toe te voegen aan IMIS blijft nog altijd een manuele stap, maar naar verwachting moet dit sneller, efficiënter en accurater verlopen. Dat moet mogelijk zijn door de heterogene zoekresultaten om te zetten in gestructureerde informatie en duplicaten op te zoeken.

1.2. Onderzoeksvraag

De probleemstelling vraagt om gedeeltelijke automatisatie. Er wordt aangenomen dat er nog steeds een manuele beslissingsstap zal zijn. Toch streeft dit onderzoek naar zoveel mogelijk automatisatie waarbij het grootste deel van de publicaties automatisch aan IMIS toegevoegd worden en waarbij slecht voor een zo klein mogelijke rest van publicaties een manuele beslissing nodig is.

Bijgevolg luidt de centrale vraag die onderzocht moet worden: “Hoe kunnen de zoekresultaten van Google Scholar automatisch toegevoegd worden aan IMIS?” Dit omvat voornamelijk 2 problemen:

- Hoe kunnen de uiteenlopende zoekresultaten van Google Scholar omgezet worden in gestructureerde data?
- Zijn alle zoekresultaten uniek identificeerbaar zodat er geen duplicaten opgeslaan worden?

De uitgewerkte oplossing zal daarom met volgende aspecten rekening moeten houden:

- Hoe kunnen ook steeds nieuwe zoekresultaten van dezelfde zoekopdracht systematisch opgezocht worden?
- Kan er een score berekend worden hoe relevant elk zoekresultaat is ten opzichte van de zoekopdracht?
- Als er geen unieke identifier is, hoe wordt er dan gecontroleerd op duplicaten?

1.3. Onderzoeksdoelstelling

Het beoogde resultaat van het onderzoek is om het toevoegen van publicaties aan IMIS zoveel mogelijk te automatiseren:

- Een proof-of-concept van de meest geschikte methode om de Google Scholar zoekresultaten om te zetten in gestructureerde data.
- Een proof-of-concept van het proces dat een score voor elk zoekresultaat berekent.
- Een proof-of-concept van het opzoeken van duplicaten.

1.4. Opzet van deze bachelorproef

De rest van deze bachelorproef is als volgt opgebouwd:

In Hoofdstuk 2 wordt een overzicht gegeven van de stand van zaken binnen het onderzoeksdomein, op basis van een literatuurstudie.

In Hoofdstuk 3 wordt de methodologie toegelicht en worden de gebruikte onderzoekstechnieken besproken om een antwoord te kunnen formuleren op de onderzoeksvragen.

Dit wordt verder uitgewerkt in Hoofdstuk 4 voor het systematisch ontvangen van de meest recente publicaties. Hoofdstuk 5 gaat in op het omzetten van de Google Scholar zoekresultaten in gestructureerde data. Hoofdstuk ?? gaat over het berekenen van een score voor elk zoekresultaat. Hoofdstuk 7 onderzoekt hoe voor elk zoekresultaat een unieke identificatie kan gezocht worden. Hoofdstuk 8 bekijkt hoe voor elke publicatie de aanwezigheid van duplicaten opgespoord wordt.

In Hoofdstuk 9 tenslotte, worden alle aspecten samengebracht. Er wordt een besluit geformuleerd in welke mate de onderzoeksvragen beantwoord zijn door het onderzoek. Daarbij wordt ook vermeld wat er verder nog moet gebeuren in de toekomst om het resultaat te verbeteren.

2

Stand van zaken

Wereldwijde onuitputtelijke digitale indexen (waaronder Google Scholar), diep maatschappelijk verankerde sociale media, het Internet Of Things, de toenemende belangstelling in AI, zorgen tezamen voor een gigantische hoeveelheid aan online data. Met deze grote berg van informatie die beschikbaar is op het world wide web, is het bijzonder relevant voor organisaties om te begrijpen hoe ze deze data kunnen ontginnen teneinde er bruikbare informatie uit te filteren (Lotfi e.a., [2021](#)). Een belangrijk begrip dat daarbij prominent op de voorgrond treedt, is 'web scraping' of 'web crawling'. «Web scraping is een geautomatiseerd data extractie proces van websites met gebruik van gespecialiseerde software» (Bhatt e.a., [2023](#)). Maar web scraping is lang niet de enige techniek die gebruikt wordt om data van een webpagina te halen (Gray e.a., [2012](#)). Andere technieken zijn het gebruik van API's, screen readers en het ontleden van online PDF documenten. Het voordeel van web scraping is dat de website niet over een API met ruwe data hoeft te beschikken. De HTML code waarmee de website opgebouwd is, wordt gebruikt door de scraper om de eigenlijke inhoud eruit te filteren.

Web scrapers kunnen gebaseerd zijn op verschillende technologieën zoals 'spidering' en 'pattern matching'. Daarnaast biedt de programmeertaal waarin ze geïmplementeerd zijn ook telkens andere mogelijkheden (Bhatt e.a., [2023](#)). Lotfi e.a. ([2021](#)) onderscheidt web scrapers zowel op basis van hun toepassing (vb. medisch, social media, financieel, marketing, onderzoek) als op basis van hun methodiek (vb. copy & paste, HTML parsing, DOM parsing, HTML DOM, reguliere expressies, XPath, vertical aggregation platform, semantic annotation recognizing, computer vision web page analyzer). Hoewel de methodologie kan wijzigen, blijft het uitgangspunt van integriteit, correctheid en betrouwbaarheid van de data steeds van primordiaal belang (Lotfi e.a., [2021](#)). Verder duidt Lotfi e.a. ([2021](#)) ook nog op het iteratieve karakter van een web scraper. Het programma wordt gevoed met een url, en zal dan op zijn beurt nieuwe urls zoeken op de pagina zodat die ook bezocht kunnen worden.

Singrodia e.a. (2019) vult bovenstaande kenmerken van web scraping nog verder aan met de systematische omzetting van ongestructureerde gegevens op webpagina's naar gestructureerde data in een databank. De gegenereerde data heeft aanleiding tot filtering of tot statistieken. Het biedt vele voordelen aangezien de data opgeschoond is en daardoor als vrij van fouten kan beschouwd worden. Andere voordelen zijn tijdswinst en centrale opslag wat verdere verwerking ten goede komt.

Daarenboven vergroot web scraping de snelheid en het volume van de data die verwerkt kan worden aanzienlijk, en vermindert het ook het aantal fouten die zouden optreden door menselijke verwerking (Bhatt e.a., 2023). Uiteindelijk zet web scraping de online informatie om in business intelligence afhankelijk van het uitgangspunt van de eindgebruiker.

Na web scrapers besproken te hebben, is het niet onbelangrijk even stil te staan bij de vraag of web scrapen wel legaal is? Volgens Castrillo-Fernández (2015) is daar geen straightforward antwoord op, de situatie is afhankelijk van het betrokken land. Maar over het algemeen zijn de meest reguleringen wel in het voordeel van web scraping. Men moet vooral opletten met het respecteren van het eigendomsrecht wanneer inhoud verworven wordt door scraping en vervolgens verder gebruikt wordt.

Tot slot van de stand van zaken omtrent web scraping nog een pleidooi voor Python. Al is het perfect mogelijk om web scraping te ontwikkelen in Php, Java, en andere, toch biedt Python de meest gebruiksvriendelijke tools aan volgens Kumar en Roy (2023). «Beautiful Soup is een van de eenvoudigste bibliotheken om data te scrapen van websites. Een simpele `find_all()` in Beautiful Soup, is krachtig genoeg om de data uit het gehele document te doorzoeken. Daarna is de taak om de data te structureren. Dat kan in Python aan de hand van Pandas die in staat is om de data in een geordend formaat te presenteren.»

Om verder gebruik te maken van deze data moet deze eerst gestructureerd opgeslaan worden. Mitchell (2015) reikt een aantal methodes aan waarop dit gedaan kan worden, steeds afhankelijk van het beoogde resultaat. De skills om al die data te beheren en ermee te interageren is zo mogelijk nog belangrijker dan het scrapen op zich. De data waarvan sprake is niet gekenmerkt door strikte relaties, maar stelt eerder een collectie semi gestructureerde data voor. Een document store database lijkt in dit geval het meest aangewezen. Lourenço e.a. (2015) vergelijkt de gangbare NoSQL databases en stelt MongoDB voor als een consistente ¹ document store database.

De website in de spotlight voor dit onderzoek is Google Scholar (GS), de grootste bron van wetenschappelijke publicaties op heden. De beta versie verscheen in 2004 en sindsdien wordt het systeem voornamelijk door academici gebruikt om een persoonlijke bibliotheek aan te leggen tezamen met statistieken omtrent cita-

¹alle clients zien steeds dezelfde data

ties en h-indexen ². Volgens Oh en Colón-Aguirre (2019) is er een symbiose tussen GS and de academische wereld waarbij academici gratis hun werk kunnen aanbieden op GS. GS op zijn beurt zorgt voor goede visibiliteit van dat werk door onder andere citaten, referenties en een oplistijng van gelijkaardige onderwerpen. In tegenstelling tot de baanbrekende vernieuwing die GS introduceerde, gaat het raadplegen van deze gegevens wel gepaard met een aantal hindernissen, in het bijzonder op grote schaal.

Web scraping is 1 van de technieken die aangewend worden om dit probleem te omzeilen. Meerdere studies ((PRATIBA e.a., 2018), (Rafsanjani, 2022), (Amin e.a., 2024), (Sulistya e.a., 2024)) gebruiken bovenstaande concepten van web scrapers en gestructureerde opslag om automatisch data te ontleden van GS aan de hand van gebruiksvriendelijke interfaces. Er zijn onderling steeds wel verschillen tussen de gebruikte tools en het gekozen formaat die beide afhankelijk zijn van het beoogde resultaat. Maar het omliggend kader is steeds hetzelfde met gebruik van web scraping en gestructureerde opslag.

Het is evident dat HTML pagina's verwerkt kunnen worden door web scrapers, maar daarnaast zijn ongeveer 70% van de feiten die op het internet gepresenteerd worden, verkregen uit PDF-documenten (Singrodia e.a., 2019). Dit verklaart de noodzaak om zowel HTML pagina's als PDF documenten te scrapen.

Yang e.a. (2017) beschrijft welke HTML elementen belangrijk zijn bij het ontleden van een GS pagina en Rahmatulloh en Gunawan (2020) toont hoe deze kunnen gemapt worden naar de custom code van de scraper.

Al deze gegevens op zich zijn waardevol, maar het is pas door ze te valideren aan de hand van Crossref dat ze ook betrouwbaar worden. Crossref is een onafhankelijke organisatie die instaat voor het beheren van de unieke links naar wetenschappelijke artikels (Hendricks e.a., 2020). Daarvoor wordt de DOI gebruikt die tezamen met alle metadata van een wetenschappelijke publicatie in Crossref bijgehouden wordt. Crossref werkt voor en door alle geregistreerde uitgevers en streeft voortdurend technologische innovatie na teneinde het publicatieproces van wetenschappelijke artikels te verbeteren.

²De h-index van een wetenschappelijk onderzoeker komt overeen met de grootste h van het aantal publicaties die minstens h keer geciteerd zijn in ander werk.

3

Methodologie

Om nieuwe publicaties die gerelateerd zijn aan wetenschappelijke projecten toe te voegen aan IMIS, moet er natuurlijk een signaal zijn wanneer dergelijke publicaties beschikbaar zijn. Indexen van academische literatuur zijn daarvoor een geschikte bron en zoals eerder beschreven is Google Scholar een interessant optie. Het is de bedoeling om enkel nieuwe resultaten te ontvangen en bijgevolg vervalt dus de optie om publicaties op te zoeken aan de hand van de Google Scholar zoekpagina. De zoekpagina houdt geen gegevens bij van vorige queries en daarom is het zeer aannemelijk dat dezelfde resultaten zullen voorkomen bij opeenvolgende zoekopdrachten. Google Scholar laat echter toe om meldingen aan te maken voor een bepaalde zoekopdracht. Die stuurt automatisch nieuwe resultaten door per e-mail onder de vorm van de Google Scholar SERP ¹. Het opstellen van een zoekopdracht en het aanmaken van een alert wordt verder uitgewerkt in Hoofdstuk 4.

Google Scholar meldingen zijn e-mails in HTML formaat. Ze bevatten inhoud die overeenkomt met de Google Scholar SERP overeenkomstig de zoekopdracht. De SERP bevat een vaste structuur. Het is een lijst met zoekresultaten die telkens dezelfde elementen bevatten:

- titel
- link naar de webpagina van de publicatie
- auteurs
- tijdschrift
- jaartal
- abstract van de publicatie of een fragment ervan

¹Search Engine Result Page

Die HTML moet omgezet worden in gestructureerde data door middel van HTML scraping technieken:

- HTML scraping door een LLM ²
 - online model
 - * OpenAI
 - * Generieke procedure onafhankelijk van het model
 - lokaal model
- HTML scraping door het parsen van de DOM
 - BeautifulSoup
 - SerpAPI

Web scraping wordt verder uitgewerkt in Hoofdstuk 5.

Elk zoekresultaat is op basis van het algoritme van Google Scholar gematched met de zoekopdracht. Maar wil dat daarom ook zeggen dat de publicatie interessant is voor IMIS? Aan de hand van Natural Language Processing (NLP) wordt een score berekend van de relevantie van het zoekresultaat voor de overeenkomstige collectie in IMIS. NLP wordt verder uitgewerkt in Hoofdstuk 6.

Om te weten of een publicatie echt nieuw is, moet ze ondubbelzinnig geïdentificeerd kunnen worden. Titel, auteurs, tijdschrift, enz. maken een publicatie echter niet uniek. Er kunnen namelijk variëte benamingen voorkomen van titels, auteurs, enz. Voor literatuur is het de DOI die een publicatie uniek maakt. Voor elk zoekresultaat wordt er op zoek gegaan naar de DOI aan de hand van een stapsgewijze procedure:

- DOI opzoeken in de link naar de webpagina van de publicatie
- DOI opzoeken in Crossref op basis van de titel
- DOI opzoeken op de webpagina van de publicatie
 - dit kan een HTML pagina zijn
 - dit kan een PDF document zijn
 - dit kan een HTML pagina zijn met een embedded PDF document

Het opzoeken van de DOI wordt verder uitgewerkt in Hoofdstuk 7.

Dan resteert de vraag of de publicatie echt nieuw is, of dat ze reeds aan IMIS toegevoegd werd? Het antwoord daarop is afhankelijk van de beschikbare informatie die tijdens de voorgaande stappen gevonden werd.

- De DOI is gevonden: er kan met 100% zekerheid opgezocht worden of de publicatie reeds in IMIS zit of niet.

²Large Language Model

- De DOI is niet gevonden: omwille van de variante benamingen kan niet met volledige zekerheid opgezocht worden of een publicatie reeds in IMIS zit of niet. Wel kan door middel van semantic search op basis van de titel de waarschijnlijkheid berekend worden dat de publicatie reeds in IMIS zit.

Semantic search wordt verder uitgewerkt in Hoofdstuk 8.

4

Google Scholar alert

4.1. Instellen van de zoekopdracht & aanmaken van het alert

Om op de hoogte te blijven van nieuwe publicaties kan een index van academische literatuur gebruikt worden. Zoals eerder beschreven is Google Scholar daarvoor een interessante bron.

Het is de bedoeling om alleen nieuwe publicaties te ontvangen die nog niet eerder opgezocht werden. Bijgevolg wordt er beter geen gebruik gemaakt van de Google Scholar zoekpagina, omdat die dezelfde publicaties kan tonen bij opeenvolgende opzoeken. Daarentegen zijn er Google Scholar alerts. Dat zijn meldingen onder de vorm van automatische e-mails waarin de nieuwe publicaties opgelijst staan. Aangezien de meldingen steeds gekoppeld zijn aan een account, houdt het systeem rekening met de zoekgeschiedenis en zijn de zoekresultaten incrementeel. Bijlage Google Scholar zoekopdracht geeft een uitvoerige uiteenzetting over het opmaken van een zoekopdracht en het instellen van een melding.

In het kader van deze opdracht werd een nieuw account **google-scholar@marineinfo.org** aangemaakt, als een gedeeld account waar meerdere gebruikers toegang tot hebben. Vervolgens werd een Google account aangemaakt met hetzelfde e-mailadres. Vervolgens werden meerdere zoekopdrachten opgemaakt telkens voor een bepaald project, en werden deze geactiveerd als meldingen.

Bijvoorbeeld om relevante publicaties over het VLIZ te vinden, worden volgende zoektermen gebruikt:

- Vlaams Instituut voor de Zee
- Vlaams Instituut van de Zee
- Flanders Marine Institute
- VLIZ

- Simon Stevin
- R/V Simon Stevin
- RV Simon Stevin
- Marine Station Ostend
- Mariene Station Oostende

Vanaf dat moment zal Google Scholar voor elke zoekopdracht dagelijks e-mails met nieuwe zoekresultaten sturen naar het e-mailadres van het account, zolang dat de melding geactiveerd is.

5

Web scraping

5.1. Inleiding

De e-mails afkomstig van Google Scholar bevatten een lijst met zoekresultaten. Het formaat van de e-mail is HTML en de opmaak van de lijst is gelijkaardig aan die van de Google Scholar resultaten pagina. Dit wordt algemeen benoemd als een SERP en is voor iedereen die vertrouwd is met het internet herkenbaar als de lijst met zoekresultaten van Google (zie figuur 5.1).

De Google Scholar SERP heeft een vaste structuur, namelijk een lijst met zoekresultaten bestaande uit:

- titel
- link naar de webpagina van de publicatie
- auteurs
- naam van het tijdschrift
- jaartal
- abstract of fragment van het abstract

Bovenstaande gegevens moeten uit de SERP gefilterd worden zodat ze opgeslaan kunnen worden voor verder gebruik in de volgende stappen.

Informatie uit HTML pagina's halen is algemeen gekend onder de naam “web scraping” of “web crawling”. Deze techniek geniet momenteel veel aandacht omdat hij de gebruiker in staat stelt om veel data te verzamelen. Die data dient dan weer als brandstof voor AI. Het punt is dat in die context meerdere technieken bestaan om dezelfde job toe doen:

- web scraping met gebruik van een LLM ¹

¹Large Language Model

Big Five personality traits and university students' academic performance: A meta-analysis

S. Chen, AOK Cheung, Z. Zeng - *Personality and Individual Differences*, 2025
Meta-analytical evidence on the linkage between university students' academic performance and personality **traits**, considering cultural and contextual factors, remains limited. This meta-analysis reports comprehensive assessment of the ...



Measuring personality Traits: Simulations Exploring (Mis) Alignment between methods and measurands

P. Durkee - *Journal of Research in Personality*, 2025
... Measuring personality **traits** requires that measurement methods align with intended measurands. There are two broad classes of personality **trait** ... The simulations highlight challenges in capturing personality **trait** measurands and offer ...



Landscape modification and species traits shape seasonal wildlife community dynamics within an arid metropolitan region

J.D. Haight, S.J. Hall, J.S. Lewis - *Landscape and Urban Planning*, 2025
... For this study, we evaluated how landscape characteristics and species **traits** influenced habitat use (i.e. patch occupancy, persistence, and ...). However, habitat use further varied according to species **traits**, with larger-bodied species exhibiting ...



[HTML] Consequences of Dietary Olive Leaf Powder Supplementation on Growth Performance, Carcass Traits, Blood Biochemical Parameters and Gut Microbiota in Broilers

M.H. Nagm, A.K. Aldhalmi, E.A. Ashour, L.A. Mohamed... - *Poultry Science*, 2025
This experiment examined the potential of olive leaf powder (OLP) as a natural growth enhancer in broiler feed and its influences on growth performance, carcass characteristics, blood parameters, and intestinal bacterial count. A total of 210 one-day-old ...



[HTML] ... -scale integration of meta-QTL and genome-wide association study identifies genomic regions and candidate genes for photosynthetic efficiency traits in bread wheat

M. Chen, T. Chen, L. Yun, Z. Che, J. Ma, B. Kong, J. Long... - *BMC Genomics*, 2025
... efficiency, a complex **trait**, through the ... in **trait** coverage, demonstrating significant advantages in data size and **trait** resolution dimensions. The increasing and intensive application of MQTL has effectively revealed the complex genetic ...



[HTML] Comprehensive evaluation of sugar beet varieties based on genotype× yield× trait (GYT) in different environments

D. Xinwang, L. yuhang, H. Xiaohang, L. yanli - *Scientific Reports*, 2025
... The aim of this study was to assess the agronomic **traits** and genotype by yield × **trait** (GYT) interactions of sugar beet across ... **trait** combinations was analyzed using the GYT biplot technique. Principal components PC1 and PC2 accounted for ...



[HTML] A Unified Method for Detecting Phylogenetic Signals in Continuous, Discrete, and Multiple Trait Combinations

L. Yan, Y. Yuan - *Ecology and Evolution*, 2025
... We used d phylo to represent phylogenetic distance and d **trait** to represent **trait** distance, and d **trait** was calculated using Gower's distance, which allows our method to handle any type of **trait**. Considering three species i, j, and k, if d [i, j] phylo < ...



... sequence repeat (SSR) markers for Chinese yam (*Dioscorea opposita* Thunb.) and genetic diversity and association analysis using phenotypic and genotypic traits

Q. Zhang, H. Tian, S. Guo, Y. Wang, S. Pei, C. Wu, Q. Wang... - *Genetic Resources and ...*, 2025
Chinese yam (*Dioscorea opposita* Thunb.) has a long history of cultivation and application, and is an important medicinal and edible plant in China. However, the lack of molecular markers limits the efficiency and accuracy of the genetic breeding ...



[HTML] Aboveground-canopy and belowground-root-trait correlations contribute


```
<h3      style=font-weight:normal;margin:0;font-size:17px;line-height:20px;><span
style=font-size:11px;font-weight:bold;color:1a0dab;vertical-align:2px>[HTML]</span>
<a href="https://scholar.google.be/scholar_url?url=https://www.sciencedirect.com/science/article/pii/
GcCSieoPqLGJyAM&amp;scisig=AFWwaebeqxEetR78Fz7JpxPtT7ui&amp;oi=scholaralrt&amp;hist=I
kY6yEQWAvLJNk68&amp;html=&amp;pos=0&amp;folt=kw-
topclass=gse_alrt_titlestyle=font-size:17px;color:1a0dab;line-height:22px>The  role
of narcissistic personality <b>traits </b>in bullying behavior in adolescence–A
systematic review and meta-analysis</a></h3>
```

Codefragment 5.1: HTML fragment van de titel van een publicatie.

- online model
- lokaal model
- web scraping met verwerking van de DOM
 - BeautifulSoup
 - SerpAPI

5.2. Web scraping met gebruik van een LLM

De meest gekende vorm van web scraping gebruikt de DOM structuur om er de inhoud uit te filteren. Daar komt dus niets van AI bij kijken. Het probleem met die aanpak is dat de custom code die de HTML structuur verwerkt, sterk afhankelijk is van de HTML zelf. Bijvoorbeeld wanneer de titel van een publicatie tussen `<h3>` tags staat die gekenmerkt worden door een `class="gse_alrt_title"` (zie codefragment 5.1), dan zal de custom code specifiek filteren op die DOM elementen. Maar wanneer de structuur van de HTML wijzigt om welke reden dan ook, dan zal de web scraper niet langer werken zolang de custom code niet werd aangepast.

Daarom begint dit hoofdstuk met meer recente technieken die gebaseerd zijn op AI en die niet strikt afhankelijk zijn van de DOM structuur.

5.2.1. Web scraping met OpenAI

LLMs zijn bijzonder goed in het beantwoorden van vragen. OpenAI is voor het brede publiek beter gekend door zijn chatbot ChatGPT. Maar hoe goed is OpenAI in het parsen van een webpagina? OpenAI biedt ook een API ("OpenAI developer platform", 2025) aan waarmee gebruikers opdrachten kunnen sturen naar een model. De opdracht/vraag in kwestie is: "Geef de titels, originele links, auteurs, naam van de tijdschriften, jaartal van de publicaties en tekstfragmenten van hetvolgende Google Scholar zoekresultaat."

Het online artikel "Web scraping experiment with AI (Parsing HTML with GPT-4 and GPT-4o)" (2025) beschrijft stap voor stap hoe de gevraagde gegevens verkregen kunnen worden met het "gpt-4-1106-preview" model van OpenAI. De code op git-

```
1      # Remove all occurrences of content between <head> and </head>
2      body_text = re.sub(r'<head.*?>.*?</head>', '', body_text,
3                          flags=re.DOTALL)
4      # Remove all occurrences of content between <script> and
5      </script>
6      body_text = re.sub(r'<script.*?>.*?</script>', '', body_text,
                          flags=re.DOTALL)
7      # Remove all occurrences of content between <style> and
8      </style>
9      body_text = re.sub(r'<style.*?>.*?</style>', '', body_text,
                          flags=re.DOTALL)
```

Codefragment 5.2: Codefragment voor het opkuisen van de HTML.

```
1      messages=[
2      {"role": "system",
3       "content": "You are a master at scraping Google Scholar
4       results data. Scrape top 10 organic results data
5       from Google Scholar search result page."},
6      {"role": "user", "content": body_text}
7      ],
```

Codefragment 5.3: Codefragment voor het opstellen van een prompt.

hub “Web scraping Google Scholar SERP using OpenAI” ([2025](#)) maakt hiervan een implementatie aangepast voor de Google Scholar alerts.

De HTML wordt eerst ontdaan van overbodige tags zodat alleen de inhoud behouden blijft zoals te zien is in codefragment 5.2. De prompt voor het model is te zien in codefragment 5.3. De verwachte parameters die het model moet zoeken zijn te zien in codefragment 5.4.

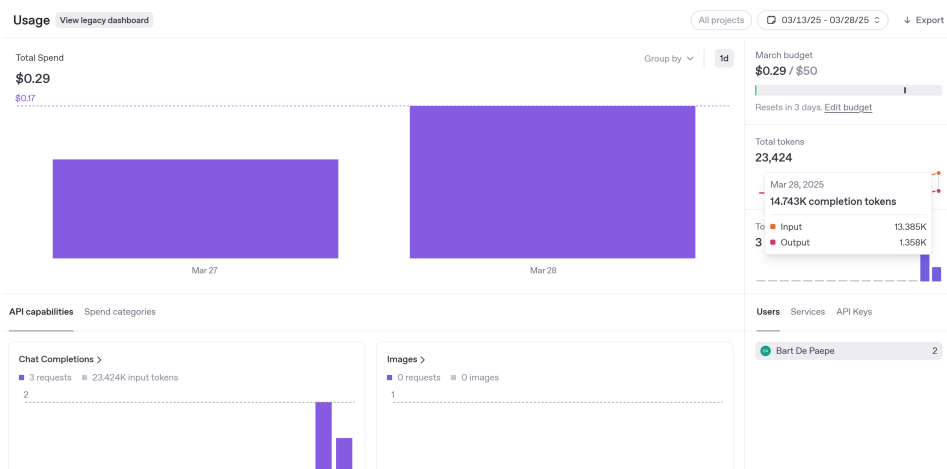
Het is vereist om een account te registreren bij OpenAI en om dit account te crediteren. Vervolgens kan er met de API gewerkt worden en wordt er betaald naargelang het verbruik. Een test met 1 Google Scholar alert met 10 zoekresultaten heeft een prijs van 0,17\$ zoals te zien is in figuur 5.2. Het model gebruikte daarvoor 14743 tokens. De test was succesvol. Het model vond de gevraagde parameters voor elk van de 10 zoekresultaten. Toch houdt het onderzoek naar web scraping hier niet op. Het resultaat tot zover is immers betalend en dat is niet noodzakelijk de beste oplossing voor het probleem. Bovendien werkt deze oplossing enkel met OpenAI modellen wat weinig flexibiliteit toelaat.

```

1      "function": {
2          "name": "parse_data",
3          "description": "Parse organic results from Google Scholar
4                          SERP raw HTML data nicely",
5          "parameters": {
6              'type': 'object',
7              'properties': {
8                  'data': {
9                      'type': 'array',
10                     'items': {
11                         'type': 'object',
12                         'properties': {
13                             'title': {'type': 'string'},
14                             'original_url': {'type': 'string'},
15                             'authors': {'type': 'string'},
16                             'year_of_publication': {'type':
17                                 'integer'},
18                             'journal_name': {'type': 'string'},
19                             'snippet': {'type': 'string'}
20                         }
21                     }
22                 }
23             }
24         }
25     }

```

Codefragment 5.4: Codefragment voor het opstellen van de zoekopties.



Figuur 5.2: OpenAI dashboard.

```
1 class FieldDefinition(BaseModel):  
2     """Define the fields to extract from the webpage."""  
3  
4     name: str = Field(..., description="The desired name for this  
5         field.")  
6     type: Literal["str", "int", "float", "bool", "list"]
```

Codefragment 5.5: Codefragment voor het opstellen van een Pydantic Field.

5.2.2. Web scraping met Mirascope en Anthropic

De vorige scraper was zeer specifiek geschreven voor OpenAI. De code zou niet werken voor een model van een andere provider. Daarom wordt er verder gezocht naar een meer generieke methode die niet afhankelijk is van de provider van het model. Mirascope “LLM abstractions that aren’t obstructions” (2025) is een gespecialiseerde code library om op uniforme wijze met verschillende LLMs te werken. Het online artikel “Extracting Structured Data from Websites with Mirascope, Anthropic, and Streamlit” (2025) beschrijft stap voor stap hoe de gevraagde gegevens verkregen kunnen worden aan de hand van Mirascope en het Anthropic Claude model. Net zoals OpenAI biedt Anthropic ook een API (“Using the API Getting started”, 2025) aan voor development. Er moet eveneens een account geregistreerd worden bij Anthropic maar hun free tier is uitgebreider dan dat van OpenAI zodat het account voor dit experiment niet gecrediteerd moest worden.

De code op github “Web scraping Google Scholar SERP using Anthropic” (2025) maakt hiervan een implementatie aangepast voor de Google Scholar alerts. De parameters die opgezocht moeten worden, dienen eerst als velden en schema gedeclareerd te worden. Dit is te zien in codefragmenten 5.5 en 5.6. Met deze structuur van veld en schema kan vervolgens voor elke opdracht een gepast model gemaakt worden aan de hand van de code in codefragment 5.7. Dan kan er geparsed worden met de resultaten volgens het gevraagde schema zoals getoond wordt in codefragment 5.8.

Ook deze test was succesvol. Het model vond de gevraagde parameters voor elk van de 10 zoekresultaten. Toch houdt het onderzoek naar web scraping ook hier niet op. De oplossing tot zover heeft nog steeds een beperkte free tier en meerdere opdrachten zouden ook aanleiding geven tot een betalende service. Daarnaast zijn tot zover beide oplossingen ook afhankelijk van een online service wat in bepaalde gevallen niet wenselijk kan zijn (bijvoorbeeld wanneer de klant een on-premise oplossing wil).

```

1      class SchemaGenerator(AnthropicExtractor[list[FieldDefinition]]):
2          """Generate a schema based on a user query."""
3
4          api_key = settings.anthropic_api_key
5
6          extract_schema: Type[list] = list[FieldDefinition]
7
8          prompt_template = """
9          Call your tool with field definitions based on this query:
10         {query}
11         """
12
13         query: str

```

Codefragment 5.6: Codefragment voor het opstellen van een Pydantic Schema.

```

1      def generate_schema(self) -> None:
2          """Sets `extract_schema` to a schema generated based on
3             `query`."""
4          field_definitions =
5              SchemaGenerator(query=self.query).extract()
6          model = create_model(
7              "ExtractedFields",
8              __doc__=DEFAULT_TOOL_DOCSTRING,
9              **{
10                  field.name.replace(" ", "_"): (field.type, ...)
11                  for field in field_definitions
12              },
13          )
14          self.extract_schema = list[model]

```

Codefragment 5.7: Codefragment voor het opstellen van een Pydantic Model.

```
1      @computed_field
2      @property
3      def webpage_content(self) -> str:
4          """Returns the text content of the webpage found at `url`."""
5          soup = BeautifulSoup(self.html_text, "html.parser",
6                               from_encoding="utf-8")
7          text = soup.get_text()
8          for link in soup.find_all("a"):
9              text += f"\n{link.get('href')}}"
10         return text
```

Codefragment 5.8: Codefragment voor het parsen.

```
docker exec -it ollama ollama pull llama3:8b docker exec -it ollama ollama run llama3
```

Codefragment 5.9: Ollama

5.2.3. Web scraping met een lokaal model

Om niet afhankelijk te zijn van een online service moet het model lokaal gehost worden. “Ollama: Get up and running with large language models.” (2025) is een platform om LLMs te hosten op je eigen systeem. Eenmaal geïnstalleerd kan een model naar keuze gedownload en uitgevoerd worden op de lokale machine. Dit kan bijvoorbeeld aan de hand van Docker zoals te zien is codefragment 5.9. Het is daarbij aan de gebruiker om te kiezen welk model het best geschikt is voor het probleem, maar ook om de vereisten van het model en de beschikbare hardware op elkaar af te stemmen. Er is een groot aanbod van beschikbare modellen, maar hieronder staan de 2 modellen die gebruikt werden voor dit onderzoek.

- LLaMA — geschikt voor een brede waaier aan toepassingen, LLaMA blinkt uit in het begrijpen en het genereren van tekst die de menselijke taal nabootst (“Introducing Meta Llama 3: The most capable openly available LLM to date”, 2025).
- Gemma — biedt een lichtere versie van een model dat ook geschikt is voor web scraping (“Gemma3: Powerful AI Within Everyone’s Reach”, 2025).

Ieder van deze modellen is beschikbaar in verschillende versies van grootte van het model. Een groot model heeft meer kennis en zal beter op de vraag kunnen antwoorden. Anderzijds zal de verwerking trager verlopen bij een groot model en zal een groot model veel meer resources nodig hebben om te kunnen werken. Daarom wordt in eerste instantie met een klein llama3:8b model (4.7GB) gewerkt dat kan uitgevoerd worden op een doorsnee laptop. Het online artikel “Scrape-graph AI Tutorial; Scrape websites easily with LLaMA AI” (2025) beschrijft stap voor

```

1      graph_config = {
2          "llm": {
3              "model": "ollama/gemma3:27b", # Specifies the large
              language model to use
4              "temperature": 0,           # Temperature controls
              randomness; 0 makes it deterministic
5              "format": "json",           # Output format is set to
              JSON
6              "base_url": "http://localhost:11434", # Base URL
              where the Ollama server is running
7              "model_tokens": 128000,
8          },
9          "embeddings": {
10             "model": "ollama/nomic-embed-text", # Specifies the
              embedding model to use
11             "temperature": 0,                 # Keeps the
              generation deterministic
12             "base_url": "http://localhost:11434", # Base URL for
              the embeddings model server
13         },
14         "verbose": True, # Enables verbose output for more
              detailed log information
15     }

```

Codefragment 5.10: Codefragment voor het configureren van Scrapegraph.

stap hoe webpagina's gescraped kunnen worden aan de hand van ScrapegraphAI en het LLaMa3 model. "ScrapeGraphAI: You Only Scrape Once" (2025) is een library die LLMs en grafen combineert om webpagina's te scrapen.

De code op github "Web scraping Google Scholar SERP using Llama3" (2025) maakt hiervan een implementatie aangepast voor de Google Scholar alerts. De code veronderstelt dat ollama draait op het systeem en dat het LLaMa3 model erop uitgevoerd wordt. Codefragment 5.10 connecteert de scraping graph met het model. De prompt voor het model is te zien in codefragment 5.11.

De test was niet succesvol. Het model vond geen enkele van de gevraagde parameters voor geen enkele van de 10 zoekresultaten. In eerste instantie werd gedacht dat het geselecteerde model te zwak is voor deze taak omdat er voor een klein model gekozen werd. Daarom werd dezelfde test uitgevoerd op een krachtiger systeem met 40GB RAM met het gemma3:27b model (17GB). Maar ook deze test was niet succesvol: de opdracht bleef gedurende 1 uur lopen zonder resultaat, waarna het commando manueel onderbroken werd.

```
1 smart_scraper_graph = SmartScraperGraph(  
2     prompt="You are an AI assistant specialized in processing  
        Google Scholar search engine result pages and returning  
        structured JSON data. Always provide your response as  
        valid, well-formatted JSON without any additional text  
        or comments. Focus on extracting and organizing the most  
        relevant information from Google Scholar search engine  
        result pages, including title, original url, authors,  
        name of journal, year of publication, snippet.", # The  
        AI prompt specifies what to extract  
3     source=body_text, # URL of the website from which to scrape  
        data  
4     config=graph_config, # Uses predefined configuration  
        settings  
5 )  
6  
7 # Execute the scraping process  
8 result = smart_scraper_graph.run()  
9
```

Codefragment 5.11: Codefragment voor het opstellen van een prompt.



Figuur 5.3: HTML structuur van de GS alert.

5.3. Web scraping via het parsen van de DOM

Tot zover is er nog geen bruikbaar resultaat. Het gebruik van een online model is betalend en kan niet gewenst zijn door de klant. Anderzijds leidt een lokaal model voorlopig niet tot een goed resultaat. Daarom wordt er teruggегреpen naar de meer klassieke benadering van web scraping waarbij de HTML structuur ontleed wordt aan de hand van custom code. De Google Scholar SERP is een gestructureerde HTML pagina. Voor het parsen wordt er gebruik gemaakt van deze structuur. Elk afzonderlijk resultaat heeft een H3-tag met als klasse 'gse_alrt_title' waarin de titel staat. De korte tekst van het zoekresultaat staat in een DIV-tag met als klasse 'gse_alrt_sni'. Tussen de titel en het snippet staan de auteurs, de uitgever en het jaartal gegroepeerd in een DIV-tag zoals te zien is in figuur 5.3.

Meerdere specifieke libraries laten toe om HTML te parsen, maar de meest gekende bibliotheek is BeautifulSoup ("Beautiful Soup", (2025)). Door de naam van de HTML-tag en de gebruikte klasse (gse_alrt_title, en gse_alrt_sni) mee te geven aan BeautifulSoup, worden de overeenkomstige elementen in het HTML-document weergegeven. BeautifulSoup laat ook toe om te navigeren in het document zodat het element met de auteurs, de uitgever en het jaartal gemakkelijk gevonden kan worden.

De code op github "Web scraping Google Scholar SERP using BeautifulSoup" (2025) maakt hiervan een implementatie aangepast voor de Google Scholar alerts. De selectie van tags aan de hand van BeautifulSoup is te zien in codefragment 5.12 De navigatie tussen tags aan de hand van BeautifulSoup is te zien in codefragment 5.13 De volledige uitwerking is te zien in bijlage Parse Google Scholar zoekopdracht. De test met een SERP met 10 zoekresultaten was geslaagd. De gevraagde parameters worden gevonden voor elk van de 10 resultaten.

```
1      ...
2      soup = BeautifulSoup(body_text, "html.parser",
3                             from_encoding="utf-8")
4      all_titles = soup.find_all("a", {"class": "gse_alrt_title"})
5      all_snippets = soup.find_all("div", {"class":
6                                     "gse_alrt_sni"})
7
8      ...
9
10     for i in range(0, len(all_titles)):
11         title = all_titles[i].get_text()
12         snippet = all_snippets[i].get_text()
13         try:
14             data = self.parse_search_result(email_id, all_titles[i],
15                                             all_snippets[i])
16         ...
```

Codefragment 5.12: Codefragment voor het selecteren van tags aan de hand van BeautifulSoup.

```
1      ...
2      media_type = title.find_previous()
3      if media_type.name.lower() == "span":
4
5      ...
6
7      author_publisher_year = snippet.find_previous()
8      author_publisher_year = author_publisher_year.get_text()
9      ...
```

Codefragment 5.13: Codefragment voor het navigeren tussen tags aan de hand van BeautifulSoup.

5.4. Meest geschikte Web scraping

Tenslotte moet ook nog SerpAPI (“Scrape Google and other search engines from our fast, easy, and complete API.”, 2025) vermeld worden. Tijdens het onderzoek naar web scraping kwam dit platform meermaals naar voren, zelfs voor de toepassingen die gebruik maken van AI. SerpAPI biedt specifieke oplossingen voor elke mogelijke search engine, waaronder ook Google Scholar. Alleen focust hun product telkens op de online search engine en niet op alerts waardoor de zoekresultaten niet incrementeel zijn.

Tot besluit van dit hoofdstuk wordt er dus verder gewerkt met het parsen van de zoekresultaten aan de hand van BeautifulSoup. De grote tekortkoming dat de code afhankelijk is van de HTML wordt verzacht doordat die HTML ook weerspiegeld wordt in de code. Aanpassingen zouden dus begrijpbaar en beperkt moeten blijven.

Op een dag gaat de code dus breken en het zou goed zijn mocht men hiervan op de hoogte gebracht worden. Daarvoor kan er een e-mail notificatie geïmplementeerd worden die een bericht stuurt wanneer BeautifulSoup de gewenste tags in het HTML document niet langer vindt.

6

Natural Language Processing

6.1. Inleiding

Google Scholar heeft feitelijk al gezocht naar relevante publicaties voor onze zoekopdracht. De sortering van de zoekresultaten is op basis van de omvang van de integrale tekst, het aanzien van het tijdschrift, het aanzien van de auteurs en het aantal citaties (inclusief de leeftijd ervan). Op die manier volgt Google Scholar de manier van werken binnen de academische wereld.

Maar dit laat nog steeds ruimte om voor elk zoekresultaat een score af te leiden die aangeeft hoe relevant de publicatie is voor IMIS. De zoekopdracht voor het VLIZ bijvoorbeeld gebruikt als 1 van de trefwoorden “Simon Stevin”, zijnde het wetenschappelijk vaartuig van het VLIZ. Dat betekent dat Google Scholar ook zoekresultaten zal geven van publicaties die verwijzen naar “Simon Stevin”, de vlaamse geleerde uit de 16de eeuw. Het is interessant om die artikels te kunnen markeren met een lagere score ten opzichte van publicaties die het trefwoord “VLIZ” bevatten.

Natural Language Processing is een verzamelnaam van een hele groep technieken die tekst omzetten in informatie. De toepassingen en technieken zijn zeer uitgebreid en daarom ook onderverdeeld in verschillende deelgebieden. Het bepalen van de relevantie van een tekst in functie van een trefwoord valt eerder onder het deelgebied van de “Natural Language Understanding”.

6.2. Relevantiescore

Het uitgangspunt is om een score te berekenen door te tellen hoe vaak een trefwoord voorkomt in een tekst. De frequentie is dan recht evenredig met de relevantie van de tekst. Een eenvoudige techniek voor het bepalen van de frequentie is aan de hand van een “Bag of Words” (BoW). Dat is een tabel met een rij voor elke tekst en een kolom voor elk uniek woord. De cellen tonen het aantal keer dat het

	I	like	this	article	medium	data
I like this article	1	1	1	1	0	0
I like medium	1	1	0	0	1	0
I like data	1	1	0	0	0	1

Figuur 6.1: Bag of Words.

$$score = tf * idf \quad (6.1)$$

where

$$tf = \text{term frequency (see above)} \quad (6.2)$$

$$idf_t = \log \left(\frac{N}{df_t} \right) \quad (6.3)$$

$$N = \text{total number of documents} \quad (6.4)$$

$$df_t = \text{the number of documents in which term } t \text{ occurs} \quad (6.5)$$

Codefragment 6.1: term frequency-inverse document frequency

woord voorkomt in de tekst. Dit is te zien in figuur 6.2. Toch is de frequentie op zich nog niet de beste waardemeter voor een teksts. Er kunnen namelijk woorden heel vaak voorkomen die op zich weinig vertellen over het onderwerp. Om daaraan tegemoet te komen bestaat er een variant van de BoW, die in plaats van gewoon te tellen de “term frequency-inverse document frequency” (Tf-Idf) geeft. De formules zijn te zien in codefragment 6.1 en het resultaat in figuur ???. Maar de relevantie van een publicatie kan niet afhangen van het aantal zoekresultaten, daarom valt de “inverse document frequency” weg. In de plaats zijn er andere parameters om de relevantie te benaderen.

- De topic-sentence ratio: Een ratio van het trefwoord count ten opzichte van het totaal aantal zinnen.
- De topic-noun ratio: Een ratio van het trefwoord count ten opzichte van het totaal aantal zelfstandige naamwoorden.

Aan de hand van die parameters kan een relevantiescore berekend worden zoals

	I	like	this	article	medium	data
I like this article	0,35	0,35	0,35	0,35	0	0
I like medium	0,35	0,35	0,35	0	0,35	0
I like data	0,35	0,35	0,35	0	0	0,35

Figuur 6.2: term frequency-inverse document frequency.

$$\text{score} = \log(\text{topiccount} * \text{topic} - \text{sentenceratio} * \text{topic} - \text{nounratio}) \quad (6.6)$$

Codefragment 6.2: relevantiescore

te zien in codefragment 6.2. De code op github “Natural Language Processing of Google Scholar search results” (2025) maakt hiervan een implementatie aangepast voor de Google Scholar alerts.

6.3. Tekstverwerking

De BoW geeft wel aanleiding voor tabellen met zeer grote dimensies. Hoe meer teksten er zijn, des te uitgebreider zal de woordenschat worden. Daarom wordt de BoW altijd voorafgegaan door tekstverwerking die stopwoorden¹ verwijdert en lemmatisering².

Voornaamwoorden kunnen verwijzen naar het trefwoord. Het is dus relevant om ze te vervangen door het trefwoord zelf zodat het van invloed is op de score. “Coreference resolution” is daar de geschikte NLP techniek voor. In een eerste fase worden alle gerelateerde voornaamwoorden en zelfstandige naamwoorden opgezocht. In een tweede fase worden alle gevonden voornaamwoorden vervangen door hun bijhorende zelfstandig naamwoord. Daarvoor gebruikt deze techniek een coreference model dat bestaat uit opeenvolgende paren van tokens met een bijhorende coreference score. Het opzoeken van paren is te zien in codefragment 6.3.

¹Stopwoorden (of, a, the, in, you, ...) komen vaak voor maar hebben geen toegevoegde waarde over het onderwerp.

²Lemmatisering zet woorden om naar hun basisvorm, maar houdt daarbij rekening met de context. Voorbeeld “caring” wordt omgezet in “care” en niet in “car”.

```
1      import spacy
2      ...
3      def __init__(self, db_service: DBService,
4                    logging_service: LoggingService):
5      self.db_service = db_service
6      self.logging_service = logging_service
7      self.nlp = spacy.load("en_coreference_web_trf")
8      ...
9      def coreference_resolution(self, text):
10     doc = self.nlp(text)
11     spans = doc.spans
12     span_array = []
13     self.logging_service.logger.debug(spans)
14     for spangroup in spans.values():
15     span_tuple = []
16     for span in spangroup:
17
18         self.logging_service.logger.debug(text[span.start_char:span.end_char])
19     span_tuple.append(text[span.start_char:span.end_char])
20     span_array.append(span_tuple)
21     return span_array
```

Codefragment 6.3: Coreference resolver

7

Linked data

7.1. Inleiding

De web scraping leverde gegevens zoals titel, auteurs, naam van tijdschrift en jaartal van de publicatie op. Deze informatie heeft echter geen garantie van uniciteit. Dat is te wijten aan mogelijke varianten (zoals bijvoorbeeld afkortingen) in de titel, auteurs en naam van tijdschriften.

Nochtans moet een publicatie uniek geïdentificeerd kunnen worden om uit te sluiten dat er reeds duplicaten in IMIS zitten. De DOI vormt de unieke indentifier van elke publicatie. Er zal getracht worden om deze voor elke publicatie op te zoeken. Hiervoor worden 3 afzonderlijke bronnen geraadpleegd in chronologische volgorde en zolang de DOI niet gevonden werd:

- In de URL van de link naar de originele publicatie
- Aan de hand van een opzoeking in Crossref op basis van de titel
- Op de originele webpagina van de publicatie

7.2. De DOI opzoeken in de URL

De web scraping vond ook de URL naar de webpagina van de oorspronkelijke publicatie. Empirisch valt het op dat veel van die URLs opgebouwd zijn met de DOI van de publicatie (voorbeeld in codefragment 7.1).

Alle DOIs hebben dezelfde structuur: ze beginnen met het cijfer 10, gevolgd door een punt en 4 tot 9 cijfers, daarna volgt een slash. Verder kan elke willekeurige opvolging van letters, cijfers, speciale tekens en slashes voorkomen. De volgende

<https://pubs.acs.org/doi/full/10.1021/acsami.4c21991>

Codefragment 7.1: Originele URL van de publicatie

lijst reguliere expressies is aanbevolen om DOIs op te zoeken “DOIs and matching regular expressions” (2025):

- `/^10.\d4,9/[-._;():A-Z0-9]+$ /i`
- `/^10.1002/[^s]+$ /i`
- `/^10.\d4/\d+-\d+X?(\d+)\d+<[\d\w]+:[\d\w]*>\d+.\d+.\w+;\d$/i`
- `/^10.1021/\w\w\d++$/i`
- `/^10.1207/[\w\d]+\&\d+_d+$/i`

Indien 1 van deze reguliere expressies matcht met de URL, dan is de DOI gevonden. De omzetting in code is te zien in codefragment 7.2

7.3. De DOI opzoeken in Crossref

Indien de voorgaande stap geen DOI opleverde, dan wordt hier op basis van de titel een opzoeking van de DOI gedaan in Crossref. Daarvoor beschikt Crossref over een API (“Tips for using the Crossref REST API”, 2025). De documentatie leert dat ondoordachte requests kunnen leiden tot zeer langdurige queries en/of ongewenste resultaten. Concreet wordt er afgeraden om meer dan 2 velden op te nemen in de query van een sample, of om meer dan 2 resultaten te vragen in het geval van een matching. Verder wordt voor matching aangeraden om te zoeken aan de hand van de bibliografische gegevens zoals te zien in codefragment 7.3. Maar Google Scholar alerts geven geen bibliografie, daarom wordt enkel de titel verder gebruikt. Het is niet nodig om de API rechtstreeks te bevragen aangezien meerdere onafhankelijke Python libraries daar een wrapper voor geschreven hebben:

- Crossref Commons for Python (“Crossref Commons for Python”, 2025)
- Habanero (“Habanero”, 2025)
- Crossrefapi (“Crossref API”, 2025)

Al deze bibliotheken bieden dezelfde tools en presteren gelijkaardig. Zonder bijzondere reden, behalve dat Crossref Commons ontwikkeld wordt door Crossref zelf, wordt er met Crossref Commons for Python gewerkt. De code om een sample op te vragen op basis van de titel is bijzonder compact zoals te zien in 7.4. Indien de Crossref API een resultaat geeft, dan is de DOI gevonden.

7.4. De DOI opzoeken op de webpagina van de publicatie

Indien de voorgaande stap geen DOI opleverde, dan wordt het DOI gezocht op de webpagina van de publicatie. In de meeste gevallen is dat op de website van de uitgever van het tijdschrift. Op die pagina staan altijd de titel, auteurs, naam

```

1      def search_in_text(text, link):
2          # find using regex
3          patterns = get_patterns()
4          doi_result = None
5          while (doi_result is None) and (len(patterns) > 0):
6              doi_result = re.search(patterns.pop(), text, re.IGNORECASE)
7
8          if doi_result is not None:
9              # update DOI
10             doi = doi_result.group(0)
11             link.doi = doi
12             link.is_doi_success = True
13             link.log_message = "DOI successfully retrieved"
14
15             ...
16
17         def get_patterns():
18             # https://www.crossref.org/blog/does-and-matching-regular-
19               expressions/
20             # /^10.\d{4,9}/[-._;()/:A-Z0-9]+$ /i
21             # /^10.1002/[^s]+$ /i
22             # /^10.\d{4}/\d+-
23               \d+X?(\d+)\d+<[\d\w]+:[\d\w]*>\d+.\d+.\w+;\d$/i
24             # /^10.1021/\w\w\d++$/i
25             # /^10.1207/[\w\d]+\&\d+_\d+$/i
26
27             patterns = [r"10.1207/[\w\d]+\&\d+_\d+", r"10.1021/\w\w\d++",
28               r"10.\d{4}/\d+-\d+X?(\d+)\d+<[\d\w]+:[\d\w]*>\d+.\d+.\w+;\d",
29               r"10.1002/[^s]+",
30               r"10.\d{4,9}/[-._;()/:A-Z0-9]+"]
31             return patterns

```

Codefragment 7.2: Codefragment voor het opzoeken van DOI in de URL van de publicatie.

[http://api.crossref.org/works?query.bibliographic="Toward a Unified Theory of High-Energy Metaphysics, Josiah Carberry 2008-08-13"rows=2](http://api.crossref.org/works?query.bibliographic=)

Codefragment 7.3: Query op basis van de bibliografie

```

1         try:
2             filter = {}
3             queries = {'query.title': title}
4             response = crossref_commons.sampling.get_sample(size=2,
                    filter=filter, queries=queries)

```

Codefragment 7.4: Codefragment voor opvragen van een sample op basis van de titel aan Crossref.

van het tijdschrift, jaartal en abstract van de publicatie. Soms staat ook de DOI op die pagina. In sommige gevallen is de integrale tekst van de publicatie hier beschikbaar als HTML of als PDF document.

Het formaat van de pagina kan verschillende zijn:

- Een gewone HTML pagina
- Een PDF document
- Een webpagina met een embedded PDF document

Voor elk van de 3 gevallen is er een andere verwerking nodig:

- In het geval van een HTML pagina wordt de inhoud geparsed met BeautifulSoup net zoals dat eerder gebeurde voor de SERP. Vervolgens wordt een DOI opgezocht in de inhoud door middel van de reguliere expressies uit hoofdstuk 7.2.
- In het geval van een PDF document is er een extra tussenstap nodig. De inhoud moet gelezen worden met gebruik van PyMuPDF (“PyMuPDF is a high-performance Python library for data extraction, analysis, conversion manipulation of PDF (and other) documents.” 2025). Daarna wordt er ook gezocht aan de hand van de reguliere expressies. Voorbeeld in codefragment 7.5.
- In het geval van een embedded PDF is de inhoud niet onmiddellijk beschikbaar. De gebruiker moet als het ware nog een handeling verrichten (vb. op een knop klikken) alvorens toegang te krijgen tot de inhoud. Dat gaat niet voor een geautomatiseerd script, maar door middel van Selenium (“Selenium with Python”, 2025) kan de gebruikersinteractie nagebootst worden zodat de embedded content toch automatisch gedownload wordt. Voorbeeld in codefragment 7.6. Eenmaal gedownload kan het bestand gewoon geopend worden en doorzocht worden naar een DOI op dezelfde manier als voor PDF documenten.

Er is geen garantie dat de DOI op de webpagina van de publicatie gevonden wordt. Anderzijds is het ook mogelijk dat er meerdere verschillende DOIs gevonden worden. Voor een mens is het vaak evident om te weten welke DOI dan juist is, maar

```
1      def search_in_pdf(pdf, link):
2          doc = pymupdf.Document(stream=pdf)
3          # Extract all Document Text
4          text = chr(12).join([page.get_text() for page in doc])
5          patterns = get_patterns()
6          doi_result = None
7          while (doi_result is None) and (len(patterns) > 0):
8              pattern = re.compile(patterns.pop(), re.IGNORECASE)
9              doi_result = pattern.search(text)
10
11          if doi_result is not None:
12              #update DOI
13              doi = doi_result.group(0)
14              link.doi = doi
15              link.is_doi_success = True
16              link.log_message = "DOI successfully retrieved"
```

Codefragment 7.5: Codefragment voor het openen van een online pdf.

voor geautomatiseerd script is daar geen context voor.

Er kan dus besloten worden dat de DOI met grote zekerheid achterhaald kan worden op basis van de URL of adhv. Crossref. In het geval van een opzoeking op de webpagina van de publicatie is de vindkans een pak kleiner.

```

1      ...
2      url = link.location_replace_url
3
4      options = webdriver.ChromeOptions()
5
6      download_folder =
7          os.path.join(str(Path(__file__).parent.parent.parent.parent),
8              "online_pdf")
9
10     logging_service.logger.debug(f"Download folder:
11         {download_folder}")
12
13     profile = {
14         "plugins.plugins_list": [{"enabled": False, "name":
15             "Chrome PDF Viewer"}],
16         "download.default_directory": download_folder,
17         "download.extensions_to_open": "",
18         "download.prompt_for_download": False,
19         "download.directory_upgrade": True,
20         "plugins.always_open_pdf_externally": True
21     }
22     options.add_experimental_option("prefs", profile)
23     options.add_argument("start-maximized") # open Browser in
24         maximized mode
25     options.add_argument("disable-infobars") # disabling
26         infobars
27     options.add_argument("--disable-extensions") # disabling
28         extensions
29     options.add_argument("--disable-gpu") # applicable to
30         windows os only
31     options.add_argument("--disable-dev-shm-usage") #
32         overcome limited resource problems
33     options.add_argument("--no-sandbox") # Bypass OS security
34         model
35     options.add_argument("--headless")
36
37     logging_service.logger.debug("Downloading file from link:
38         {}".format(link.location_replace_url))
39     print("Downloading file from link:
40         {}".format(link.location_replace_url))
41
42     driver = webdriver.Chrome(options=options)
43     driver.get(url)
44
45     logging_service.logger.debug("Status: Download
46         Complete.")
47     print("Status: Download Complete.")

```

8

Semantic search

8.1. Inleiding

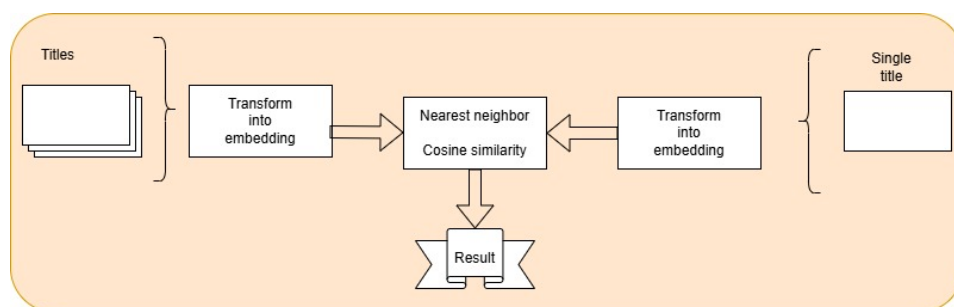
Alle publicaties in IMIS moeten uniek zijn. Op basis van voorgaande stappen is dat met 100% zekerheid te bepalen in het geval de DOI van de publicatie gevonden werd. Als de DOI niet gevonden werd, kan er niet met zekerheid gesproken worden, maar er kan wel een score berekend worden dat de publicatie reeds in IMIS zit. Dat is mogelijk aan de hand van semantic search.

Op basis van de titel van het artikel wordt een embedding berekend. Vervolgens wordt die vergeleken met alle embeddings van de titels in IMIS om zo tot een score te komen van de gelijkenis tussen beide. Dit wordt voorgesteld in figuur 8.1

Indien de score hoog is dan is er veel waarschijnlijkheid dat de publicatie reeds in IMIS zit. In het geval van een lage score is de waarschijnlijkheid eerder klein.

8.2. Chroma

Embeddings worden opgeslaan in een vector databank. Er zijn talloze producten beschikbaar:



Figuur 8.1: Semantic search.

```
1      ...
2      chroma_client = chromadb.Client()
3      # https://docs.trychroma.com/docs/collections/configure
4      self.collection = chroma_client.create_collection(
5          name="my_collection",
6          metadata={
7              "hnsw:space": "cosine"
8          }
9      )
10     ...
```

Codefragment 8.1: Codefragment voor het connecteren met Chroma.

- Pinecone
- Chroma
- Weaviate
- Qdrant
- Milvus
- Vespa
- SingleStore
- Redis
- Elastic Stack
- Mongo
- ...
- (Every database that can store an n-array of numbers)

Chroma ("Chroma", [2025](#)) is een gebruiksvriendelijke database die toelaat om semantic search lokaal te testen zonder extra kosten, zonder account, en zonder installatie van andere software.

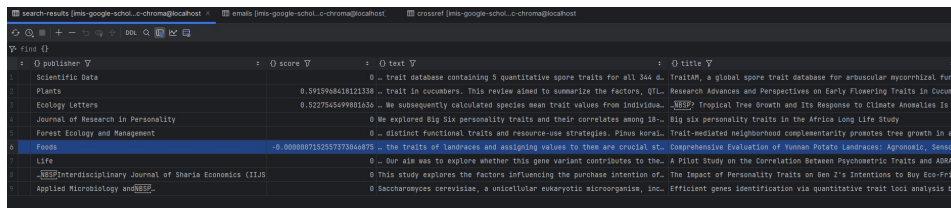
Het online artikel "How to Use Chroma to Build Your First Similarity Search" ([2025](#)) beschrijft stap voor stap hoe een semantic search met Chroma uitgevoerd kan worden. De code op github "Semantic search using Chroma" ([2025](#)) maakt hiervan een implementatie aangepast voor de Google Scholar alerts. De connectie met de Chroma database is te zien in codefragment [8.1](#). Vervolgens worden alle titels uit IMIS opgevraagd en worden de overeenkomstige embeddings berekend aan de hand van de code in codefragment [8.2](#). Dan kan een titel opgezocht worden zoals

```
1         def initialize_embeddings(self):
2             result = re-
3                 quests.get(`https://www.vliz.be/nl/imis?count=2000&module=ref&searchs
4             publications = result.json()
5             documents = []
6             ids = []
7             count = 1
8             for publication in publications:
9                 #self.logging_service.logger.debug(publication)
10                documents.append(publication['StandardTitle'])
11                ids.append(f"id{count}")
12                count+=1
13
14            self.collection.add(
15                documents=documents,
16                ids=ids
17            )
```

Codefragment 8.2: Codefragment voor het berekenen van de embeddings.

```
1         def do_semantic_search(self, title):
2             results = self.collection.query(
3                 query_texts=[title], # Chroma will embed this for you
4                 n_results=2 # how many results to return
5             )
6             return results['distances'][0][0]
```

Codefragment 8.3: Codefragment voor het opzoeken van een titel.



Figuur 8.2: Chroma resultaat.

```

...
self.model =
    SentenceTransformer("nomic-ai/nomic-embed-text-v1",
                        trust_remote_code=True)
...

```

Codefragment 8.4: Codefragment voor het configureren van het model met MongoDB.

getoond wordt in codefragment 8.3.

De test was succesvol. Bij wijze van proef werd de titel van 1 van de zoekresultaten toegevoegd aan de lijst van titels uit IMIS. Het resultaat is te zien in figuur 8.2. Deze oplossing voegt Chroma toe als database naast MongoDB die nodig is voor de Web scraping. Het zou beter zijn indien de embeddings ook in MongoDB opgeslaan kunnen worden. Dat wordt verder onderzocht.

8.3. MongoDB

Atlas Vector Search ("Vector Search", 2025) laat toe om de embeddings tezamen op te slaan met de originele data waar ze van afgeleid zijn.

De online documentatie "Atlas Vector Search Quick Start" (2025) beschrijft stap voor stap hoe een semantic search met Atlas Vector Search lokaal uitgevoerd kan worden. De code op github "Semantic search using MongoDB" (2025) maakt hiervan een implementatie aangepast voor de Google Scholar alerts. De selectie van het model voor de embeddings is te zien in codefragment 8.4. Vervolgens worden alle titels uit IMIS opgevraagd en worden de overeenkomstige embeddings berekend aan de hand van de code in codefragment 8.5 Vervolgens wordt de lijst met embedding geïndexeerd zoals getoond wordt in codefragment 8.6. Dan kan een titel opgezocht worden zoals getoond wordt in codefragment 8.7.

De test was succesvol. Bij wijze van proef werd de titel van 1 van de zoekresultaten toegevoegd aan de lijst van titels uit IMIS. Het resultaat is te zien in figuur 8.3. Deze oplossing breidt de bestaande MongoDB databank uit met embeddings. Voor artikels zonder indentificatie kan een score berekend worden die duidt op duplicaten in IMIS.

```

1         def initialize_embeddings(self):
2             docs = []
3             data = []
4             embeddings = []
5             result = requests.get(IMIS)
6             publications = result.json()
7             for publication in publications:
8                 data.append(publication['StandardTitle'])
9             embed-
                dings.append(self.get_embedding(publication['StandardTitle']))
10
11         for i, (embedding, title) in enumerate(zip(embeddings,
12             data)):
13             doc = {
14                 "_id": i,
15                 "title": title,
16                 "embedding": embedding,
17             }
18             docs.append(doc)
19             self.db_service.set_collection("embeddings")
20             self.db_service.insert_many(docs)

```

Codefragment 8.5: Codefragment voor het berekenen van de embeddings.

publisher	score	text	title
Scientific Data	0	trait database containing 5 quantitative spore traits for all 144 descr.	TraitAM, a global spore trait database for
Plants	0	trait in cucumbers. This review aimed to summarize the factors, QTL map.	Research Advances and Perspectives on Early
Ecology Letters	0.7465339303016663	We subsequently calculated species mean trait values from individual tr.	Tropical Tree Growth and Its Respons
Journal of Research in Personality	0.7654287219647546	We explored Big Six personality traits and their correlates among 18-year.	Big six personality traits in the Africa lo
Forest Ecology and Management	0	distinct functional traits and resource-use strategies. Pinus koraiensi.	Trait-mediated neighborhood complementarity
Foods	0	the traits of landraces and assigning values to them are crucial steps	Comprehensive Evaluation of Yunnan Potato L
Life	1	Our aim was to explore whether this gene variant contributes to the psy.	A Pilot Study on the Correlation Between Ps
Interdisciplinary Journal of Sharia Economics (IISJ)	0	This study explores the factors influencing the purchase intention of env.	The Impact of Personality Traits on Gen Z's
Applied Microbiology and	0	Saccharomyces cerevisiae, a unicellular eukaryotic microorganism, include.	Efficient genes identification via quantita

Figuur 8.3: Mongoddb resultaat.

```
1         def create_search_index(self):
2             search_index_model = SearchIndexModel(
3             definition={
4                 "fields": [
5                     {
6                         "type": "vector",
7                         "path": "embedding",
8                         "similarity": "dotProduct",
9                         "numDimensions": 768
10                    }
11                ]
12            },
13            name="vector_index",
14            type="vectorSearch"
15        )
16        self.db_service.set_collection("embeddings")
17        self.db_service.create_search_index(model=search_index_model)
```

Codefragment 8.6: Codefragment voor het indexeren van de embeddings.

```
1         def do_semantic_search(self, title):
2             query_embedding = self.get_embedding(title)
3             pipeline = [
4                 {
5                     "$vectorSearch": {
6                         "index": "vector_index",
7                         "queryVector": query_embedding,
8                         "path": "embedding",
9                         "exact": True,
10                        "limit": 5
11                    }
12                },
13                {
14                    "$project": {
15                        "_id": 0,
16                        "title": 1,
17                        "score": {
18                            "$meta": "vectorSearchScore"
19                        }
20                    }
21                }
22            ]
23             self.db_service.set_collection("embeddings")
24             result = self.db_service.aggregate(pipeline)
```

Codefragment 8.7: Codefragment voor het opzoeken van een titel.

9

Conclusie

De IMIS collecties van het VLIZ kunnen voortaan semi-automatisch uitgebreid worden door het onderzoek van deze bachelorproef. Elke nieuwe publicatie afkomstig van een Google Scholar zoekopdracht krijgt een score die aangeeft of de publicatie reeds in IMIS zit. Voor alle publicaties die uniek identificeerbaar zijn, krijgen deze een score van 100% of 0% en kunnen zo in het eerste geval automatisch aan IMIS toegevoegd worden. Voor alle publicaties die niet uniek identificeerbaar zijn zal de score tussen 0% en 100% liggen. In die gevallen is manuele tussenkomst nodig om te beslissen of de publicatie al dan niet aan IMIS toegevoegd moet worden. Curabitur nunc magna, posuere eget, venenatis eu, vehicula ac, velit. Aenean ornare, massa a accumsan pulvinar, quam lorem laoreet purus, eu sodales magna risus molestie lorem. Nunc erat velit, hendrerit quis, malesuada ut, aliquam vitae, wisi. Sed posuere. Suspendisse ipsum arcu, scelerisque nec, aliquam eu, molestie tincidunt, justo. Phasellus iaculis. Sed posuere lorem non ipsum. Pellentesque dapibus. Suspendisse quam libero, laoreet a, tincidunt eget, consequat at, est. Nullam ut lectus non enim consequat facilisis. Mauris leo. Quisque pede ligula, auctor vel, pellentesque vel, posuere id, turpis. Cras ipsum sem, cursus et, facilisis ut, tempus euismod, quam. Suspendisse tristique dolor eu orci. Mauris mattis. Aenean semper. Vivamus tortor magna, facilisis id, varius mattis, hendrerit in, justo. Integer purus.

Vivamus adipiscing. Curabitur imperdiet tempus turpis. Vivamus sapien dolor, congue venenatis, euismod eget, porta rhoncus, magna. Proin condimentum pretium enim. Fusce fringilla, libero et venenatis facilisis, eros enim cursus arcu, vitae facilisis odio augue vitae orci. Aliquam varius nibh ut odio. Sed condimentum condimentum nunc. Pellentesque eget massa. Pellentesque quis mauris. Donec ut ligula ac pede pulvinar lobortis. Pellentesque euismod. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent elit. Ut laoreet ornare est. Phasellus gravida vulputate nulla. Donec sit amet arcu ut sem tempor

malesuada. Praesent hendrerit augue in urna. Proin enim ante, ornare vel, consequat ut, blandit in, justo. Donec felis elit, dignissim sed, sagittis ut, ullamcorper a, nulla. Aenean pharetra vulputate odio.

Quisque enim. Proin velit neque, tristique eu, eleifend eget, vestibulum nec, lacus. Vivamus odio. Duis odio urna, vehicula in, elementum aliquam, aliquet laoreet, tellus. Sed velit. Sed vel mi ac elit aliquet interdum. Etiam sapien neque, convallis et, aliquet vel, auctor non, arcu. Aliquam suscipit aliquam lectus. Proin tincidunt magna sed wisi. Integer blandit lacus ut lorem. Sed luctus justo sed enim.

Morbi malesuada hendrerit dui. Nunc mauris leo, dapibus sit amet, vestibulum et, commodo id, est. Pellentesque purus. Pellentesque tristique, nunc ac pulvinar adipiscing, justo eros consequat lectus, sit amet posuere lectus neque vel augue. Cras consectetur libero ac eros. Ut eget massa. Fusce sit amet enim eleifend sem dictum auctor. In eget risus luctus wisi convallis pulvinar. Vivamus sapien risus, tempor in, viverra in, aliquet pellentesque, eros. Aliquam euismod libero a sem.

Nunc velit augue, scelerisque dignissim, lobortis et, aliquam in, risus. In eu eros. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Curabitur vulputate elit viverra augue. Mauris fringilla, tortor sit amet malesuada mollis, sapien mi dapibus odio, ac imperdiet ligula enim eget nisl. Quisque vitae pede a pede aliquet suscipit. Phasellus tellus pede, viverra vestibulum, gravida id, laoreet in, justo. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Integer commodo luctus lectus. Mauris justo. Duis varius eros. Sed quam. Cras lacus eros, rutrum eget, varius quis, convallis iaculis, velit. Mauris imperdiet, metus at tristique venenatis, purus neque pellentesque mauris, a ultrices elit lacus nec tortor. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent malesuada. Nam lacus lectus, auctor sit amet, malesuada vel, elementum eget, metus. Duis neque pede, facilisis eget, egestas elementum, nonummy id, neque.

10

Google Scholar zoekopdracht

10.1. Zoeken in Google Scholar (“Google Scholar Guide”, 2025)

10.1.1. Basis zoeken

Google Scholar biedt een vertrouwde gebruikersinterface met een inputveld waar de relevante zoektermen ingevuld moeten worden. Default wordt er in elke taal gezocht, maar de gebruiker kan dit beperken tot zijn eigen taal.

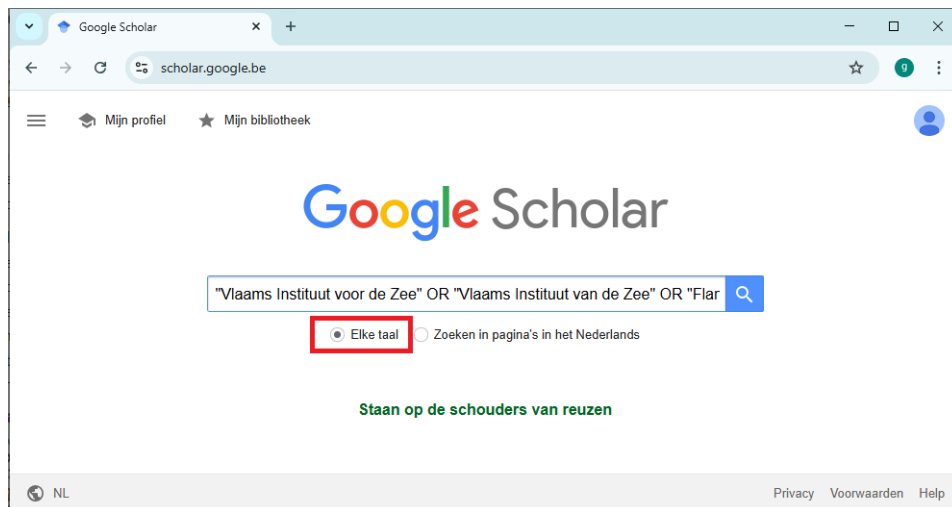
Wanneer een zoekopdracht verzonden wordt, dan is er een antwoord binnen de 3 seconden. Het resultaat kan vervolgens verder gefilterd worden:

- **Elke periode:** Dit is de default filter zodat alle resultaten getoond worden ongeacht hun publicatiedatum.
- **Sinds jaar:** Hierbij worden enkel resultaten gefilterd die sinds het gespecificeerde jaar gepubliceerd werden.
- **Aangepast bereik:** Hierbij worden enkel resultaten gefilterd waarvan de publicatiedatum binnen het gespecificeerde bereik ligt.
- **Sorteren op relevantie:** Dit is de default filter tezamen met 'Elke periode' die de resultaten sorteert op basis van hun belangrijkheid.¹
- **Sorteren op datum:** Hierbij worden de resultaten gesorteerd op publicatiedatum.
- **Reviewartikelen:** Hierbij worden enkel state of the art publicaties gefilterd.²

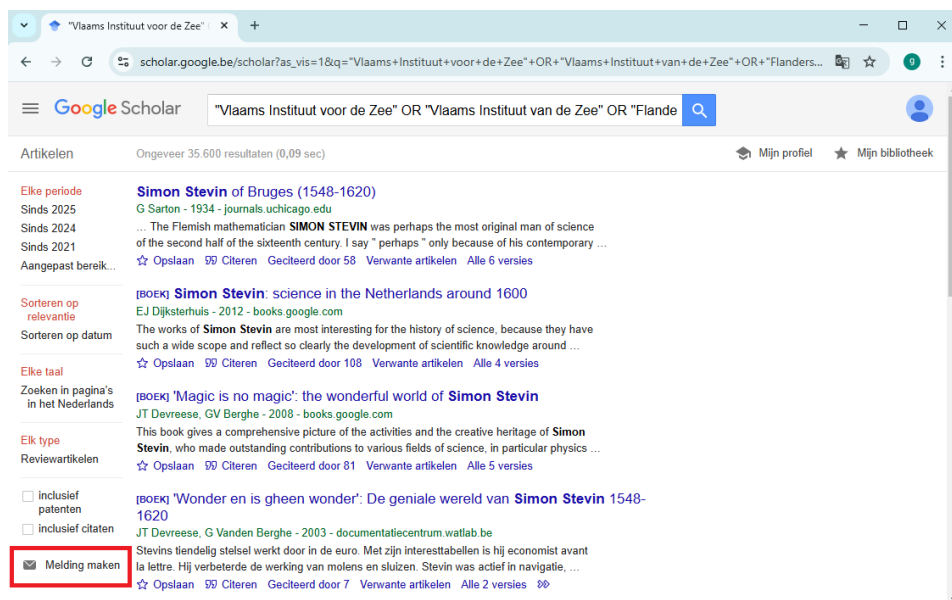
Elk resultaat kan verder uitgediept worden:

¹De relevantie van elke publicatie wordt in de eerste plaats bepaald door het aantal citaties ((Beel & Gipp, 2009))

²Een reviewartikel ondergaat een systematische review door een groep van experts volgens de op dat moment geldende 'State of the art' ((Sataloff e.a., 2021))



Figuur 10.1: Google Scholar user interface voor het basis zoeken van publicaties op basis van de ingevoerde zoektermen.



Figuur 10.2: Google Scholar zoekresultaten op basis van een zoekopdracht.

- **Geciteerd door:** Een oplistijng van publicaties die zelf het artikel citeren. Dit kan leiden tot andere relevante artikels.
- **Verwante artikelen:** Andere artikels in hetzelfde thema. Dit kan leiden tot andere relevante artikels.
- **Alle versies:** Alternatieve locaties waar dezelfde informatie kan teruggevonden worden. Dit kan leiden tot een breder beeld van organisaties, instituten en uitgevers.

10.1.2. Geëvanceerd zoeken

Google Scholar heeft ook een meerdere geëvanceerde zoekopties:

- **Zoek artikels met alle termen:** Combineert zoektermen. Zoekt publicaties die alle termen bevatten.
- **Zoek artikels met de exacte zoekterm:** Zoekt publicaties waar de zoekterm of zin exact in terug te vinden is.
- **Zoek artikels met op zijn minst 1 van de zoektermen:** Zoekt publicaties waar alle of minstens 1 van de zoektermen in voorkomen.
- **Zoek publicaties zonder de zoektermen:** Matcht publicaties waar geen enkele van de zoektermen in voorkomen.

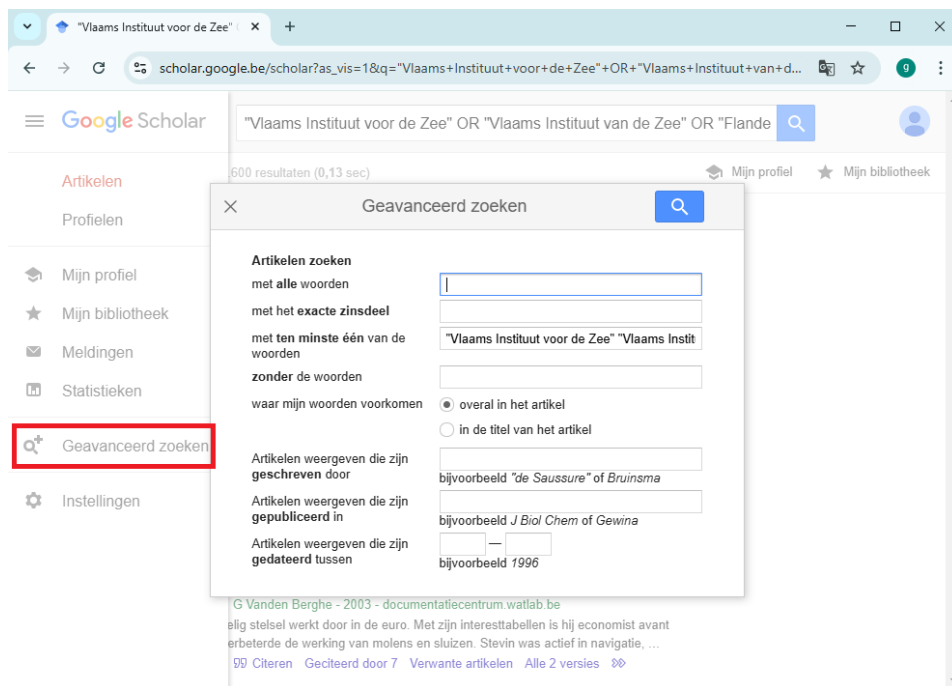
Voor alle bovenstaande filters kan ingesteld worden of er enkel in de titel of overal in de tekst mag gezocht worden.

Daarnaast zijn er 3 bijkomende geëvanceerde filters:

- **Zoek artikels op basis van auteurs:** Zoekt publicaties die geschreven zijn door een bepaalde auteur.
- **Zoek artikels op basis van de uitgever:** Zoekt publicaties die uitgegeven zijn door een bepaalde uitgever.
- **Zoek artikels op basis van publicatiedatum:** Zoekt publicaties die gepubliceerd zijn tussen 2 opgegeven datums.

Alle geëvanceerde filters kunnen verder gespecificeerd worden door middel van logische operatoren: (AND, OR, NOT, AROUND).

- **AND:** Zoekt beide zoektermen in de publicatie.
- **OR:** Zoekt 1 of beide zoektermen in de publicatie.
- **NOT:** Sluit ongewenste tekst uit van het zoekresultaat.
- **AROUND:** Zoekt zoektermen in de ingestelde nabijheid van de opgegeven zoekterm.



Figuur 10.3: Google Scholar user interface voor het geavanceerd zoeken van publicaties op basis van de ingevoerde zoektermen en filters.

Alle geavanceerde filters kunnen verder gespecificeerd worden door middel van hulpwoorden:

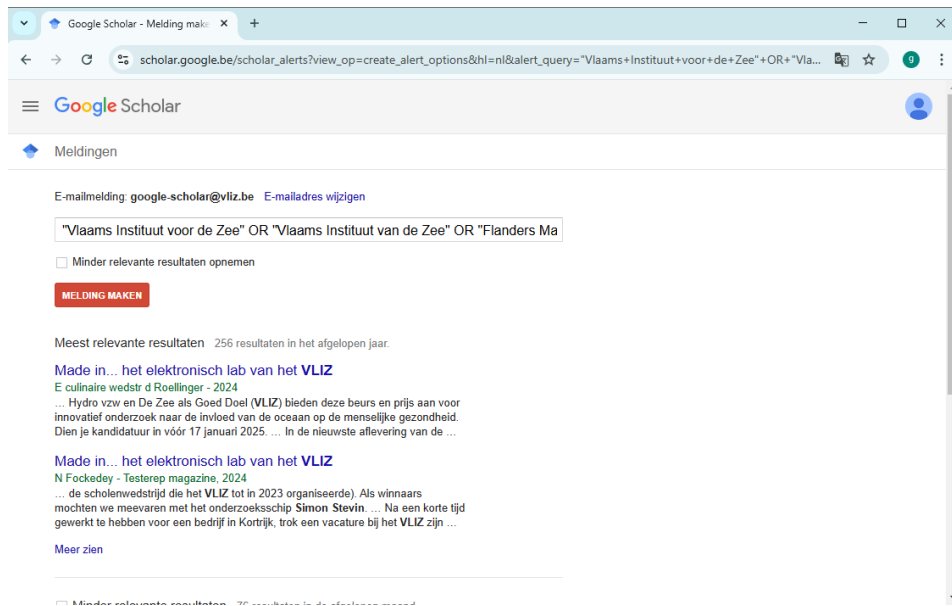
- **intitle:** De zoekresultaten bevatten de opgegeven zoekterm in de titel.
- **intext:** De zoekresultaten bevatten de opgegeven zoekterm in de tekst.
- **author:** De zoekresultaten bevatten de opgegeven auteur.
- **source:** De zoekresultaten bevatten de opgegeven uitgever.

Alle geavanceerde filters kunnen verder gespecificeerd worden door middel van enkele leestekens:

- **aanhalingstekens (""):** De zoekresultaten bevatten de exacte tekst tussen aanhalingstekens.
- **liggend streepje (A-B):** Om aan te tonen dat 2 zoektermen sterk verbonden zijn.
- **liggend streepje (A -B):** Om de tweede zoekterm uit te sluiten van de resultaten.

10.2. E-mail alerts

Het is mogelijk om de zoekresultaten te personaliseren. Voor elke zoekopdracht die wordt aangemaakt, kan een overeenkomstige alert ingesteld worden door te



Figuur 10.4: Google Scholar user interface voor het aanmaken van een e-mail alert voor de ingevoerde zoekopdracht.

klikken op 'Melding maken' zoals getoond op figuur 10.2. Het volstaat om het e-mailadres in te vullen naar waar de alerts verstuurd moeten worden. Dit genereert een verificatie e-mail en na bevestiging is de alert geactiveerd.

Vanaf dan worden nieuwe publicaties die voldoen aan de filtercriteria systematisch doorgestuurd naar het e-mailadres. Overeenkomstig werd een nieuw account **google-scholar@marineinfo.org** aangemaakt, als een gedeeld account waar meerdere gebruikers toegang tot hebben. Vervolgens werd een Google account aangemaakt met hetzelfde e-mailadres.

10.3. Opstellen van een zoekopdracht

Om relevante publicaties over het VLIZ te vinden, worden volgende zoektermen gebruikt:

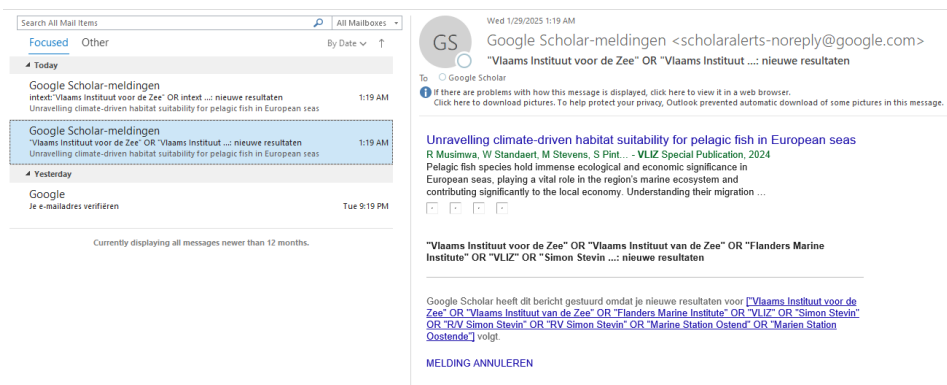
- Vlaams Instituut voor de Zee
- Vlaams Instituut van de Zee
- Flanders Marine Institute
- VLIZ
- Simon Stevin
- R/V Simon Stevin
- RV Simon Stevin
- Marine Station Ostend
- Mariene Station Oostende

De zoekopdracht moet resultaten geven wanneer minstens 1 of meerdere zoektermen voorkomen in de titel of de tekst van het artikel. Dit kan bereikt worden door de zoektermen tussen aanhalingstekens te schrijven en door ze te verbinden met de OR operator.

```

``Vlaams Instituut voor de Zee''
    OR
``Vlaams Instituut van de Zee''
    OR
``Flanders Marine Institute''
    OR
``VLIZ''
    OR
``Simon Stevin''
    OR
``R/V Simon Stevin''
    OR
``RV Simon Stevin''
    OR
``Marine Station Ostend''
    OR
``Mariene Station Oostende''

```



Figuur 10.5: Google Scholar e-mail alert met de nieuwe resultaten sinds het aanmaken van de melding en sinds de vorige melding.

11

Parsen van Google Scholar zoekresultaten

11.1. Doelstelling

De GS alert geeft aanleiding tot e-mails met de meest recente zoekresultaten. Voor elke e-mail en voor elk zoekresultaat wordt de DOI opgezocht. Daarvoor worden 3 stappen doorlopen:

- E-mail lezen
- Body van de e-mail parsen
- DOI opzoeken in de link die aanwezig is in het geparseerde resultaat

Er moet een geautomatiseerd systeem komen dat

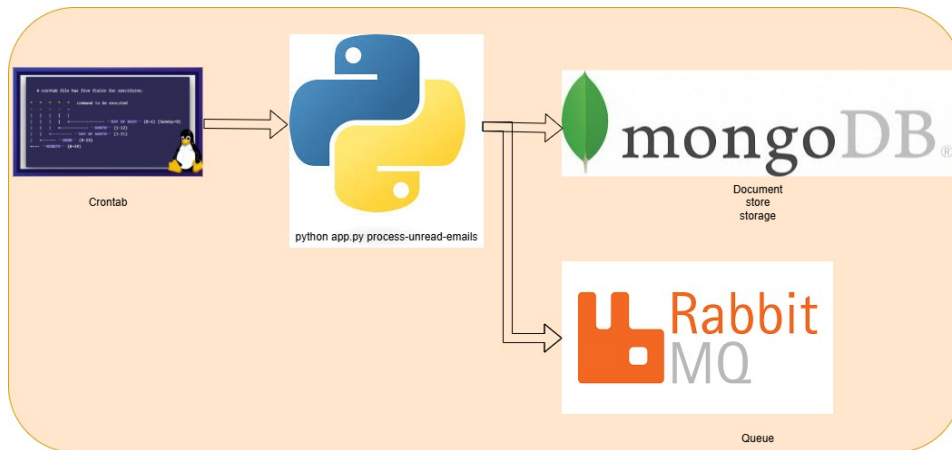
- de binnenkomende e-mail alerts in een mailbox leest
- vervolgens de zoekresultaten in de body parset

11.1.1. E-mail lezen

Op gezette tijdstippen wordt de inbox van het account `google-scholar@marineinfo.org` geopend op zoek naar ongelezen e-mails. Elke e-mail wordt gelezen en de informatie van de e-mail wordt opgeslaan. De e-mail wordt verplaatst van de inbox naar een werkmap overeenkomstig het onderwerp van de e-mail. Dit wordt schematisch voorgesteld in figuur 11.1.

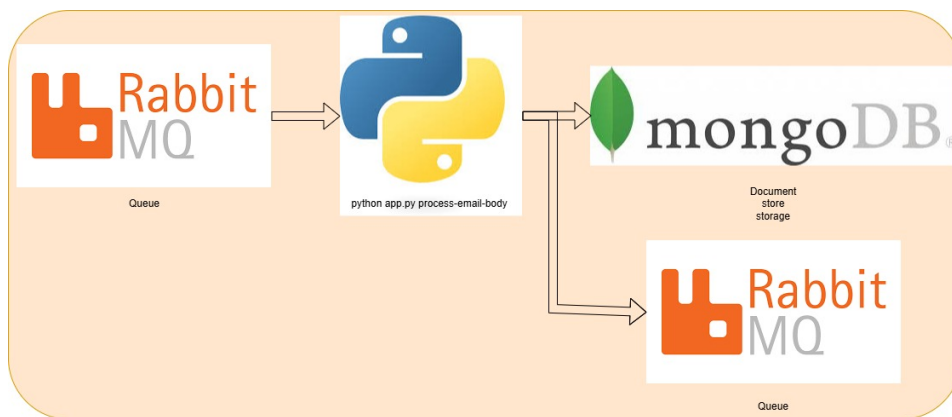
11.1.2. Body van de e-mail parsen

Elke e-mail bevat 1 of meerdere zoekresultaten afkomstig van GS. In deze stap wordt elk zoekresultaat uit de body van de e-mail gefilterd en afzonderlijk opgeslaan. De



Figuur 11.1: Flow e-mail lezen.

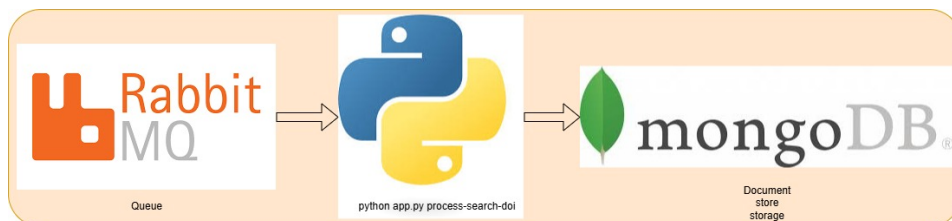
informatie van het zoekresultaat wordt opgeslaan. Dit wordt schematisch voorgesteld in figuur 11.2.



Figuur 11.2: Flow e-mail body parsen.

11.1.3. DOI opzoeken in de link van het geparseerde resultaat

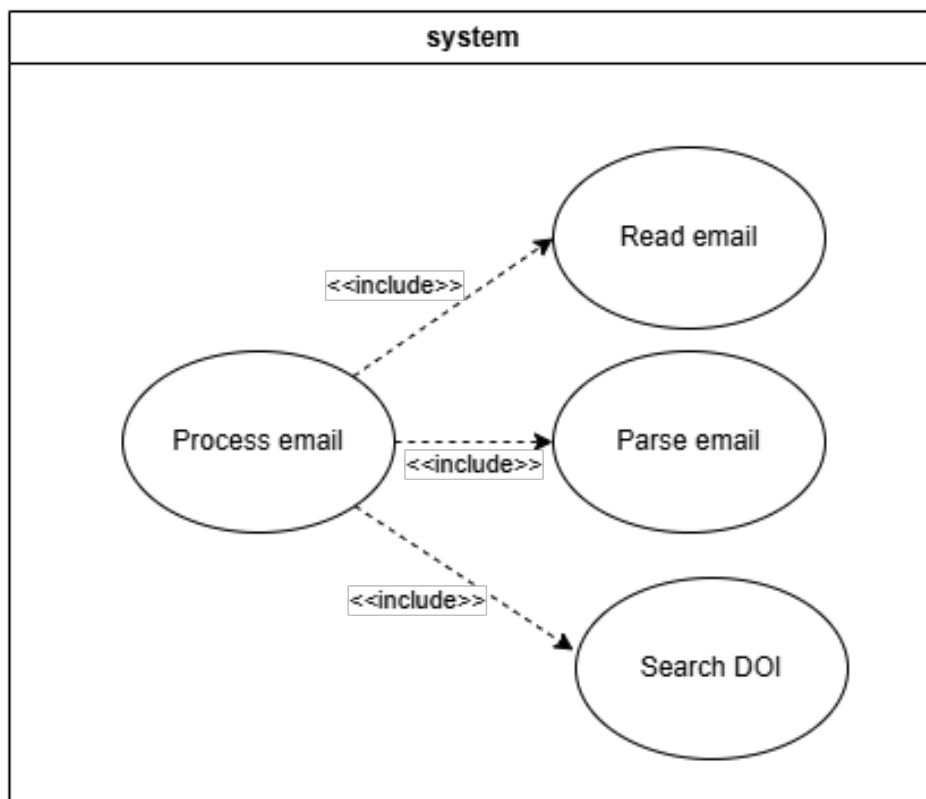
Het GS zoekresultaat zelf bevat geen DOI. Maar elk zoekresultaat heeft wel een weblink die werd opgeslaan in de vorige stap. Hier wordt een request gestuurd naar de link. Het response wordt vervolgens verwerkt met als doel om het DOI erin op te zoeken. De informatie van het DOI wordt opgeslaan. Dit wordt schematisch voorgesteld in figuur 11.3.



Figuur 11.3: DOI opzoeken in de link van het geparseerde resultaat.

11.2. Analyse

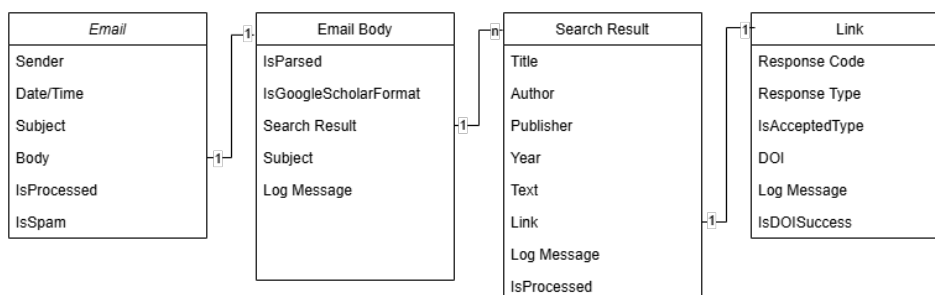
11.2.1. Use Cases



Figuur 11.4: Use case diagram.

11.2.2. Domeinmodel

De 3 bovenstaande stappen geven aanleiding tot 4 objecten die de kern van het domeinmodel uitmaken en de nodige properties bevatten om de tussenresultaten bij te houden alsook de uiteindelijke DOI.



Figuur 11.5: Domeinmodel.

Use Case Read email	GS stuurt e-mails naar google-scholar@marineinfo.org. Deze e-mails moeten automatisch verwerkt worden.
primary actor	-
stakeholders	-
pre-condities	onverwerkte e-mail
post-condities	verwerkte e-mails
normaal verloop	<ol style="list-style-type: none"> 1. systeem opent e-mail 2. systeem leest afzender 3. systeem evalueert afzender (DR-email) 4. systeem leest datum / tijdstip 5. systeem leest onderwerp 6. systeem leest body 7. systeem slaat deze gegevens op 8. systeem verplaatst e-mail van inbox naar mailbox folder op basis van onderwerp
alternatief verloop	<p>3A. de afzender is verkeerd (DR-email)</p> <p>3A1. systeem verplaatst e-mail van inbox naar SPAM folder</p>
domeinregels	DR-email: afzender moet zijn: scholaralerts-noreply@google.com

Tabel 11.1: GS stuurt e-mails naar google-scholar@marineinfo.org. Deze e-mails moeten automatisch verwerkt worden.

Use Case Parse email	De body van een e-mail wordt ontleed en de informatie wordt gestructureerd opgeslaan.
primary actor	-
stakeholders	-
pre-condities	onverwerkte body van e-mail
post-condities	verwerkte body van e-mail
normaal verloop	<ol style="list-style-type: none"> 1. systeem ontvangt de body van een e-mail 2. systeem ontleedt de body en zoekt volgende onderdelen (DR-Google-Scholar formaat) <ol style="list-style-type: none"> a. lijst <ol style="list-style-type: none"> i. titel (DR-uniek) ii. auteur iii. uitgever iv. jaartal v. tekst (DR-uniek) vi. link (DR-uniek) b. onderwerp 3. systeem slaat bovenstaande informatie gestructureerd op
alternatief verloop	<ol style="list-style-type: none"> 2A. de verwachte informatie wordt niet gevonden (DR-Google-Scholar formaat) <ol style="list-style-type: none"> 2A1. systeem logt het probleem 2A2. use case eindigt zonder bereiken van de post-conditie 2B. de body is onleesbaar <ol style="list-style-type: none"> 2B1. systeem logt het probleem 2B2. use case eindigt zonder bereiken van de post-conditie 2C. titel en/of tekst en/of link zijn niet uniek (zijn reeds eerder opgeslaan) <ol style="list-style-type: none"> 2C1. systeem logt het voorval 2C2. systeem gaat verder met stap 2 van het normale verloop
domeinregels	<ul style="list-style-type: none"> • DR-Google-Scholar-formaat • DR-uniek

Tabel 11.2: De body van een e-mail wordt ontleed en de informatie wordt gestructureerd opgeslaan.

Use Case Search DOI	Zoek een DOI in een tekstueel document.
primary actor	-
stakeholders	-
pre-condities	systeem beschikt over een online-url
post-condities	het DOI is gevonden
normaal verloop	<ol style="list-style-type: none"> 1. systeem stuurt een request naar de link 2. systeem ontvangt response 3. systeem bepaalt type van de response 4. systeem systeem kiest de juiste tools in functie van het type (DR-type) 5. systeem zoekt DOI in de response 6. systeem slaat het gevonden DOI op
alternatief verloop	<ol style="list-style-type: none"> 2A. geen response 2A1. het systeem logt het probleem 2A2. use case eindigt zonder bereiken van de post-conditie 4A. type is niet ondersteund (DR-type) 4A1. het systeem logt het probleem 4A2. use case eindigt zonder bereiken van de post-conditie 5A. DOI wordt niet gevonden 5A1. het systeem logt het probleem 5A2. use case eindigt zonder bereiken van de post-conditie 5B. meerdere DOIs worden gevonden 5A1. het systeem logt het probleem 5A2. use case eindigt zonder bereiken van de post-conditie
domeinregels	DR-type: pdf, html

□

Tabel 11.3: Zoek een DOI in een tekstueel document.

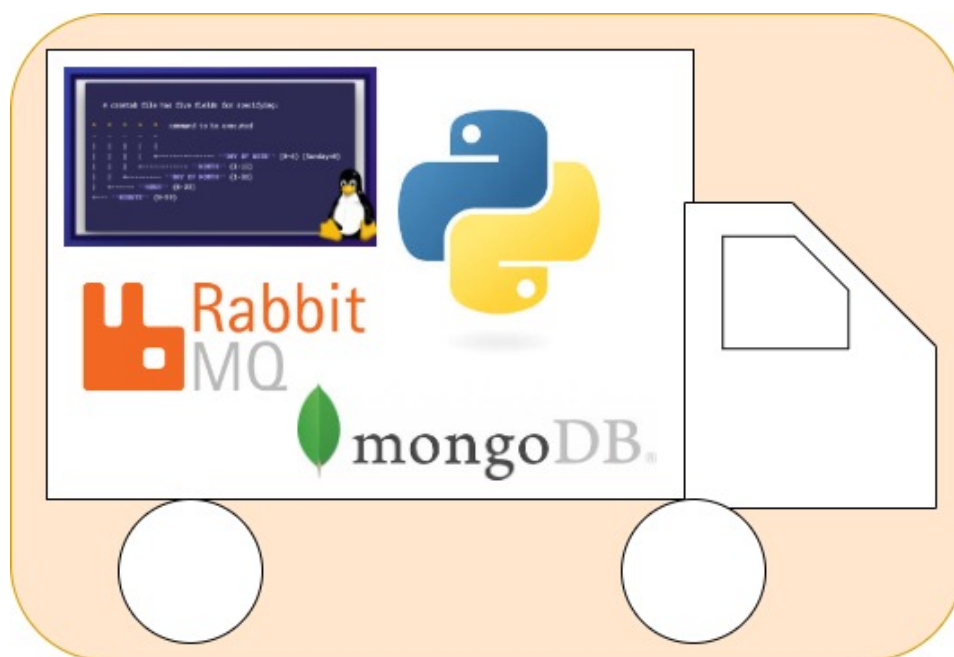
11.3. Implementatie

Daarvoor moet er custom code geschreven worden en bijhorend moet de vraag gesteld worden welke technology stack er gebruikt zal worden? Daarbij zijn 2 aspecten belangrijk:

- Wat is de bestaande technology stack van de klant en kan de te ontwikkelen software daarin ondergebracht worden?
- Welke technologie is het meest geschikt om het gestelde probleem op te lossen?

Voor het eerste punt zijn er bij de klant 2 pijlers:

- Alle websites en datasystemen zijn geschreven in Php, al dan niet gebruik makend van het Symfony framework. Voor hun data steunen ze voornamelijk op SqlServer en op PostgreSQL. Voor hun user interface steunen ze hoofdzakelijk op Twig en op NextJS.
- Veel data processing scripts draaien in hun eigen docker container. Ze gebruiken vooral R en Python als programmeertaal.



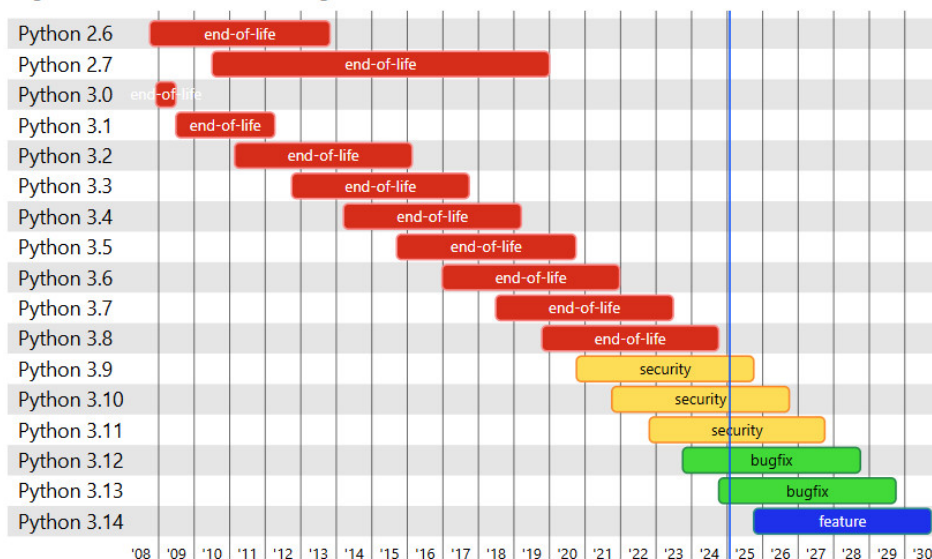
Figuur 11.6: Technology Stack.

Onze opdracht past het best binnen de tweede pijler. Het is niet opportuun om een webserver te belasten met het systematisch herhalen van een opdracht die los staat van de websites die hij host. De GS alerts zijn op zich ook data die in hun eigen container verwerkt zullen worden.

Voor het tweede punt is de technologie zeer duidelijk. Python (maar ook R) beschikken over hele grote bibliotheken die toelaten om uiteenlopende soorten data te verwerken.

Het programmeerwerk zal uitgevoerd worden in Python. Op het moment van schrijven is Python versie 3.13.1 de meest recente stabiele versie.

Python release cycle



Figuur 11.7: Python release cycle. ("Status of Python versions", 2025)

De geparseerde e-mail alerts worden opgeslagen in een MongoDB database, dit is een NoSQL document store databank. Een RabbitMQ broker zorgt ervoor dat voor elke e-mail iedere stap in de flow doorlopen wordt. Tenslotte wordt alles ondergebracht in zijn eigen Docker container die binnen het project met elkaar communiceren. De GS alert geeft aanleiding tot e-mails met de meest recente zoekresultaten. Voor elke e-mail en voor elk zoekresultaat wordt de DOI opgezocht. Daarvoor worden 3 stappen doorlopen:

- E-mail lezen
- Body van de e-mail parsen
- DOI opzoeken in de link die aanwezig is in het geparseerde resultaat

11.3.1. E-mail lezen

Meerdere specifieke Python libraries ondersteunen alle mogelijke bewerkingen met e-mail:

- **imaplib - IMAP4 protocol client** ("imaplib - IMAP4 protocol client", 2025):
Legt een verbinding met een IMAP server en implementeert de functionaliteit

ten van het IMAP¹ protocol om e-mails op te halen. ("Internet Message Access Protocol, Uit Wikipedia, de vrije encyclopedie", 2025)

- **email - An email and MIME handling package ("email - An email and MIME handling package", 2025):** Laat toe om emails te verwerken adhv. 3 modules.
 - **message:** Alle bewerkingen met een e-mail.
 - **parser:** Om een e-mail om te zetten in tekst.
 - **generator:** Om tekst om te zetten in een e-mail.

De gelezen e-mail wordt opgeslaan in een document store database ("Document-oriented database", 2025). Dit type databank is geschikt voor het bijhouden van dynamische data die als geheel (als één document) opgeslaan wordt. MongoDB ("MongoDB", 2025) is de meest gekende technologie van dit type databank. Py-mongo ("PyMongo", 2025) is de specifieke python library die de verbinding met de Mongoddb server mogelijk maakt.

Het custom script dat de e-mail leest wordt opgeroepen door de Windows task scheduler. Click_ ("Click_■", 2025) laat toe om meerdere commando's te definiëren in hetzelfde script.

Tenslotte wordt er gebruik gemaakt van de Dependency Injector ("Dependency Injector", 2025) om al deze services door middel van dependency injection ("Dependency injection", 2025) ter beschikking te stellen aan het script.

Aan het einde van het process dat de e-mail leest, wordt de body van de e-mail in een wachtrij geplaatst zodat deze opgepikt wordt door de volgende stap.

11.3.2. Body van de e-mail parsen

De body van de e-mail komt overeen met de Search Engine Result Page (SERP) en is een gestructureerde html pagina. Voor het parsen wordt er gebruik gemaakt van deze structuur. Elk afzonderlijk resultaat heeft een H3-tag met als klasse 'gse_alrt_title' waarin de titel staat. De korte tekst van het zoekresultaat staat in een DIV-tag met als klasse 'gse_alrt_sni'. Tussen de titel en het snippet staan de auteurs, de uitgever en het jaartal gegroepeerd in een DIV-tag zoals te zien is in figuur 11.8.

¹Internet Message Access Protocol



11.3.3. DOI opzoeken in de link van het geparseerde resultaat

De url die in het zoekresultaat van GS zit, verwijst zelf ook naar de GS website, zie codefragment 11.2. Daarin valt te zien dat de eigenlijke url van het zoekresultaat verpakt zit binnen het GS adres. Inderdaad, het response heeft een HTTP 302 status code en redirect het de request naar de eigenlijke url. Het verwerkende script kan de uiteindelijke url opzoeken in de javascript van het response aan de hand van de **location.replace** instructie. Voorbeeld in codefragment 11.3. Met de url uit de **location.replace** kan er nu opnieuw een request gestuurd worden. Het response

```
headers = "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/62.0.3202.94 Safari/537.36", "Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/apng,*/*;q=0.8", "Accept-Encoding": "gzip, deflate, br", "Accept-Language": "en-US,en;q=0.9,it;q=0.8,et;q=0.7,de;q=0.6",
```

Codefragment 11.1: Voorbeeld headers.

```
https://scholar.google.be/scholar_url?url=https://www.nature.com/articles/s41598-024-83657-0&hl=nl&sa=X&d=1146959663452366258&ei=tRGsZ4eaLIC96rQP29mI6AY&scisig=AFWwaKYGyEQWAVLJNk68&html=&pos=0&folt=kW-top
```

Codefragment 11.2: Embedded url.

```
1      pattern = r"location\.replace\(['\"](?:[^\\""]+|'\\.')*\(['\"])\)"
2      match = re.search(pattern, js_code)
3      if match:
4          location_replace_url = match.group(1)
5          self.logging_service.logger.debug(f"Extracted
           location.replace URL for search result link:
           {location_replace_url}")
6      else:
7          link.log_message = "No location.replace url found for search
           result link"
8          self.logging_service.logger.debug("No location.replace url
           found for search result link")
```

Codefragment 11.3: Codefragment voor het vinden van de location.replace.

```

1      doc = pymupdf.Document(stream=pdf)
2      # Extract all Document Text
3      text = chr(12).join([page.get_text() for page in doc])

```

Codefragment 11.4: Codefragment voor het openen van een online pdf.

dat daar teruggestuurd wordt is verschillend voor elke zoekresultaat. Het is dus niet mogelijk om het DOI voor elk zoekresultaat op dezelfde manier te achterhalen. Daarom wordt er een stapsgewijze werkwijze gehanteerd, van zodra het DOI gevonden is worden de volgende stappen overgeslaan:

- **het DOI opzoeken in de link:** De link zelf bevat soms zelf reeds het DOI. Dit wordt getoets aan de hand van enkele reguliere expressies.
- **het DOI opzoeken in de content** Het antwoord komt overeen met een html-pagina of een pdf-document. Andere return types worden buiten beschouwing gelaten. De inhoud van het antwoord wordt gelezen en met gebruik van dezelfde reguliere expressie kan de aanwezigheid van het DOI opgespoord worden. Voor pdf-documenten is er een extra tussenstap nodig: de inhoud moet gelezen worden met gebruik van ("PyMuPDF is a high-performance Python library for data extraction, analysis, conversion manipulation of PDF (and other) documents." [2025](#)). Voorbeeld in codefragment [11.4](#).
- **het DOI opzoeken in de embedded content:** Als tot zover er nog steeds geen DOI gevonden is, dan kan het zijn dat het response een embedded document bevat. Op een website kan embedded content ontsloten worden door erop te klikken. Dat gaat hier niet. Door middel van ("Selenium with Python", [2025](#)) kan de embedded content automatisch gedownload worden. Eenmaal gedownload kan het bestand gewoon geopend worden en doorzocht worden naar een DOI op dezelfde manier als voor pdf-documenten in de voorgaande stap.

Alle DOIs hebben dezelfde structuur: ze beginnen met het cijfer 10, gevolgd door een punt en 4 tot 9 cijfers, daarna volgt een slash. Verder kan elke willekeurige opeenvolging van letters, cijfers, speciale tekens en slashes voorkomen. De volgende lijst reguliere expressies is aanbevolen om DOIs op te zoeken "DOIs and matching regular expressions" ([2025](#)):

- `/^10.\d{4,9}/[-._;()/:A-Z0-9]+$ /i`
- `/^10.1002/[^s]+$ /i`
- `/^10.\d4/\d+-\d+X?(\d+)\d+<[\d\w]+:[\d\w]*>\d+.\d+.\w+;\d$/i`
- `/^10.1021/\w\w\d++$/i`

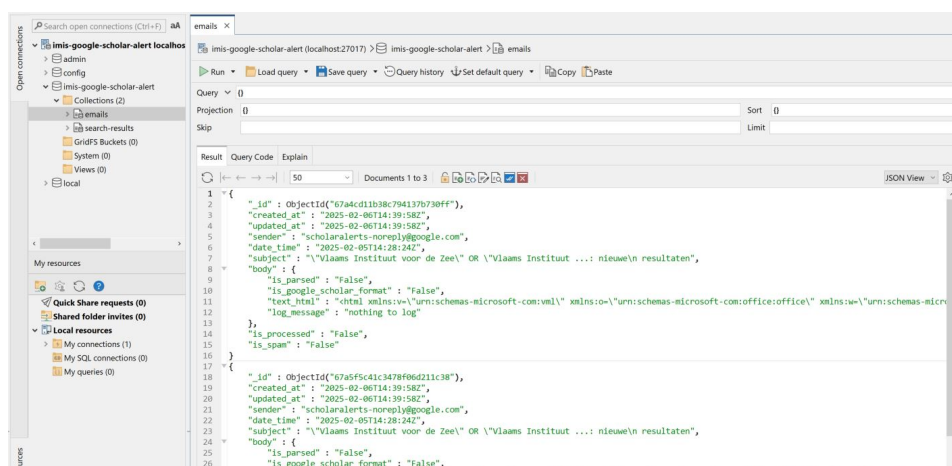
• `/^10.1207/[\w\d]+\s\d+_\d+$/i`

Het DOI wordt opgeslaan in de databank.

11.4. Resultaat

11.4.1. E-mail lezen

Het resultaat van deze stap is eerder een tussenresultaat waarbij verschillende onderdelen van de e-mail zoals datum, onderwerp, en andere opgeslaan worden in de databank. Deze opgeslaan informatie vormt de bron voor de verdere verwerking in stap 2.



Figuur 11.9: Databank JSON view.

The screenshot shows the Databank interface with the table view selected. The data is displayed in a table with the following columns: `_id`, `created_at`, `updated_at`, `sender`, `date_time`, `subject`, `body`, `is_processed`, and `is_spam`.

_id	created_at	updated_at	sender	date_time	subject	body	is_processed	is_spam
67a4cd11b38c79...	2025-02-06T14...	2025-02-06T14...	scholaralerts-nor...	2025-02-05T14...	Vlaams Instituut...	[4 fields]	False	False
67a5f5c41c3478f...	2025-02-06T14...	2025-02-06T14...	scholaralerts-nor...	2025-02-05T14...	Vlaams Instituut...	[4 fields]	False	False
67a5d0901123d...	2025-02-07T13...	2025-02-07T13...	Bart De Paep...	2025-02-05T14...	Vlaams Instituut...	[4 fields]	False	False

Figuur 11.10: Databank table view.

Een overzicht van de gegevens die opgeslaan worden zijn te zien in figuur 11.11

- **_id**: De unieke ID van de record.
- **created_at**: De datum en het tijdstip waarop de e-mail in de databank werd opgeslaan.
- **updated_at**: De datum en het tijdstip waarop de e-mail in de databank het laatst werd gewijzigd.

```
{
  "_id" : ObjectId("67b3425b3ac6e357ee41b781"),
  "created_at" : "2025-02-17T14:06:19Z",
  "updated_at" : "2025-02-17T14:06:19Z",
  "sender" : "Google Scholar-meldingen <scholaralerts-noreply@google.com>",
  "date_time" : "2025-02-12T20:43:52Z",
  "subject" : "traits: nieuwe resultaten",
  "body" : {
    "text_html" : "<!doctype html><html xmlns=\"http://www.w3.org/1999/xhtml\" xmlns:o=\"urn:schemas-microsoft-com:office:office>"
  },
  "is_processed" : false,
  "is_spam" : false,
  "log_message" : "Email read successfully."
}
```

Figuur 11.11: Databank structuur e-mail.

- **sender:** De afzender van de e-mail.
- **date_time:** De datum en het tijdstip waarop de e-mail verstuurd werd.
- **subject:** Het onderwerp van de e-mail.
- **body:**
 - **text_html:** De body van de e-mail.
- **is_processed:** Een vlag die aangeeft of de body van de e-mail verwerkt werd.
- **is_spam:** Een vlag die aangeeft of de e-mail niet van GS afkomstig is.
- **log_message:** De logberichten tijdens het lezen van de e-mail worden hier bewaard.

11.4.2. Body van de e-mail parsen

Het resultaat van deze stap is eerder een tussenresultaat waarbij voor alle zoekresultaten in de body van de e-mail, verschillende onderdelen zoals titel, snippet, auteurs, uitgever, jaartal, en link opgeslaan worden in de databank. Deze opgeslaan informatie vormt de bron voor de verdere verwerking in stap 3. De gegevens van de e-mail worden bijgewerkt om aan te duiden dat de body verwerkt is, zoals te zien is in figuur 11.12.

```
{
  "_id" : ObjectId("67b5cf80ba5d6fa0d55d98eb"),
  "created_at" : "2025-02-19T12:33:04Z",
  "updated_at" : "2025-02-19T12:33:44Z",
  "sender" : "Google Scholar-meldingen <scholaralerts-noreply@google.com>",
  "date_time" : "2025-02-11T19:12:53Z",
  "subject" : "traits: nieuwe resultaten",
  "body" : {
    "text_html" : "<!doctype html><html xmlns=\"http://www.w3.org/1999/xhtml\" xmlns:o=\"urn:schemas-microsoft-com:office:office>"
  },
  "is_parsed" : true,
  "is_google_scholar_format" : true,
  "log_message" : "Body successfully parsed. "
},
  "is_processed" : true,
  "is_spam" : false,
  "log_message" : "Email read successfully."
}
```

Figuur 11.12: Databank structuur body.

- **_id:** De unieke ID van de record.

- **created_at**: De datum en het tijdstip waarop de e-mail in de databank werd opgeslaan.
- **updated_at**: De datum en het tijdstip waarop de e-mail in de databank het laatst werd gewijzigd.
- **sender**: De afzender van de e-mail.
- **date_time**: De datum en het tijdstip waarop de e-mail verstuurd werd.
- **subject**: Het onderwerp van de e-mail.
- **body**:
 - **text_html**: De body van de e-mail.
 - **is_parsed**: Een vlag die aangeeft of het parsen van de body succesvol was.
 - **is_google_scholar_format**: Een vlag die aangeeft of de body overeenkomt met de verwachte html structuur.
 - **log_message**: De logberichten tijdens het parsen van de body worden hier bewaard.
- **is_processed**: Een vlag die aangeeft of de body van de e-mail verwerkt werd.
- **is_spam**: Een vlag die aangeeft of de e-mail niet van GS afkomstig is.
- **log_message**: De logberichten tijdens het lezen van de e-mail worden hier bewaard.

Een overzicht van de gegevens die opgeslaan worden zijn te zien in figuur 11.13.

```
{
  "_id" : ObjectId("67b5cfa84c19ccd3d1df01d0"),
  "created_at" : "2025-02-19T12:33:44Z",
  "updated_at" : "2025-02-19T12:33:44Z",
  "email" : ObjectId("67b5cf80ba5d6fa0d55d98eb"),
  "title" : "Genetic Variability and Association of Morpho-Agronomic Traits Among Ethiopian Barley (Hordeum vu",
  "author" : "A Zewodu, W Mohammed, E Shiferaw ",
  "publisher" : " Scientifica",
  "year" : " 2025",
  "text" : "... traits among 49 barley accessions. The experiment was conducted in 2021 using \r\na 7 x 7 simple",
  "link" : {
    "url" : "https://scholar.google.be/scholar_url?url=https://onlinelibrary.wiley.com/doi/pdf/10.1155/sci5/"
  },
  "log_message" : "Search result parsed successfully.",
  "is_processed" : false
}
```

Figuur 11.13: Databank structuur zoekresultaat.

- **_id**: De unieke ID van de record.
- **created_at**: De datum en het tijdstip waarop het zoekresultaat in de databank werd opgeslaan.
- **updated_at**: De datum en het tijdstip waarop het zoekresultaat in de databank het laatst werd gewijzigd.

- **email:** Referentie naar de overeenkomstige e-mail waar dit zoekresultaat in staat.
- **title:** Titel van het zoekresultaat.
- **author:** Auteur(s) van het zoekresultaat.
- **publisher:** Uitgever van het zoekresultaat.
- **year:** Jaar van het zoekresultaat.
- **text:** Snippet van het zoekresultaat.
- **link:**
 - **url:** De link naar de online bron van het zoekresultaat.
- **log_message:** De logberichten tijdens het verwerken van het zoekresultaat worden hier bewaard.
- **is_processed:** Een vlag die aanduidt of de online bron van het zoekresultaat verwerkt werd.

11.4.3. DOI opzoeken in de link van het geparseerde resultaat

Het DOI, het resultaat van deze stap, is een unieke identifier voor elk artikel. Vanaf dit punt en voor alle volgende bewerkingen kan elk artikel ondubbelzinnig geïdentificeerd worden en kunnen de bijhorende gegevens op gezocht worden. Een overzicht van de gegevens die opgeslaan worden is te zien in figuur 11.14.

```
{
  "_id" : ObjectId("67b8621c2adedfb8d1920e85"),
  "created_at" : "2025-02-21T11:23:08Z",
  "updated_at" : "2025-02-21T11:23:31Z",
  "email" : ObjectId("67b85f7e139ca121611fcf11"),
  "title" : "Natural selection and adaptive traits in the Maniq, a nomadic hunter-gatherer society from Mainland Southeast Asia",
  "author" : "T Herzog, M Larena, W Kutanan, H Lukas, M Fieder...",
  "publisher" : "Scientific Reports",
  "year" : "2025",
  "text" : "Asia is home to diverse hunter-gatherer populations characterized by significant \r\nmorphological, anthropological,
  \"link\" : {
    \"url\" : \"https://scholar.google.be/scholar_url?url=https://www.nature.com/articles/s41598-024-83657-0&hl=nl&sa=X&d=1146959\",
    \"response_code\" : NumberInt(200),
    \"response_type\" : \"text/html\",
    \"is_accepted_type\" : true,
    \"DOI\" : \"10.1002/humu.22595\",
    \"log_message\" : \"DOI successfully retrieved\",
    \"is_Doi_success\" : true
  },
  \"log_message\" : \"Search result parsed successfully.\",
  \"is_processed\" : true
}
```

Figuur 11.14: Databank structuur link.

- **_id:** De unieke ID van de record.
- **created_at:** De datum en het tijdstip waarop het zoekresultaat in de databank werd opgeslaan.
- **updated_at:** De datum en het tijdstip waarop het zoekresultaat in de databank het laatst werd gewijzigd.

- **email:** Referentie naar de overeenkomstige e-mail waar dit zoekresultaat in staat.
- **title:** Titel van het zoekresultaat.
- **author:** Auteur(s) van het zoekresultaat.
- **publisher:** Uitgever van het zoekresultaat.
- **year:** Jaar van het zoekresultaat.
- **text:** Snippet van het zoekresultaat.
- **link:**
 - **url:** De link naar de online bron van het zoekresultaat.
 - **location_replace_url:** De link naar de eigenlijke online bron van het zoekresultaat na redirect door GS.
 - **response_code:** De http status code van het response.
 - **response_type:** Het type document van het response.
 - **is_accepted_type:** Een vlag die aangeeft of het type van het response verder beschouwd zal worden.
 - **DOI:** De Digital Object Identifier.
 - **log_message:** De logberichten tijdens het verwerken van de link worden hier bewaard.
 - **is_DOI_succes:** Een vlag die aangeeft of de verwerking tijdens deze stap al dan niet succesvol was.
- **media_type:** Het type van de online bron zoals het aangegeven wordt in het zoekresultaat.
- **log_message:** De logberichten tijdens het verwerken van het zoekresultaat worden hier bewaard.
- **is_processed:** Een vlag die aangeeft of de online bron van het zoekresultaat verwerkt werd.

12

Crossref

12.1. Doelstelling

De voorgaande stappen leverden gegevens zoals titel en auteur, en ook de DOI van de publicatie op. Tot nu toe zijn al die resultaten gebaseerd op tekstverwerking en linked data. Aan de hand van Crossref kan er gevalideerd worden of de gevonden gegevens correct zijn. Dat is cruciaal voor de kwaliteit van de nieuwe publicaties ten opzichte van IMIS.

12.2. Analyse

12.3. Implementatie

De API's van Crossref worden ondersteund door meerdere onafhankelijke Python bibliotheken:

- Crossref Commons for Python ("Crossref Commons for Python", [2025](#))
- Habanero ("Habanero", [2025](#))
- Crossrefapi ("Crossref API", [2025](#))

Met als doel om de metadata van een publicatie op te vragen op basis van de DOI zijn alle bibliotheken aan elkaar gewaagd. Zonder bijzondere reden, behalve dat Crossref Commons ontwikkeld wordt door Crossref zelf, wordt er met Crossref Commons for Python gewerkt. De code om een entity op te vragen is bijzonder compact zoals te zien in [12.1](#)

12.4. Resultaat

Een overzicht van de gegevens die opgeslaan worden zijn te zien in figuur [12.1](#)

Use Case Crossref.	
primary actor	-
stakeholders	-
pre-condities	systeem beschikt over een DOI
post-condities	de gegevens van de publicatie zijn gevalideerd
normaal verloop	<ol style="list-style-type: none"> 1. systeem stuurt een request naar Crossref voor de betreffende DOI 2. systeem ontvangt response 3. systeem zoekt titel, auteur, uitgever, jaar in de response 4. systeem valideert deze gegevens met de informatie uit de voorgaande stappen 5. systeem slaat de gegevens op
alternatief verloop	<ol style="list-style-type: none"> 2A. geen response 2A1. het systeem logt het probleem 2A2. use case eindigt zonder bereiken van de post-conditie 2B. DOI is niet gevonden 2B1. het systeem logt het probleem 2B2. het systeem kan niet valideren 2B3. het systeem slaat dit slechte resultaat op

Tabel 12.1: Crossref.

```

1 response =
  crossref_commons.retrieval.get_publication_as_json(doi)

```

Codefragment 12.1: Codefragment voor opvragen van de metadata van een publicatie aan Crossref.

Record	Value	Aggregates
_id	67d055d19aa65d914744acab	
api_url		
author	Yugang Guo	
created_at	2025-03-11T15:25:05Z	
doi_url	https://doi.org/10.1038/s41598-025-93038-w	
log_message	Crossref retrieved successfully.	
publisher	Springer Science and Business Media LLC	
search_result	67d05554710151a213dde092	
title	Gender differences in the relationship between big five personality traits and aggression	
updated_at	2025-03-11T15:25:05Z	
year	2025	

Figuur 12.1: Databank structuur Crossref.

- **_id:** De unieke ID van de record.
- **created_at:** De datum en het tijdstip waarop de Crossref gegevens in de databank werden opgeslaan.
- **updated_at:** De datum en het tijdstip waarop de Crossref gegevens in de databank het laatst werden gewijzigd.
- **search_result:** Referentie naar de overeenkomstige search result waar deze Crossref gegevens op slaan.

- **title:** De titel volgens Crossref.
- **author:** De auteurs volgens Crossref.
- **publisher:** De uitgever volgens Crossref.
- **year:** Het jaartal volgens Crossref.
- **subject:** Het onderwerp van de e-mail.
- **doi_url:** Url van de DOI.
- **api_url:** Endpoint van de API voor deze DOI.
- **log_message:** De logberichten tijdens het raadplegen van Crossref worden hier bewaard.

13

Natural language processing

13.1. Doelstelling

LLMs zijn gemaakt met als doel om tekst te produceren en om tekst te begrijpen net zoals mensen dat doen. Ze zijn toepasbaar voor uiteenlopende doeleinden zoals het afleiden van de juiste context en het beantwoorden van vragen. Ze kunnen tekst samenvatten of zelf produceren. Titel, originele link, auteurs, naam van de uitgever, jaar van publicatie en snippet zijn vormen van context in een SERP. LLMs komen in verschillende geuren en kleuren, hier enkele verschillende setups. Er zijn betalende modellen. Er zijn niet betalende modellen. Zelf niet betalende modellen draaien met Ollama. Doen ze de job? Leiden ze de context af uit de SERP? Doen de verschillende setups het even goed? Wat zijn de gebreken? En wat kiezen we voor onze klant?

Kan NLP gebruikt worden om de Google Scholar SERP om te zetten in gestructureerde data? Machine learning modellen zijn sterk in het herkennen van patronen. Een Google Scholar SERP heeft ook een vast patroon:

- titel
- originele link
- auteurs
- naam van de uitgever
- jaar van publicatie
- snippet

Hoe goed zijn de LLMs in het parsen van rauwe HTML, en in het bijzonder van Google Scholar SERP? Die pagina heeft al een meer complexe structuur, cryptische CSS klassen en een heleboel verbose data in de HTML.

Big Five personality traits and university students' academic performance: A meta-analysis

S. Chen, AOK Cheung, Z. Zeng - *Personality and Individual Differences*, 2025
Meta-analytical evidence on the linkage between university students' academic performance and personality **traits**, considering cultural and contextual factors, remains limited. This meta-analysis reports comprehensive assessment of the ...



Measuring personality Traits: Simulations Exploring (Mis) Alignment between methods and measurands

P. Durkee - *Journal of Research in Personality*, 2025
... Measuring personality **traits** requires that measurement methods align with intended measurands. There are two broad classes of personality **trait** ... The simulations highlight challenges in capturing personality **trait** measurands and offer ...



Landscape modification and species traits shape seasonal wildlife community dynamics within an arid metropolitan region

J.D. Haight, S.J. Hall, J.S. Lewis - *Landscape and Urban Planning*, 2025
... For this study, we evaluated how landscape characteristics and species **traits** influenced habitat use (i.e. patch occupancy, persistence, and ...). However, habitat use further varied according to species **traits**, with larger-bodied species exhibiting ...



[HTML] Consequences of Dietary Olive Leaf Powder Supplementation on Growth Performance, Carcass Traits, Blood Biochemical Parameters and Gut Microbiota in Broilers

M.H. Nagm, A.K. Aldhalmi, E.A. Ashour, L.A. Mohamed... - *Poultry Science*, 2025
This experiment examined the potential of olive leaf powder (OLP) as a natural growth enhancer in broiler feed and its influences on growth performance, carcass characteristics, blood parameters, and intestinal bacterial count. A total of 210 one-day-old ...



[HTML] ...-scale integration of meta-QTL and genome-wide association study identifies genomic regions and candidate genes for photosynthetic efficiency traits in bread wheat

M. Chen, T. Chen, L. Yun, Z. Che, J. Ma, B. Kong, J. Long... - *BMC Genomics*, 2025
... efficiency, a complex **trait**, through the ... in **trait** coverage, demonstrating significant advantages in data size and **trait** resolution dimensions. The increasing and intensive application of MQTL has effectively revealed the complex genetic ...



[HTML] Comprehensive evaluation of sugar beet varieties based on genotype× yield× trait (GYT) in different environments

D. Xinwang, L. yuhang, H. Xiaohang, L. yanli - *Scientific Reports*, 2025
... The aim of this study was to assess the agronomic **traits** and genotype by yield × **trait** (GYT) interactions of sugar beet across ... **trait** combinations was analyzed using the GYT biplot technique. Principal components PC1 and PC2 accounted for ...



[HTML] A Unified Method for Detecting Phylogenetic Signals in Continuous, Discrete, and Multiple Trait Combinations

L. Yan, Y. Yuan - *Ecology and Evolution*, 2025
... We used d phylo to represent phylogenetic distance and d **trait** to represent **trait** distance, and d **trait** was calculated using Gower's distance, which allows our method to handle any type of **trait**. Considering three species i, j, and k, if d [i, j] phylo < ...



... sequence repeat (SSR) markers for Chinese yam (*Dioscorea opposita* Thunb.) and genetic diversity and association analysis using phenotypic and genotypic traits

Q. Zhang, H. Tian, S. Guo, Y. Wang, S. Pei, C. Wu, Q. Wang... - *Genetic Resources and ...*, 2025
Chinese yam (*Dioscorea opposita* Thunb.) has a long history of cultivation and application, and is an important medicinal and edible plant in China. However, the lack of molecular markers limits the efficiency and accuracy of the genetic breeding ...



[HTML] Aboveground-canopy and belowground-root-trait correlations contribute

```
1     messages=[  
2     {"role": "system",  
3         "content": "You are a master at scraping Google Scholar  
         results data. Scrape top 10 organic results data  
         from Google Scholar search result page."},  
4     {"role": "user", "content": body_text}  
5     ],
```

Codefragment 13.1: Codefragment voor het opstellen van een prompt.

Een betalend model van OpenAI wordt gebruikt en een niet betalend model van Anthropic.

13.2. Analyse

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

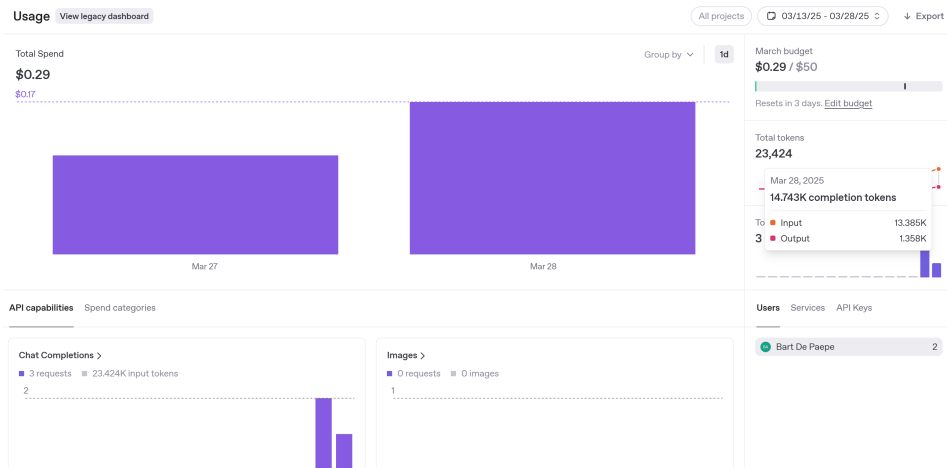
Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

13.3. Implementatie

De rauwe HTML wordt eerst lichtjes verwerkt om de <head>, <script> en <style> tags uit de HTML te verwijderen. Deze bevatten toch geen content en zouden alleen maar de rekentijd, het aantal tokens, en overeenkomstig de prijs doen toenemen. Er wordt een prompt opgesteld: Er worden parse opties opgesteld:

```
1      "function": {
2          "name": "parse_data",
3          "description": "Parse organic results from Google Scholar
4              SERP raw HTML data nicely",
5          "parameters": {
6              'type': 'object',
7              'properties': {
8                  'data': {
9                      'type': 'array',
10                     'items': {
11                         'type': 'object',
12                         'properties': {
13                             'title': {'type': 'string'},
14                             'original_url': {'type': 'string'},
15                             'authors': {'type': 'string'},
16                             'year_of_publication': {'type':
17                                 'integer'},
18                             'journal_name': {'type': 'string'},
19                             'snippet': {'type': 'string'}
20                         }
21                     }
22                 }
23             }
24         }
```

Codefragment 13.2: Codefragment voor het opstellen van parse opties.



Figuur 13.2: OpenAI dashboard.

Error code: 529 - 'type': 'error', 'error': 'type': 'overloaded', 'message' : 'Overloaded'

Codefragment 13.3: Codefragment Anthropic overloaded.

13.4. Resultaat

Het betalend model van OpenAI heeft 100% succes. Het vindt titel, originele url, auteurs, naam van de uitgever, jaar van publicatie, en snippet voor elk van de 10 zoekresultaten in de SERP.

De kostprijs van de opzoeking bedraagt 0,17\$ (dat is voor 1 Google Scholar alert met 10 zoekresultaten). Het model gebruikte daarvoor 14743 tokens. Het niet betalend model van Anthropic heeft ook 100% success. Er moet wel vermeldt worden dat niet alle jobs slaagden. In sommige gevallen gaf de Anthropic API een overloaded error terug ("Anthropic: using the API: Errors", 2025). Om daar wat aan te doen kunnen we een model lokaal draaien met ollama. Onderstaande tabel geeft een overzicht van de zoekresultaten met OpenAI ten opzichte van de zoekresultaten met Anthropic. Om een model lokaal te draaien is de limiet bepaald door het geheugen van de computer. Daarom wordt er gekozen voor het llama3.2 model dat 2GB inneemt. Het llama3 model neemt 40GB in en kan niet gebruikt worden op een laptop. Het lokale llama3.2 model is niet in staat om de SERP te scrapen. Ook niet wanneer we slechts 1 zoekresultaat als input gegeven wordt. Het is pas nadat de verbose HTML tags verwijderd zijn, dat het model in staat is om de overgebleven tekst te parsen. Dit moet verder onderzocht worden aan de hand van grotere modellen op krachtigere servers.

OpenAI	Anthropic
Auteurs	-
L Schunck, U von Döring, J Wilke	["L Schunck", "U von Döring", "J Wilke"]
FZ Ben Debbane, A Baidani, M Aarbaoui	["FZ Ben Debbane", "A Baidani", "M Aarbaoui"]
C Brittenham, AS DiCriscio, V Troiani, Y Hu, JB Wagner	["C Brittenham", "AS DiCriscio", "V Troiani", "Y Hu", "JB Wagner"]
B Tator, M Le Sage, RA Relyea	["B Tator", "M Le Sage", "RA Relyea"]
S Kang, B Yoon, K Kim, J Gratch, W Woo	["S Kang", "B Yoon", "K Kim", "J Gratch", "W Woo"]
J Wei, X Tang, J Liu, T Luo, Y Wu, J Duan, S Xiao	["J Wei", "X Tang", "J Liu", "T Luo", "Y Wu", "J Duan", "S Xiao"]
M Rakszegi, V Tóth, S Tömösközi, E Jak-sics, A Cseh	["M Rakszegi", "V Tóth", "S Tömösközi", "E Jak-sics", "A Cseh"]
XY Zhang, WL Chen	["XY Zhang", "WL Chen"]
B Larue, F Pelletier, M Festa-Bianchet, S Hamel	["B Larue", "F Pelletier", "M Festa-Bianchet", "S Hamel"]
JW Whitlock, PM Orwin, ZR Stahlschmidt	["JW Whitlock", "PM Orwin", "ZR Stahlschmidt"]
Link	-
"url": "https://www.sciencedirect.com/science/article/pii/S0191386925001308"	"url": "https://www.sciencedirect.com/science/article/pii/S0191386925001308"
"url": "https://link.springer.com/article/10.1007/s42729-025-02366-3"	"url": "https://link.springer.com/article/10.1007/s42729-025-02366-3"
"url": "https://www.nature.com/articles/s41598-025-92904-x"	"url": "https://www.nature.com/articles/s41598-025-92904-x"
"url": "https://www.sciencedirect.com/science/article/pii/S0269749125004750"	"url": "https://www.sciencedirect.com/science/article/pii/S0269749125004750"
"url": "https://ieeexplore.ieee.org/abstract/document/109355921"	"url": "https://ieeexplore.ieee.org/abstract/document/109355921"
"url": "https://www.sciencedirect.com/science/article/pii/S2666154725002091"	"url": "https://www.sciencedirect.com/science/article/pii/S2666154725002091"
"url": "https://www.sciencedirect.com/science/article/pii/S073521025000608"	"url": "https://www.sciencedirect.com/science/article/pii/S073521025000608"
"url": "https://academic.oup.com/botlinnean/advance-article/doi/10.1093/botlinnean/boaf012/8089970"	"url": "https://academic.oup.com/botlinnean/advance-article/doi/10.1093/botlinnean/boaf012/8089970"
"url": "https://www.pnas.org/doi/full/10.1073/pnas.2417158122"	"url": "https://www.pnas.org/doi/full/10.1073/pnas.2417158122"
"url": "https://journals.biologists.com/jeb/article/doi/10.1246/jeb.250210/367452"	"url": "https://journals.biologists.com/jeb/article/doi/10.1246/jeb.250210/367452"
Publisher	-
Personality and Individual Differences	Personality and Individual Differences
Journal of Soil Science and	Journal of Soil Science and Plant Nutri-tion
Scientific Reports	Scientific Reports
Environmental Pollution	Environmental Pollution
IEEE Transactions on Visualization and	IEEE Transactions on Visualization and Computer Graphics
	Journal of Agriculture and Food Engi-



Onderzoeksvoorstel

Het onderwerp van deze bachelorproef is gebaseerd op een onderzoeksvoorstel dat vooraf werd beoordeeld door de promotor. Dat voorstel is opgenomen in deze bijlage.

A.1. Inleiding

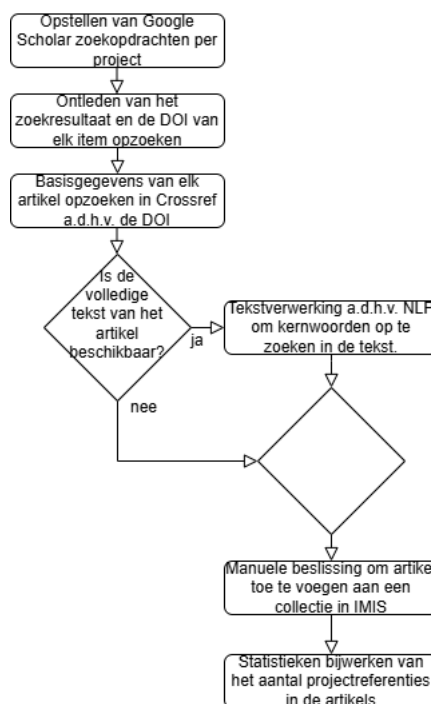
Het Vlaams Instituut voor de Zee (VLIZ) ("Vlaams Instituut voor de Zee", [2024](#)) is een pionier in zeekennis. Dit wetenschappelijk instituut gelegen in Oostende heeft onder andere een mandaat om een complete en geactualiseerde catalogus bij te houden van alle wetenschappelijke publicaties in de mariene sector. Al meer dan 20 jaar bouwt het Integrated Marine Information System (IMIS) aan deze catalogus die intussen beschikt over meerdere collecties van marien wetenschappelijk referentiemateriaal.

Naast wetenschappelijke literatuur zitten er ook collecties van mariene wetenschappelijke projecten in het systeem. Het is belangrijk voor het VLIZ om te weten door hoeveel publicaties er naar een project (vb. het World Register of Marine Species (WoRMS) ("World Register of Marine Species", [2024](#))) verwezen wordt. Deze maatstaf geeft een indicatie van het draagvlak van elk project binnen de wetenschappelijke gemeenschap en is één van de belangrijkste criteria tijdens projectevaluaties. Daarom is het cruciaal om continu nieuwe publicaties waarin verwezen wordt naar die projecten op te zoeken. Daarvoor wordt op heden Google Scholar gebruikt die bekend staat als de meest uitgebreide en geactualiseerde index. Op die manier blijft IMIS up-to-date en zijn de projectreferenties steeds geactualiseerd.

Momenteel verloopt dit proces binnen het VLIZ grotendeels handmatig. De zoekresultaten, volgens een bepaalde zoekfilter per project, afkomstig van Google Scholar worden manueel gefilterd en de Digital Object Identifier (DOI) van de geselecteerde artikels wordt opgezocht. Vervolgens worden de basisgegevens van elk arti-

kel zoals titel, auteurs, datum en uitgever opgevraagd op basis van de DOI in Crossref ("Crossref", 2024). Met deze informatie wordt manueel beslist om het artikel al dan niet toe te voegen aan een collectie binnen IMIS.

Dit is een tijdrovend proces. Daarom is er vraag naar automatisatie die de zoekresultaten verwerkt en gestructureerd opslaat. De beslissing om een artikel toe te voegen aan IMIS blijft nog altijd een manuele stap, maar naar verwachting moet dit sneller, efficiënter en accurater verlopen. Dat moet mogelijk zijn doordat de heterogene zoekresultaten omgezet worden in gestructureerde informatie.



Figuur A.1: Chronologische oplistijng van de uit te voeren stappen

De centrale vraag die opgelost moet worden is: 'Hoe kunnen de uiteenlopende zoekresultaten van Google Scholar automatisch verwerkt worden om de DOI van elk item op te zoeken?'. Daarbij moet ook rekening gehouden worden met duplicaten, relevantie en diverse media types voor elk item van het zoekresultaat.

Om bovenstaande vraag te verduidelijken, moeten er een paar deelvragen gesteld worden:

- De zoekresultaten zijn afhankelijk van de zoekopdracht. In welke mate beïnvloeden kleine variaties van de zoekopdracht de resultaten?
- De zoekresultaten worden standaard gesorteerd op relevantie door Google Scholar, maar komt dat overeen met onze verwachtingen?
- De zoekopdracht kan het aantal resultaten beperken naar wens. Wat is de optimale threshold?

Om bovenstaande vraag op te lossen, moeten er een paar deelvragen beantwoord worden:

- Op welke manier kan de Google Scholar zoekopdracht op regelmatige basis automatisch uitgevoerd worden?
- Hoe worden de zoekresultaten van Google Scholar het best geparsed?
- Zoals gesteld zijn de zoekresultaten divers van aard. Kan het resultaat mits kwaliteitscontrole verbeterd (gefilterd) worden?
- De DOI zelf ontbreekt in het resultaat. Hoe zal de DOI voor elk item opgezocht worden?

Verder moet er gekeken worden welke databanktechnologie het meest geschikt is om de verwerkte informatie van elk item op te slaan. Tenslotte moet er ook nagedacht worden hoe het manuele proces met gebruik van de gestructureerde informatie het best softwarematig ondersteund kan worden. De volledige pijplijn zal geïmplementeerd worden als een Proof of Concept (PoC) die kan getest worden op basis van echte zoekresultaten. Als vervolg op dit onderzoek kan dan overgegaan worden tot het aanpassen en uitrollen van het systeem in functie van de feedback.

A.2. Literatuurstudie

IMIS (Haspeslagh & Vanden Berghe, 2024) heeft als doel het mariene onderzoek in Vlaanderen te coördineren. Die rol wordt vervuld door op te treden als centraal kenniscentrum voor en door de mariene sector in Vlaanderen. De taken van IMIS gaan verder dan louter het uitbouwen van collecties met publicaties, ze omvatten ook referenties naar wetenschappers en naar projecten. In deze bachelorproef zijn de publicaties en de wetenschappelijke projecten belangrijk.

Google Scholar is een zoekmachine die wetenschappelijke literatuur indexeert ("Google Scholar, Uit Wikipedia, de vrije encyclopedie", 2024). Google Scholar laat toe om te gaan zoeken op sleutelwoorden, titel, auteur, domein en combinaties van deze (Noruzi, 2005). Op die manier zijn er enkel zoekresultaten die overeenkomen met onze zoekcriteria.

Er zijn verschillende kenmerken die bijdragen aan het succes van Google Scholar ten opzichte van de andere grote collecties. Google Scholar is beschikbaar zonder kosten, het bevat de grootste bibliografische collectie ter wereld, en de data is afkomstig van zowel publieke als niet publieke bronnen (Aguillo, 2011).

Google Scholar wordt meer en meer gebruikt als een bibliometrische tool om informatie te verzamelen over de impact van citaties van en naar afzonderlijke artikels (Moed e.a., 2016). Araújo e.a. (2021) legt uit hoe wetenschappelijke artikels andere publicaties citeren. Dat doen ze door middel van bibliografische referenties naar andere documenten doorheen de tekst, en een uitgewerkte lijst van referenties

aan het einde van het artikel. Een index van citaties zoals in Google Scholar, is een databank die deze referenties tussen documenten opslaat.

In dergelijk digitaal netwerk van referenties, moeten de afzonderlijke entiteiten blijvend, betrouwbaar en onderscheidbaar geïdentificeerd kunnen worden. Chandrakar (2006) legt uit hoe de DOI voldoet aan deze criteria en de standaard identificatie is geworden van intellectuele eigendom op het internet.

Aan de hand van de bekomen DOI op basis van de zoekresultaten kan alle informatie met betrekking tot het artikel opgevraagd worden. In sommige gevallen is zelfs een pdf-versie beschikbaar. Vanuit ons standpunt zijn pdf-versies van wetenschappelijke artikels een vorm van ongestructureerde informatie. Automatisch berekende voorstellingen van de inhoud kunnen bijdragen om het zoekproces te vereenvoudigen. Dat is ook het uitgangspunt van het werk van (Jehangir e.a., 2023). Natural language processing (NLP) is van groot belang voor het verwerken van deze ongestructureerde data. Met behulp van NLP kunnen er conclusies en samenvattingen gemaakt worden. Binnen het domein van NLP, staat Named Entity Recognition (NER) centraal voor het opstellen van structuur op basis van ongestructureerde tekst (Pakhale, 2023). Palshikar (z.d.) en Jehangir e.a. (2023) lijsten de verschillende NER technieken op die kunnen gebruikt worden. Koning e.a. (2005) tenslotte past NER toe om taxonomische entiteiten op te zoeken in een doorlopende tekst. Dit lijkt dusdanig ook geschikt voor het opzoeken van projectentiteiten in een artikel. Schäfer en Kiefer (2011) merkt verder ook terecht op dat wetenschappers vandaag geconfronteerd worden met een dagelijkse vloedgolf van nieuwe publicaties. Daarboven worden ook meer en meer oude artikels gedigitaliseerd. Zowel de verwerking van de platte tekst als de vloedgolf aan publicaties geven aanleiding tot het toepassen van NLP voor het beheren van wetenschappelijke artikels (Schäfer & Kiefer, 2011).

A.3. Methodologie

Ons onderzoek wordt opgesplitst in 7 afzonderlijke fases:

- opstellen van de meest accurate zoekopdrachten in Google Scholar
- parsen van de zoekresultaten om de DOI van elk artikel op te zoeken
- basisgegevens van elke publicatie opvragen in Crossref op basis van de DOI
- ontwikkelen van een tekstverwerkingsproces op basis van NER
- ontwerpen van de databank voor de opslag van de gestructureerde data
- installeren van een pijplijn die de afzonderlijke stappen combineert
- ontwerpen van een gebruikersinterface die de gestructureerde data toont

Achtereenvolgend wordt er nu dieper ingegaan op de verschillende fases. In de context van dit voorstel, wordt elke stap nog niet volledig uitgewerkt, maar is het de bedoeling om een richting aan te geven. Een uitvoerige beschrijving van alle taken zal aan bod komen in de bachelorproef.

De eerste fase gaat over het correct opstellen van de Google Scholar zoekopdrachten voor de verschillende projecten met als doel om als resultaat de meest relevante artikels te verkrijgen.

De tweede fase komt overeen met de hoofdpdracht van ons onderzoek, namelijk een systeem ontwikkelen dat de heterogene zoekresultaten kan parsen om de DOI op te zoeken. Meerdere zaken komen hier aan bod zoals ten eerste het aggregeren van zoekresultaten en groeperen per titel. Vervolgens het opzoeken van de DOI via de linked data, of via Crossref, of op de full-text versie van het artikel. Python lijkt de voor de hand liggende programmeertaal voor deze taak waarbij het draait om dataverwerking.

In de derde fase zal er met de API van Crossref gewerkt worden om de basisgegevens van een artikel op te vragen aan de hand van de DOI. Ook hier zal Python gebruikt worden.

De vierde fase onderzoekt of er met behulp van NER bruikbare informatie gefilterd kan worden uit de pdf-versie van een artikel. Op die manier kunnen de basisgegevens verrijkt worden om de manuele moderatie accurater te maken. Een zoektocht op het internet levert snel een achttal oplossingen op die ruwe tekst verwerken tot gestructureerde informatie:

- ScienceParse (“Science Parse parses scientific papers (in PDF form) and returns them in structured form.” [2024](#))
- ScienceParse API
 (“science-parse-api 1.0.1”, [2024](#))
- Grobid (“GROBID: a machine learning software for extracting information from scholarly documents”, [2024](#))
- Sci-pdf parser (“Python PDF parser for scientific publications: content and figures”, [2024](#))
- Paper AI (“Semantic search and workflows for medical/scientific papers”, [2024](#))
- Paper ETL (“ETL processes for medical and scientific papers”, [2024](#))
- Delph-in (“DELPH-IN Overview”, [2024](#))
- TaxonGrab (“TaxonGrab”, [2024](#))

Om onze specifieke business logica te implementeren, wordt gewerkt met één oplossing, of een combinatie van meerdere oplossingen, eventueel met gebruik van

een geschikte NER methode. Ook hier is Python de gekozen programmeertaal voor de dataverwerking.

De vijfde fase evalueert welke databanktechnologie het best geschikt is om onze gestructureerde data in op te slaan. Daar zijn bijvoorbeeld relationele, grafische en document store databanken.

De zesde fase focust zich op de automatisatie van het gehele proces: van zoekresultaat tot opgeslagen gestructureerde data. Er zal een pijplijn opgesteld worden die de voortgang van elke stap observeert en de status van het hele proces update. Dit wordt gebouwd met Symfony rekening houdend met de technology stack binnen het VLIZ.

De zevende en laatste fase tenslotte moet ervoor zorgen dat alle voorgaande stappen ook opbrengen voor de eindgebruiker. De opgeslagen gestructureerde data wordt ontsloten via een API (API Platform) en een custom front-end (NextJS) applicatie zorgt ervoor dat deze overzichtelijk voorgesteld wordt.

A.4. Verwacht resultaat, conclusie

Naar verwachting zal deze pijplijn de verwerking van de Google Scholar zoekresultaten versnellen. Op die manier draagt ons onderzoek bij aan de continue uitbreiding van de collecties van IMIS. Er wordt ook verwacht dat het aantal projectreferenties beter gemonitord zal worden. Zo helpt ons onderzoek bij het beter in kaart brengen van het draagvlak voor de wetenschappelijke projecten van het VLIZ.

Bibliografie

- Aguillo, I. F. (2011). Is Google Scholar useful for bibliometrics? A webometric analysis. *Scientometrics*, 91(2), 343–351. <https://doi.org/10.1007/s11192-011-0582-8>
- Amin, M. M., Sutrisman, A., & Dwitayanti, Y. (2024). Google Scholar Crawling for Constructing Research Database. In *Proceedings of the 7th FIRST 2023 International Conference on Global Innovations (FIRST-ESCSI 2023)* (pp. 331–337). Atlantis Press International BV. https://doi.org/10.2991/978-94-6463-386-3_36
- Anthropic: using the API: Errors*. (2025, april 1). <https://docs.anthropic.com/en/api/errors>
- Araújo, P. C. d., Gutierrez, R. C., & Hjørland, B. (2021). Citation Indexing and Indexes. *KNOWLEDGE ORGANIZATION*, 48(1), 72–101. <https://doi.org/10.5771/0943-7444-2021-1-72>
- Atlas Vector Search Quick Start*. (2025, mei 1). <https://www.mongodb.com/docs/atlas/atlas-vector-search/tutorials/vector-search-quick-start/>
- Beautiful Soup*. (2025, maart 1). <https://www.crummy.com/software/BeautifulSoup/>
- Beel, J., & Gipp, B. (2009). Google Scholar's Ranking Algorithm: An Introductory Overview. *12th Int. Conf. Scientometrics Informetrics (ISSI)*, 439–446.
- Bhatt, C., Bisht, A., Chauhan, R., Vishvakarma, A., Kumar, M., & Sharma, S. (2023). Web Scraping Techniques and Its Applications: A Review. *2023 3rd International Conference on Innovative Sustainable Computational Technologies (CISCT)*, 1–8. <https://doi.org/10.1109/cisct57197.2023.10351298>
- Castrillo-Fernández, O. (2015, december 1). *Web Scraping: Applications and Tools* (tech. rap.). ePSIplatform Topic Report No. 2015 / 10 , December 2015.
- Chandrakar, R. (2006). Digital object identifier system: an overview. *The Electronic Library*, 24(4), 445–452. <https://doi.org/10.1108/02640470610689151>
- Chroma*. (2025, mei 1). <https://docs.trychroma.com/docs/overview/getting-started>
- Click*. (2025, maart 1). <https://click.palletsprojects.com/en/stable/>
- Crossref*. (2024, november 1). <https://www.crossref.org/>
- Crossref API*. (2025, april 1). <https://github.com/fabiobatalha/crossrefapi>
- Crossref Commons for Python*. (2025, april 1). https://gitlab.com/crossref/crossref_commons_py
- DELPH-IN Overview*. (2024, november 1). <https://github.com/delph-in/docs/wiki>
- Dependency injection*. (2025, maart 1). https://en.wikipedia.org/wiki/Dependency_injection

- Dependency Injector*. (2025, maart 1). <https://python-dependency-injector.ets-labs.org/index.html>
- Document-oriented database*. (2025, maart 1). https://en.wikipedia.org/wiki/Document-oriented_database
- DOIs and matching regular expressions*. (2025, maart 1). <https://www.crossref.org/blog/does-and-matching-regular-expressions/>
- email - An email and MIME handling package*. (2025, maart 1). <https://docs.python.org/3.13/library/email.html>
- ETL processes for medical and scientific papers*. (2024, november 1). <https://github.com/neuml/paperetl>
- Extracting Structured Data from Websites with Mirascope, Anthropic, and Streamlit*. (2025, mei 1). <https://ploomber.io/blog/mirascope-url-extractor/>
- Gemma3: Powerful AI Within Everyone's Reach*. (2025, mei 1). <https://gemma3.org/>
- Google Scholar Guide*. (2025, maart 1). <https://library.acg.edu/how-to-guides/google-scholar/overview>
- Google Scholar, Uit Wikipedia, de vrije encyclopedie*. (2024, november 1). https://nl.wikipedia.org/wiki/Google_Scholar
- Gray, J., Bounegru, L., & Chambers, L. (2012, augustus). *The Data Journalism Handbook*.
- GROBID: a machine learning software for extracting information from scholarly documents*. (2024, november 1). <https://github.com/kermitt2/grobid>
- Habanero*. (2025, april 1). <https://github.com/sckott/habanero>
- Haspeslagh, J., & Vanden Berghe, E. (2024, november 1). *IMIS: Integrated Marine Information System*. Vlaams Instituut voor de Zee. https://www.researchgate.net/publication/33549102_IMIS_Integrated_Marine_Information_System
- Hendricks, G., Tkaczyk, D., Lin, J., & Feeney, P. (2020). Crossref: The sustainable source of community-owned scholarly metadata. *Quantitative Science Studies*, 1(1), 414–427. https://doi.org/10.1162/qss_a_00022
- How to Use Chroma to Build Your First Similarity Search*. (2025, mei 1). <https://iamajithkumar.medium.com/how-to-use-chroma-to-build-your-first-similarity-search-5c054bfd5add>
- imaplib - IMAP4 protocol client*. (2025, maart 1). <https://docs.python.org/3/library/imaplib.html>
- Internet Message Access Protocol, Uit Wikipedia, de vrije encyclopedie*. (2025, maart 1). https://en.wikipedia.org/wiki/Internet_Message_Access_Protocol
- Introducing Meta Llama 3: The most capable openly available LLM to date*. (2025, mei 1). <https://ai.meta.com/blog/meta-llama-3/>

- Jehangir, B., Radhakrishnan, S., & Agarwal, R. (2023). A survey on Named Entity Recognition — datasets, tools, and methodologies. *Natural Language Processing Journal*, 3, 100017. <https://doi.org/10.1016/j.nlp.2023.100017>
- Koning, D., Sarkar, I. N., & Moritz, T. (2005). TaxonGrab: Extracting Taxonomic Names From Text. *Biodiversity Informatics*, 2(0). <https://doi.org/10.17161/bi.v2i0.17>
- Kumar, S., & Roy, U. B. (2023). A technique of data collection. In *Statistical Modeling in Machine Learning* (pp. 23–36). Elsevier. <https://doi.org/10.1016/b978-0-323-91776-6.00011-7>
- LLM abstractions that aren't obstructions.* (2025, mei 1). <https://mirascope.com/>
- Lotfi, C., Srinivasan, S., Ertz, M., & Latrous, I. (2021). Web Scraping Techniques and Applications: A Literature Review. In *SCRS CONFERENCE PROCEEDINGS ON INTELLIGENT SYSTEMS* (pp. 381–394). Soft Computing Research Society. <https://doi.org/10.52458/978-93-91842-08-6-38>
- Lourenço, J. R., Cabral, B., Carreiro, P., Vieira, M., & Bernardino, J. (2015). Choosing the right NoSQL database for the job: a quality attribute evaluation. *Journal of Big Data*, 2(1). <https://doi.org/10.1186/s40537-015-0025-0>
- Mitchell, R. (2015). *Web Scraping with Python: Collecting Data from the Modern Web* (1st). O'Reilly Media, Inc.
- Moed, H. F., Bar-Ilan, J., & Halevi, G. (2016). A new methodology for comparing Google Scholar and Scopus. *Journal of Informetrics*, 10(2), 533–551. <https://doi.org/10.1016/j.joi.2016.04.017>
- MongoDB. (2025, maart 1). <https://www.mongodb.com/resources/basics/databases/document-databases>
- Natural Language Processing of Google Scholar search results.* (2025, mei 1). https://github.com/bart-de-paepe/google-scholar-natural-language-processing/blob/main/app/src/services/natural_language_processing_service.py
- Noruzi, A. (2005). Google Scholar: The New Generation of Citation Indexes. *Libri*, 55(4). <https://doi.org/10.1515/libr.2005.170>
- Oh, K., & Colón-Aguirre, M. (2019). A Comparative Study of Perceptions and Use of Google Scholar and Academic Library Discovery Systems. *College amp; Research Libraries*, 80(6), 876–891. <https://doi.org/10.5860/crl.80.6.876>
- Ollama: Get up and running with large language models.* (2025, mei 1). <https://github.com/ollama/ollama>
- OpenAI developer platform.* (2025, mei 1). <https://platform.openai.com/docs/overview>
- Pakhale, K. (2023). Comprehensive Overview of Named Entity Recognition: Models, Domain-Specific Applications and Challenges. <https://doi.org/10.48550/ARXIV.2309.14084>

- Palshikar, G. K. (z.d.). Techniques for Named Entity Recognition: A Survey. In *Bioinformatics* (pp. 400–426). IGI Global. <https://doi.org/10.4018/978-1-4666-3604-0.ch022>
- PRATIBA, D., M.S., A., DUA, A., SHANBHAG, G. K., BHANDARI, N., & SINGH, U. (2018). Web Scraping And Data Acquisition Using Google Scholar. *2018 3rd International Conference on Computational Systems and Information Technology for Sustainable Solutions (CSITSS)*, 277–281. <https://doi.org/10.1109/csitss.2018.8768777>
- PyMongo. (2025, maart 1). <https://pymongo.readthedocs.io/en/stable/>
- PyMuPDF is a high-performance Python library for data extraction, analysis, conversion manipulation of PDF (and other) documents. (2025, maart 1). <https://pymupdf.readthedocs.io/en/latest/how-to-open-a-file.html#opening-remote-files>
- Python PDF parser for scientific publications: content and figures. (2024, november 1). https://github.com/titipata/scipdf_parser
- Rafsanjani, M. R. (2022). ScrapPaper: A web scrapping method to extract journal information from PubMed and Google Scholar search result using Python. <https://doi.org/10.1101/2022.03.08.483427>
- Rahmatulloh, A., & Gunawan, R. (2020). Web Scraping with HTML DOM Method for Data Collection of Scientific Articles from Google Scholar. *Indonesian Journal of Information Systems*, 2(2), 95–104. <https://doi.org/10.24002/ijis.v2i2.3029>
- Requests is a simple, yet elegant, HTTP library. (2025, maart 1). <https://pypi.org/project/requests/>
- Sataloff, R. T., Bush, M. L., Chandra, R., Chepeha, D., Rotenberg, B., Fisher, E. W., Goldenberg, D., Hanna, E. Y., Kerschner, J. E., Kraus, D. H., Krouse, J. H., Li, D., Link, M., Lustig, L. R., Selesnick, S. H., Sindwani, R., Smith, R. J., Tysome, J. R., Weber, P. C., & Welling, D. B. (2021). Systematic and other reviews: Criteria and complexities. *World Journal of Otorhinolaryngology - Head and Neck Surgery*, 7(3), 236–239. <https://doi.org/10.1016/j.wjorl.2021.04.007>
- Schäfer, U., & Kiefer, B. (2011). Advances in Deep Parsing of Scholarly Paper Content. In *Advanced Language Technologies for Digital Libraries* (pp. 135–153). Springer Berlin Heidelberg. https://doi.org/10.1007/978-3-642-23160-5_9
- Science Parse parses scientific papers (in PDF form) and returns them in structured form. (2024, november 1). <https://github.com/allenai/science-parse>
- science-parse-api 1.0.1. (2024, november 1). <https://pypi.org/project/science-parse-api>
- Scrape Google and other search engines from our fast, easy, and complete API. (2025, mei 1). <https://serpapi.com/google-scholar-api>

- Scrapegraph AI Tutorial; Scrape websites easily with LLaMA AI.* (2025, mei 1). <https://www.scrapingbee.com/blog/scrapegraph-ai-tutorial-scrape-websites-easily-with-llama-ai/>
- ScrapeGraphAI: You Only Scrape Once.* (2025, mei 1). <https://github.com/ScrapeGraphAI/Scrapegraph-ai>
- Selenium with Python.* (2025, maart 1). <https://selenium-python.readthedocs.io/>
- Semantic search and workflows for medical/scientific papers.* (2024, november 1). <https://github.com/neuml/paperai>
- Semantic search using Chroma.* (2025, mei 1). <https://github.com/bart-de-paepe/google-scholar-semantic-chroma>
- Semantic search using MongoDB.* (2025, mei 1). <https://github.com/bart-de-paepe/google-scholar-semantic-mongodb>
- Singrodia, V., Mitra, A., & Paul, S. (2019). A Review on Web Scrapping and its Applications. *2019 International Conference on Computer Communication and Informatics (ICCCI)*. <https://doi.org/10.1109/iccci.2019.8821809>
- Status of Python versions.* (2025, maart 1). <https://devguide.python.org/versions/>
- Sulistya, Y. I., Wardhana, A. C., Istighosah, M., & Riyandi, A. (2024). Automated Google Scholar Crawling with a Web-Based Tool for Publication Data Management. *Jurnal Teknologi Dan Sistem Informasi Bisnis*, 6(4), 768–773. <https://doi.org/10.47233/jteksis.v6i4.1604>
- TaxonGrab.* (2024, november 1). <https://sourceforge.net/projects/taxongrab/>
- Tips for using the Crossref REST API.* (2025, mei 1). <https://www.crossref.org/documentation/retrieve-metadata/rest-api/tips-for-using-the-crossref-rest-api/>
- Using the API Getting started.* (2025, mei 1). <https://docs.anthropic.com/en/api/getting-started>
- Vector Search.* (2025, mei 1). <https://www.mongodb.com/products/platform/atlas-vector-search>
- Vlaams Instituut voor de Zee.* (2024, november 1). <https://www.vliz.be>
- Web scraping experiment with AI (Parsing HTML with GPT-4 and GPT-4o).* (2025, mei 1). <https://serpapi.com/blog/web-scraping-and-parsing-experiment-with-ai-openai/>
- Web scraping Google Scholar SERP using Anthropic.* (2025, mei 1). https://github.com/bart-de-paepe/google-scholar-anthropic/blob/main/app/src/services/parse_service.py
- Web scraping Google Scholar SERP using BeautifulSoup.* (2025, mei 1). https://github.com/bart-de-paepe/google-scholar-beautifulsoup/blob/main/app/src/services/parse_service.py

- Web scraping Google Scholar SERP using Llama3.* (2025, mei 1). https://github.com/bart-de-paepe/google-scholar-ollama/blob/main/app/src/services/parse_service.py
- Web scraping Google Scholar SERP using OpenAI.* (2025, mei 1). https://github.com/bart-de-paepe/google-scholar-openai/blob/main/app/src/services/parse_service.py
- World Register of Marine Species.* (2024, november 1). <https://www.marinespecies.org>
- Yang, D., Zhang, A. N., & Yan, W. (2017). Performing literature review using text mining, Part I: Retrieving technology infrastructure using Google Scholar and APIs. *2017 IEEE International Conference on Big Data (Big Data)*, 3290–3296. <https://doi.org/10.1109/bigdata.2017.8258313>