

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

---

**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

Aan de basis van deze bachelorproef ligt een JIRA ticket dat al een tijdje oud is. Dat ticket was aangemaakt door het team van IMIS met de bedoeling om een bestaande procedure te verbeteren, doch zonder daar enige prioriteit aan te koppelen. Bij mijn zoektocht naar een onderwerp kwam ik al snel bij mijn werkgever, het VLIZ, terecht. JIRA werd erbij gehaald en het ticket waarvan sprake stak met kop en schouders boven andere onderwerpen uit omwille van de toepasbaarheid, de mate van uitdaging en de haalbaarheid.

Wat volgde was een boeiende verdieping in de wereld van de academische literatuur, in het gigantische net van de web scraping, en niet in het minst van het onmetelijke universum van de LLMs.

Ik wil eerst en vooral mijn manager Bart Vanhoorne bedanken die een omgeving creëert waarin het mogelijk is om aan een bachelorproef te werken.

Bijzondere dank ook voor mijn collega Milan Lamote voor zijn aanstekelijke positiviteit en omdat hij zonder aarzelen het co-promotorschap van deze bachelorproef aanvaardde.

Bijzondere dank ook voor mijn promotor Jan Claes voor het opzetten van een transparant kader voor deze bachelorproef en voor zijn kritische feedback.

Oprechte dank voor mijn collega Fons Verheyde voor zijn enthousiasme bij het bespreken van deze bachelorproef.

Oprechte dank voor mijn collega Cedric Decruw voor zijn inspiratie over LLMs.

Tenslotte aan het einde van deze bachelorproef, maar vooral aan het einde van deze opleiding, eeuwige dank aan mijn echtgenote Lies Knockaert. Veel van de resultaten in mijn studies zijn onrechtstreeks ook haar verdienste.

Bart De Paepe,  
Sint-Baafs-Vijve, 1 mei 2025

# Samenvatting

Zoekresultaten afkomstig van Google Scholar moeten aan het Integrated Marine Information System (IMIS) toegevoegd worden. Dit onderzoek bekijkt hoe dit proces volledig of grotendeels geautomatiseerd kan worden. Daarbij komen verschillende aspecten kijken die allemaal stuk voor stuk afzonderlijk behandeld worden.

Zo worden Google Scholar alerts gebruikt om continu nieuwe zoekresultaten te ontvangen. Er wordt uitvoerig toegelicht hoe een zoekopdracht aangemaakt moet worden en hoe daarvoor een melding ingesteld kan worden.

Google Scholar zoekresultaten komen in de vorm van een HTML pagina. HTML aan IMIS toevoegen is niet interessant. Er wordt dieper ingegaan hoe de lijst met zoekresultaten van de HTML pagina gescraped kan worden. Daarbij worden verschillende technieken uitgeprobeerd waaronder Large Language Models (LLMs). Uiteindelijk wordt een klassieke benadering gekozen die de HTML parset aan de hand van BeautifulSoup.

Na het scrapen zou de gevonden informatie in IMIS opgeslagen kunnen worden. Er zijn echter nog bijkomende stappen nodig om ervoor te zorgen dat enkel kwalitatieve gegevens naar IMIS vloeien. Behalve het feit dat Google Scholar de zoekresultaten selecteert voor de zoekopdracht en rangschikt volgens een bepaalde volgorde, zijn er verder geen criteria die aangeven hoe relevant een zoekresultaat is voor IMIS. Natural Language Processing (NLP) biedt daar oplossingen voor. Er wordt een relevantiescore berekend van de mate waarin de zoekopdracht aanwezig is in het zoekresultaat. Algemeen wordt daarvoor aangenomen dat de frequentie van de zoekopdracht in het zoekresultaat recht evenredig is met de relevantie van het resultaat voor IMIS.

Er mogen ook geen duplicaten in IMIS opgeslagen worden. Duplicaten kunnen gedetecteerd worden aan de hand van de Digital Object Identifier (DOI) van de publicatie. Die kan stapsgewijs opgezocht worden in de link van de publicatie, in Crossref, of op de webpagina van de publicatie. Maar er is geen garantie dat de DOI gevonden zal worden. Voor die gevallen moet er minstens een score gegeven worden of een publicatie al dan niet een duplicaat is. Dat kan gedaan worden aan de hand van "Semantic search". Daarvoor worden embeddings berekend voor alle titels in IMIS. Vervolgens wordt de embedding van de titel van het zoekresultaat daarmee vergeleken. De mate van gelijkenis is een score voor duplicaten van de publicatie.

Uiteindelijk worden al deze stappen geïntegreerd in een pijplijn voor het semi-automatisch toevoegen van publicaties aan IMIS. In het geval de relevantiescore voldoende is

en indien de DOI gevonden wordt, kan de publicatie zonder tussenkomst aan IMIS toegevoegd worden. Wanneer er geen DOI gevonden wordt, is er nog steeds een manuele stap nodig om te beslissen op basis van de duplicatenscore of de publicatie toegevoegd mag worden.

# Inhoudsopgave

<b>Lijst van figuren</b>	<b>viii</b>
<b>Lijst van tabellen</b>	<b>ix</b>
<b>Lijst van codefragmenten</b>	<b>x</b>
<b>1 Inleiding</b>	<b>1</b>
1.1 Probleemstelling . . . . .	1
1.2 Onderzoeksvraag . . . . .	2
1.3 Onderzoeksdoelstelling . . . . .	2
1.4 Opzet van deze bachelorproef . . . . .	3
<b>2 Stand van zaken</b>	<b>4</b>
2.1 IMIS . . . . .	4
2.2 Google Scholar . . . . .	4
2.3 Web scraping . . . . .	4
2.4 Natural Language Processing . . . . .	6
2.5 Linked data . . . . .	6
2.6 Semantic Search . . . . .	7
<b>3 Methodologie</b>	<b>9</b>
<b>4 Google Scholar alert</b>	<b>12</b>
4.1 Instellen van de zoekopdracht & aanmaken van het alert . . . . .	12
<b>5 Web scraping</b>	<b>14</b>
5.1 Inleiding . . . . .	14
5.2 Web scraping met gebruik van een LLM . . . . .	16
5.2.1 Web scraping met OpenAI . . . . .	16
5.2.2 Web scraping met Mirascope en Anthropic . . . . .	18
5.2.3 Web scraping met een lokaal model . . . . .	20
5.3 Web scraping via het parsen van de DOM . . . . .	22
5.4 Meest geschikte Web scraping . . . . .	24
<b>6 Natural Language Processing</b>	<b>26</b>
6.1 Inleiding . . . . .	26
6.2 Relevantiescore . . . . .	26
6.3 Tekstverwerking . . . . .	28

<b>7</b>	<b>Linked data</b>	<b>30</b>
7.1	Inleiding . . . . .	30
7.2	De DOI opzoeken in de URL . . . . .	30
7.3	De DOI opzoeken in Crossref . . . . .	31
7.4	De DOI opzoeken op de webpagina van de publicatie . . . . .	33
<b>8</b>	<b>Semantic search</b>	<b>36</b>
8.1	Inleiding . . . . .	36
8.2	Chroma . . . . .	37
8.3	MongoDB Atlas . . . . .	39
<b>9</b>	<b>Conclusie</b>	<b>43</b>
<b>A</b>	<b>Google Scholar zoekopdracht</b>	<b>45</b>
A.1	Zoeken in Google Scholar ("Google Scholar Guide", 2025). . . . .	45
A.1.1	Basis zoeken . . . . .	45
A.1.2	Geavanceerd zoeken . . . . .	47
A.2	E-mail alerts . . . . .	48
<b>B</b>	<b>Onderzoeksvoorstel</b>	<b>50</b>
B.1	Inleiding . . . . .	50
B.2	Literatuurstudie . . . . .	52
B.3	Methodologie . . . . .	53
B.4	Verwacht resultaat, conclusie . . . . .	55
	<b>Bibliografie</b>	<b>56</b>

# Lijst van figuren

2.1	Illustratie van embeddings. . . . .	7
5.1	Google Scholar SERP. . . . .	15
5.2	OpenAI dashboard. . . . .	19
5.3	HTML structuur van de GS alert. . . . .	24
6.1	Bag of Words. . . . .	27
6.2	term frequency-inverse document frequency. . . . .	27
8.1	Semantic search. . . . .	36
8.2	Chroma resultaat. . . . .	39
8.3	Mongodb resultaat. . . . .	40
A.1	Google Scholar basis zoeken. . . . .	46
A.2	Google Scholar zoekresultaten. . . . .	46
A.3	Google Scholar geavanceerd zoeken. . . . .	48
A.4	Google Scholar melding maken. . . . .	49
A.5	Google Scholar email alert. . . . .	49
B.1	Chronologische oplijsting van de uit te voeren stappen . . . . .	51



# Lijst van tabellen

3.1	HTML scraping technieken. . . . .	10
4.1	VLIZ zoektermen. . . . .	13
6.1	Voor en na coreference resolution. . . . .	28
7.1	Reguliere expressies om een DOI te matchen. . . . .	31

# Lijst van codefragmenten

5.1	Prompt HTML fragment	16
5.2	HTML cleaning	17
5.3	OpenAI prompt	17
5.4	OpenAI parameters	18
5.5	Pydantic field	19
5.6	Pydantic schema	20
5.7	Pydantic model	20
5.8	Beautiful Soup HTML parser	21
5.9	Ollama met Docker	21
5.10	Scrapegraph configuratie	22
5.11	Scrapegraph Prompt	23
5.12	Beautiful Soup tag selection	25
5.13	Beautiful Soup navigatie	25
6.1	term frequency-inverse document frequency	27
6.2	Relevantiescore	28
6.3	Coreference resolver	29
7.1	URL met DOI	30
7.2	Opzoeken van DOI in de URL van de publicatie	32
7.3	Crossref query op basis van de bibliografie	32
7.4	Crossref commons	32
7.5	Pymupdf	34
7.6	Selenium	35
8.1	Connect Chroma	37
8.2	Chroma embeddings	38
8.3	Query Chroma	38
8.4	MongoDB set model	39
8.5	MongoDB embeddings	40
8.6	MongoDB calculate index	41
8.7	MongoDB query	42

# 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 dat 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 opgeslagen 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. Hoewel 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 wordt en waarbij slechts 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 opgeslagen 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 6 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

### 2.1. IMIS

Het Integrated Marine Information System (IMIS) (Haspeslagh & Vanden Berghe, [2002](#)) 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.

### 2.2. Google Scholar

Google Scholar (GS) is een online index van academische literatuur. 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 citaties en h-indexen <sup>1</sup>. GS laat toe om te gaan zoeken op sleutelwoorden, titel, auteur, domein en combinaties van deze (Noruzi, [2005](#)).

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

### 2.3. Web scraping

Web scraping of web crawling is een verzamelnaam voor een groep technieken waarmee gegevens onttrokken kunnen worden van webpagina's (Bhatt e.a., [2023](#)).

---

<sup>1</sup>De 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. vb. Een h-index van 3 betekent dat de auteur minstens 3 publicaties heeft die elk minstens 3 keer geciteerd werden.

Web scrapers kunnen gebaseerd zijn op verschillende technologieën zoals 'spidering' en 'pattern matching'. Lotfi e.a. (2021) onderscheidt web scrapers zowel op basis van hun toepassing als op basis van hun methodiek. Daarnaast biedt de programmeertaal waarin ze geïmplementeerd zijn ook telkens andere mogelijkheden (Bhatt e.a., 2023). Python lijkt de voor de hand liggende taal voor web scraping door de beschikbaarheid van een groot aantal libraries voor het beheren van data (Kumar & Roy, 2023).

Met de opkomst van AI blijkt het ook mogelijk om een webpagina te scrapen met gebruik van een Large Language Model (LLM) (Ahluwalia & Wani, 2024). AI biedt voordelen voor het omgaan met dynamische websites in vergelijking met de klassieke methodes en verbetert de efficiëntie van het scrapen (Ayuso e.a., 2024).

Web scraping vergroot de snelheid en het volume van de data die verwerkt kunnen worden aanzienlijk, en vermindert 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.

PRATIBA e.a. (2018), Rafsanjani (2022), Amin e.a. (2024), Sulistya e.a. (2024) scrapen onafhankelijk van elkaar GS aan de hand van het parsen van de DOM <sup>2</sup>.

Voorbeelden die gebruik maken van AI, hebben nog niet tot een publicatie geleid, maar blijven beperkt tot blog posts ("Web scraping experiment with AI (Parsing HTML with GPT-4 and GPT-4o)", 2025), ("Extracting Structured Data from Websites with Mirascope, Anthropic, and Streamlit", 2025)).

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.

Los van de manier waarop, blijft het doel stevast de omzetting van ongestructureerde gegevens op webpagina's naar gestructureerde data in een databank (Singrodia e.a., 2019).

Mitchell (2015) reikt een aantal methodes aan waarmee 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 belangrijk dan het scrapen op zich. Gescrapte data zijn niet gekenmerkt door strikte relaties, maar stellen eerder een collectie semi gestructureerde data voor. Een document store database lijkt in dit geval het meest aangewezen omdat er geen schema voor nodig is maar de data toch gestructureerd opgeslagen kunnen worden (Meier & Kaufmann, 2019). Lourenço e.a. (2015) vergelijkt de gangbare NoSQL databases en stelt MongoDB voor als document store database.

Het is belangrijk even stil te staan bij de vraag of web scraping legaal is, aangezien er constant met data van anderen gewerkt wordt? Volgens Castrillo-Fernández

---

<sup>2</sup>Document Object Model

(2015) is daar geen eenvoudig antwoord op, omdat de situatie afhankelijk is van het betrokken land maar over het algemeen zijn de meeste 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.

Tenslotte is het interessant om te beseffen dat een groot deel van de informatie op het internet gepresenteerd wordt onder de vorm van PDF-documenten (Singrodia e.a., 2019). Dit geeft aanleiding tot de beslissing om zowel HTML pagina's als PDF documenten te scrapen.

## 2.4. Natural Language Processing

Natural Language Processing (NLP) is een tak van de informatietechnologie die in staat is voor het omzetten van gesproken of geschreven menselijke taal in gestructureerde en verwerkbare data (Fanni e.a., 2023). Het heeft als doel om een machine tekst te laten lezen met volledig begrip van alle complexiteit die eigen is aan de menselijke taal.

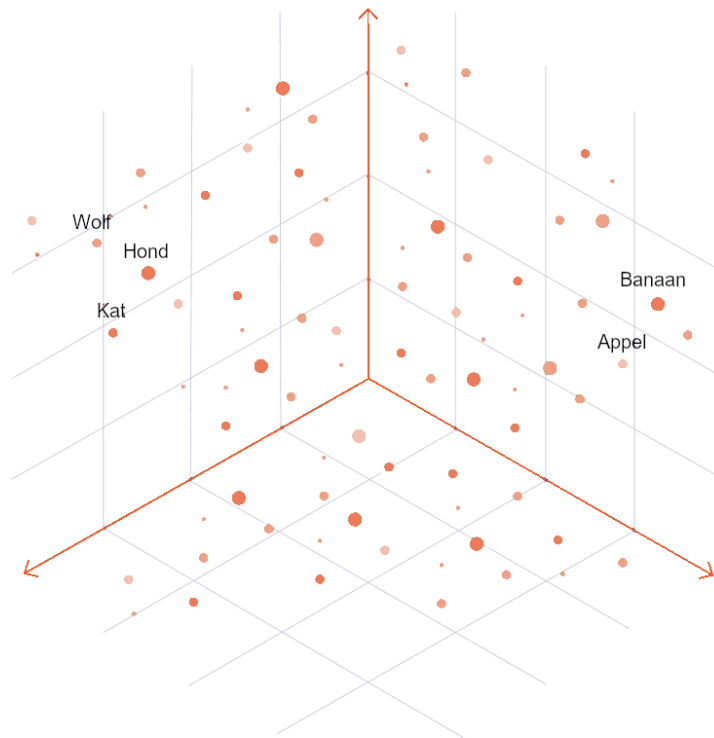
NLP is een verzamelnaam die meerdere toepassingen heeft waaronder het maken van samenvattingen, de classificatie van tekst, en het begrijpen van de boodschap (Khurana e.a., 2022). NLP is bijzonder populair als voorverwerking van AI toepassingen omdat het de ruwe taal omzet in betekenisvolle data waarmee modellen getraind worden.

Een bepaalde groep van toepassingen gaat over het analyseren van de betekenis van teksten en is gekend onder de naam van Natural Language Understanding. Daaronder vallen ook de technieken om de relevantie van een tekst voor een bepaald onderwerp te achterhalen. Een courante methode daarvoor is de "term frequency-inverse document frequency" (Havrlant & Kreinovich, 2017). Deze methode berekent een score die aangeeft hoe representatief een trefwoord is voor een bepaalde tekst.

## 2.5. Linked data

Alle academische publicaties worden wereldwijd beheerd door Crossref aan de hand van een metadata record en een Digital Object Identifier (DOI) (Hendricks e.a., 2020). Crossref vertegenwoordigt de gemeenschap van academische tijdschriften en streeft voortdurend technologische innovatie na teneinde het beheer van wetenschappelijke publicaties te verbeteren. In een dergelijk systeem, 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.





**Figuur 2.1:** Illustratie van embeddings.

Crossref laat toe om publicaties op te zoeken aan de hand van hun metadata, of om de metadata van een publicatie op te zoeken. Dit is mogelijk doordat deze gegevens met elkaar gelinkt zijn.

## 2.6. Semantic Search

Semantic search interpreteert de betekenis van de zoekterm om aan de hand daarvan aan alle mogelijke resultaten een score te geven in functie van hun gelijkheid met de zoekterm, in plaats van een exacte match tussen de teksten en de zoekterm (Bast e.a., 2016). Het beste resultaat is bijgevolg de tekst met de hoogste score. Die tekst leunt het dichtste aan bij de zoekterm.

Elke tekst wordt omgezet in een numerieke waarde, een vector, die een embedding genoemd wordt (Almeida & Xexéo, 2019). Dit is geïllustreerd in figuur 2.1. De embeddings zelf worden berekend aan de hand van een taalmodel (vb. Word2vec ("Word2vec", 2025), GloVe ("GloVe", 2025), ..). Een taalmodel is op zijn beurt gemaakt door een zeer grote verzameling teksten te trainen met machine learning technieken. Aan de hand van deze embeddings kan het systeem dan de gelijkheid beoordelen tussen de zoekterm en de teksten.

Dankzij semantic search verbetert het zoeken omdat het systeem niet alleen de woordelijke gelijkheid opzoekt, maar de virtuele multidimensionale positie van de zoekterm berekent en vervolgens vergelijkt met de posities van de teksten. Bijgevolg kan het resultaat ook gevonden worden in het geval van kleine afwijkingen in

de tekst (vb. afkortingen, schrijffouten, ..).

# 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 (GS) een interessant alternatief. 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 GS zoekpagina. De zoekpagina houdt geen gegevens bij van vorige queries en daarom is het zeer aannemelijk dat dezelfde resultaten zullen voorkomen bij opeenvolgende zoekopdrachten. GS 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 GS SERP <sup>1</sup>. Het opstellen van een zoekopdracht en het aanmaken van een alert worden verder uitgewerkt in hoofdstuk 4.

GS meldingen zijn e-mails in HTML formaat. Ze bevatten inhoud die overeenkomt met de GS 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

Die HTML moet omgezet worden in gestructureerde data door middel van HTML scraping technieken zoals te zien zijn in tabel 3.1.

---

<sup>1</sup>Search Engine Result Page

HTML scraping door een LLM	<ul style="list-style-type: none"> <li>• online model <ul style="list-style-type: none"> <li>– OpenAI</li> <li>– Generieke procedure onafhankelijk van het model</li> </ul> </li> <li>• lokaal model</li> </ul>
HTML scraping door het parsen van de DOM	<ul style="list-style-type: none"> <li>• BeautifulSoup</li> <li>• SerpAPI</li> </ul>

**Tabel 3.1:** HTML scraping technieken.

Web scraping wordt verder uitgewerkt in hoofdstuk 5.

Elk zoekresultaat is op basis van het algoritme van GS 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 variëteit 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.

- 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 (GS) 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 GS zoekpagina, omdat die dezelfde publicaties kan tonen bij opeenvolgende zoekingen. Daarentegen zijn er GS 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"<sup>A</sup> 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 toe hebben. Daarna 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 de zoektermen gebruikt uit tabel 4.1.

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

- 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

**Tabel 4.1:** VLIZ zoektermen.

# 5

## Web scraping

### 5.1. Inleiding

De e-mails afkomstig van Google Scholar (GS) 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 GS 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 GS SERP heeft een vaste structuur, namelijk een lijst met zoekresultaten bestaande uit:

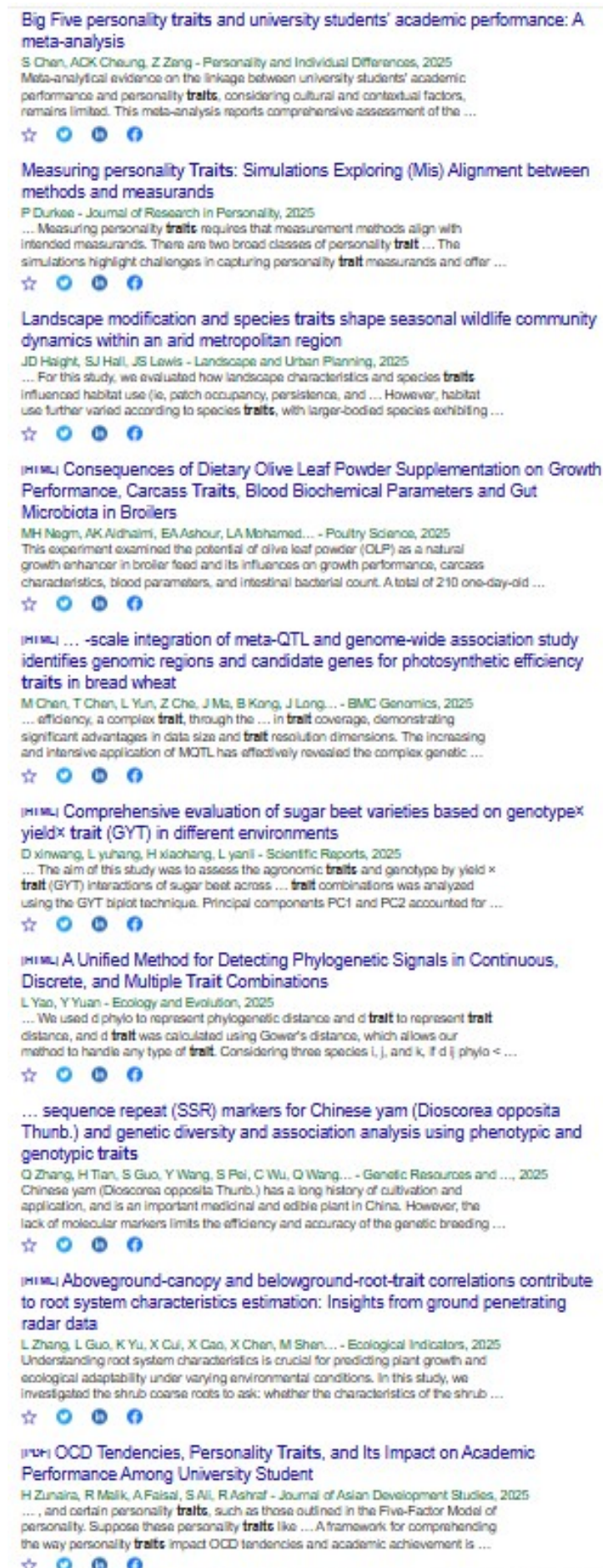
- titel
- link naar de webpagina van de publicatie
- auteurs
- naam van het tijdschrift
- jaartal
- abstract of fragment van het abstract

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

Informatie uit HTML pagina's halen is algemeen bekend 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 dienen dan weer als brandstof voor AI. Het punt is dat in die context meerdere technieken bestaan om dezelfde job te doen:

- web scraping met gebruik van een LLM





Figuur 5.1: Google Scholar SERP.

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

## 5.2. Web scraping met gebruik van een LLM

De meest bekende 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

```
<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://www.sciencedirect.com/science/article/pii/S0191886925001308" class=gse_alrt_titleš
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.

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 github "Web scraping Google Scholar SERP using OpenAI" (2025) maakt hiervan een implementatie aangepast voor de GS 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

```
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:** Opkuisen van de HTML.

in codefragment 5.3. De verwachte parameters die het model moet zoeken zijn te

```
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:** Opstellen van een prompt.

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 GS 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     }

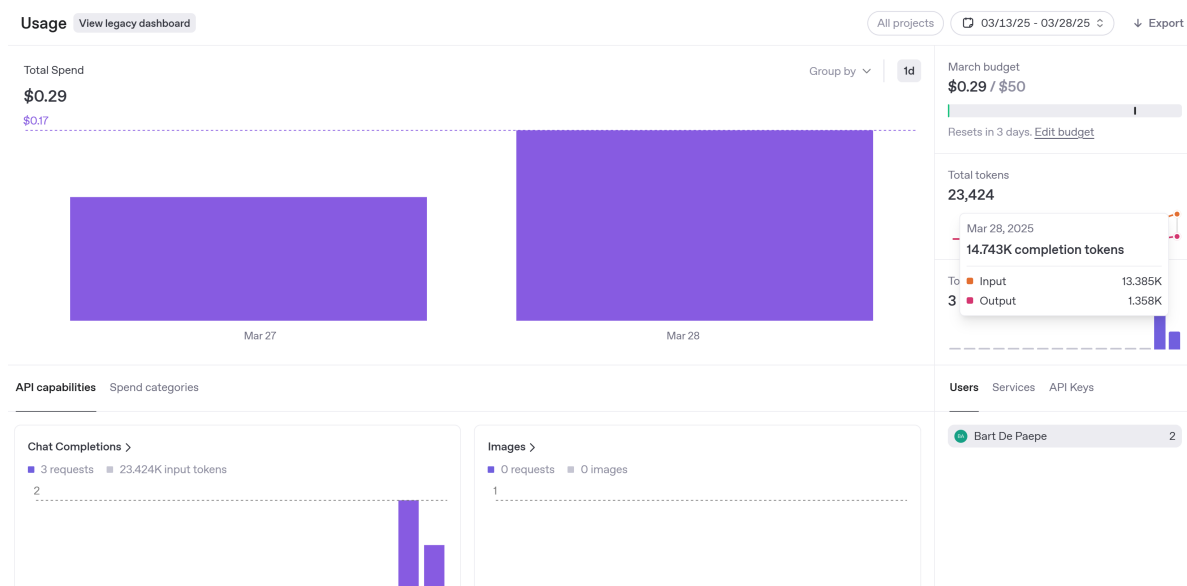
```

---

**Codefragment 5.4:** Opstellen van de zoekopties.

### 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.



**Figuur 5.2:** OpenAI dashboard.

De code op github “Web scraping Google Scholar SERP using Anthropic” (2025) maakt hiervan een implementatie aangepast voor de GS alerts. De parameters die opgezocht moeten worden, dienen eerst als velden en schema gedeclareerd te worden. Dit is mogelijk met Pydantic <sup>1</sup> zoals te zien is in codefragmenten 5.5 en 5.6. Met deze structuur van veld en schema kan vervolgens voor elke opdracht een

```

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

```

**Codefragment 5.5:** Aanmaken van een Pydantic field.

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).

<sup>1</sup>Pydantic laat toe om data te structureren en te valideren (“Introducing PydanticAI”, 2025).

---

```

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

```

---

**Codefragment 5.6:** Aanmaken 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:** Aanmaken van een Pydantic model.

### 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.

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

**Codefragment 5.8:** HTML parsen aan de hand van BeautifulSoup.

- 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-

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

**Codefragment 5.9:** Ollama uitvoeren met Docker.

graph AI Tutorial; Scrape websites easily with LLaMA AI” ([2025](#)) beschrijft stap voor 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 GS 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 parame-



---

```

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:** Configureren van Scrapegraph.

ters 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. Bij grote modellen moeten de tokens van de prompt opgezocht worden in veel meer data. Dat geeft aanleiding tot een zeer groot aantal instructies voor de CPU waardoor specifieke hardware noodzakelijk is, in het bijzonder grafische kaarten.

### 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 teruggegrepen naar



---

```
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:** Opstellen van een prompt voor Scrapegraph.

de meer klassieke benadering van web scraping waarbij de HTML structuur ontleed wordt aan de hand van custom code. De GS 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.



**Figuur 5.3:** HTML structuur van de GS alert.

Meerdere specifieke libraries laten toe om HTML te parsen, maar de meest bekende 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 GS 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 test met een SERP met 10 zoekresultaten was geslaagd. De gevraagde parameters worden gevonden voor elk van de 10 resultaten.

## 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 GS. Alleen focust hun product telkens op de online search engine en niet op alerts waardoor de zoekresultaten niet incre-

---

```
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:** 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:** Navigeren tussen tags aan de hand van BeautifulSoup.

menteel 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

Na de web scraping zijn de gegevens van elk zoekresultaat gestructureerd opgeslagen. Per definitie zijn al deze resultaten relevant voor IMIS want Google Scholar (GS) heeft ze opgezocht in functie van de zoekopdracht. De sortering van de zoekresultaten gebeurt op basis van de omvang van de integrale tekst, het aanzien van het tijdschrift en van de auteurs en het aantal citaties (inclusief de gedateerdheid ervan). Op die manier volgt GS de manier van werken binnen de academische wereld.

Maar daarnaast is het nog steeds mogelijk om voor elk zoekresultaat een score af te leiden die aangeeft hoe relevant de publicatie is voor de zoekopdracht aan de hand van de frequentie van elke zoekterm in de tekst. Algemeen wordt aangenomen dat hoe meer een zoekterm voorkomt in de tekst, des te relevanter die tekst zal zijn.

Daar bestaan Natural Language Processing (NLP) technieken voor. NLP is een verzamelnaam van een hele groep toepassingen en algoritmes die tekst omzetten in informatie. Deze zijn zeer divers 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” die alle technieken bundelt die toelaten om de betekenis van een tekst beter te begrijpen.

### 6.2. Relevantiescore

Het uitgangspunt is om een score te berekenen door te tellen hoe vaak een trefwoord voorkomt in de tekst. De frequentie is dan recht evenredig met de relevantie van die 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 woord voorkomt in de tekst. Dit is te zien in figuur 6.1. Toch is de frequentie op

	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.

zich nog niet de beste waardemeter voor een tekst. 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 score is dan ook omgekeerd evenredig met de frequentie van de zoekterm in alle documenten. De formules om de Tf-Idf te berekenen, zijn te zien in codefragment 6.1 en het resultaat in figuur 6.2. Uit de formules valt te zien dat de “inverse document

$$score = tf * idf$$

where

$tf$  = term frequency (see above)

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

$N$  = total number of documents

$df_t$  = the number of documents in which term  $t$  occurs

**Codefragment 6.1:** Berekening van de term frequency-inverse document frequency.

	I	like	this	article	medium	data
I like this article	0	0	0,48	0,48	0	0
I like medium	0	0	0	0	0,48	0
I like data	0	0	0	0	0	0,48

**Figuur 6.2:** Term frequency-inverse document frequency.

frequency” afhangt van het aantal documenten. 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 te zien in codefragment 6.2. De code op github “Natural Language Processing of

$$\text{score} = \text{topiccount} * \log(\text{“topic - sentence”ratio}) * \log(\text{“topic - noun”ratio})$$

**Codefragment 6.2:** Berekenen van de relevantiescore.

Google Scholar search results” (2025) maakt hiervan een implementatie aangepast voor de GS alerts.

### 6.3. Tekstverwerking

De BoW geeft wel aanleiding tot 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 stopwoorden <sup>1</sup> verwijdert en lemmatisering <sup>2</sup>.

Verder kunnen ook voornaamwoorden in de tekst verwijzen naar het trefwoord. Het is dus interessant om ze te vervangen door het trefwoord zelf zodat ze mee een invloed hebben op de relevantiescore (vb. tabel 6.1). “Coreference resolution”

voor	<p>“Het VLIZ is pionier in zeekennis. Het heeft de Simon Stevin als onderzoeksvaartuig. Daarmee voert het marien onderzoek uit.”</p>
na	<p>“Het VLIZ is pionier in zeekennis. Het VLIZ heeft de Simon Stevin als onderzoeksvaartuig. Simon Stevin als onderzoeksvaartuig voert het VLIZ marien onderzoek uit.”</p>

**Tabel 6.1:** Voor en na 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

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

<sup>2</sup>Lemmatisering 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, logging_service:
4          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         for spangroup in spans.values():
14             span_tuple = []
15             for span in spangroup:
16                 span_tuple.append(text[span.start_char:span.end_char])
17                 span_array.append(span_tuple)
18         return span_array
```

---

**Codefragment 6.3:** Coreference resolver

(Lee e.a., [2017](#)). Het volledige algoritme (“End-to-end Neural Coreference Resolution in spaCy”, [2025](#)) gebruikt deze scores om de tokens met de hoogste waarden te clusteren. Het opzoeken van paren aan de hand van de Spacy Coreference resolver (“Industrial-Strength Natural LanguageProcessing”, [2025](#)) is te zien in codefragment [6.3](#) en is eveneens geïmplementeerd in “Natural Language Processing of Google Scholar search results” ([2025](#)).

# 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 (vb. 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 lijst met

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

**Codefragment 7.1:** Originele URL van de publicatie



reguliere expressies die volgens “DOIs and matching regular expressions” (2025) gebruikt wordt om DOIs te matchen is te zien in tabel 7.1. Indien één van deze re-

```

• /^10.\d{4,9}/[-._;()/:A-Z0-9]+$ /i
• /^10.1002/[^s]+$ /i
• /^10.\d{4}/\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

```

**Tabel 7.1:** Reguliere expressies om een DOI te matchen.

guliere 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.

---

```

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         patterns = [r"10.1207/[\w\d]+\&\d+\_d+",
19                    r"10.1021/\w\w\d++",
20                    r"10.\d{4}/\d+-\d+X?(\d+)\d+<[\d\w]+:[\d\w]*>\d+.\d+.\w+;\d",
21                    r"10.1002/[^s]+",
22                    r"10.\d{4,9}/[-._;()/:A-Z0-9]+"]
23         return patterns

```

---

**Codefragment 7.2:** 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:** Opvragen van een sample op basis van de titel aan Crossref.

## 7.4. De DOI opzoeken op de webpagina van de publicatie

Indien de voorgaande stap geen DOI opleverde, dan wordt de 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 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 verschillend 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 Beautiful Soup 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 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 aan de hand van Crossref. In het geval van een opzoeking op de webpagina van de publicatie is de vindkans een pak kleiner.

---

```
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(),
9                                   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:** Openen van een online pdf.

---

```
1      ...
2      url = link.location_replace_url
3      options = webdriver.ChromeOptions()
4      download_folder =
5          os.path.join(str(Path(__file__).parent.parent.parent.parent),
6              "online_pdf")
7      profile = {
8          "plugins.plugins_list": [{"enabled": False, "name":
9              "Chrome PDF Viewer"}],
10         "download.default_directory": download_folder,
11         "download.extensions_to_open": "",
12         "download.prompt_for_download": False,
13         "download.directory_upgrade": True,
14         "plugins.always_open_pdf_externally": True
15     }
16     options.add_experimental_option("prefs", profile)
17     options.add_argument("start-maximized") # open Browser in
18         maximized mode
19     options.add_argument("disable-infobars") # disabling
20         infobars
21     options.add_argument("--disable-extensions") # disabling
22         extensions
23     options.add_argument("--disable-gpu") # applicable to
24         windows os only
25     options.add_argument("--disable-dev-shm-usage") #
26         overcome limited resource problems
27     options.add_argument("--no-sandbox") # Bypass OS security
28         model
29     options.add_argument("--headless")
30     driver = webdriver.Chrome(options=options)
31     driver.get(url)
32     driver.close()
33     ...
```

---

**Codefragment 7.6:** Nabootsen van gebruikersinteractie met Selenium.

# 8

## Semantic search

### 8.1. Inleiding

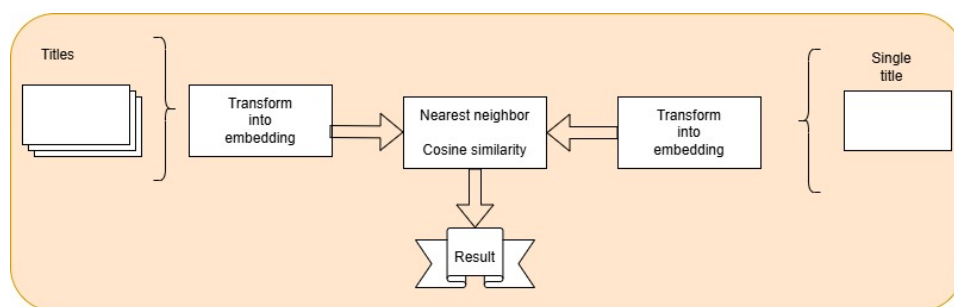
Alle publicaties in IMIS moeten uniek zijn. Op basis van voorgaande stappen is dat met 100% zekerheid te bepalen in het geval de DOI van de publicatie gevonden werd. Als de DOI niet gevonden werd, kan er niet met zekerheid gesproken worden, maar kan er wel een score berekend worden die aangeeft wat de kans is 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 gelijkheid tussen beide. Dit wordt voorgesteld in figuur 8.1.

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

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

- Pinecone
- Chroma



**Figuur 8.1:** Semantic search overzicht.

---

```
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:** Connecteren met Chroma.

- Weaviate
- Qdrant
- Milvus
- Vespa
- SingleStore
- Redis
- Elastic Stack
- Mongo
- ...
- (Elke database die een n-array van cijfers kan opslaan)

## 8.2. Chroma

Chroma ("Chroma", [2025](#)) is een gebruiksvriendelijke vector 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 getoond wordt in codefragment [8.3](#).

---

```
1     def initialize_embeddings(self):
2         result = requests.get(`https://www.vliz.be/nl/
3         imis?count=2000&module=ref&searchspcollist=39&show=json'`)
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:** 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:** Opzoeken van een titel.



	publisher	score	text	title
1	Scientific Data		0 - trait database containing 5 quantitative spore traits for all 344 d.	TraitAM, a global spore trait database for arbuscular mycorrhizal fun...
2	Plants	0.5915968418121338	0 - trait in cucumbers. This review aimed to summarize the factors, QTL	Research Advances and Perspectives on Early Flowering Traits in Cucum...
3	Ecology Letters	0.5227545499801636	0 - We subsequently calculated species mean trait values from individua...	WESEP: Tropical Tree Growth and Its Response to Climate Anomalies Is M...
4	Journal of Research in Personality		0 - We explored Big Six personality traits and their correlates among 18-	Big six personality traits in the Africa Long Life Study
5	Forest Ecology and Management		0 - distinct functional traits and resource-use strategies. Pinus korai...	Trait-mediated neighbourhood complementarity promotes tree growth in a...
6	Food5	-0.8008007152557373946875	0 - the traits of landraces and assigning values to them are crucial st...	Comprehensive Evaluation of Yunnan Potato Landraces: Agronomic, Sens...
7	Life		0 - Our aim was to explore whether this gene variant contributes to the...	A Pilot Study on the Correlation Between Psychometric Traits and ADRA...
8	Interdisciplinary Journal of Sharia Economics (IJSE)		0 - This study explores the factors influencing the purchase intention of...	The Impact of Personality Traits on Gen Z's Intentions to Buy Eco-Fri...
9	Applied Microbiology and		0 - Saccharomyces cerevisiae, a unicellular eukaryotic microorganism, inc...	Efficient genes identification via quantitative trait loci analysis b...

Figuur 8.2: Chroma resultaat.

```

1      ...
2      self.model =
          SentenceTransformer("nomic-ai/nomic-embed-text-v1",
                              trust_remote_code=True)
3      ...

```

**Codefragment 8.4:** Configureren van het model met MongoDB.

De test was succesvol. Bij wijze van proef werd de titel van één 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 opgeslagen kunnen worden. Dat wordt verder onderzocht.

### 8.3. MongoDB Atlas

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 embeddings geïndexeerd zoals getoond wordt in codefragment 8.6. Dan kan een titel opgezocht worden zoals getoond wordt in codefragment 8.7.

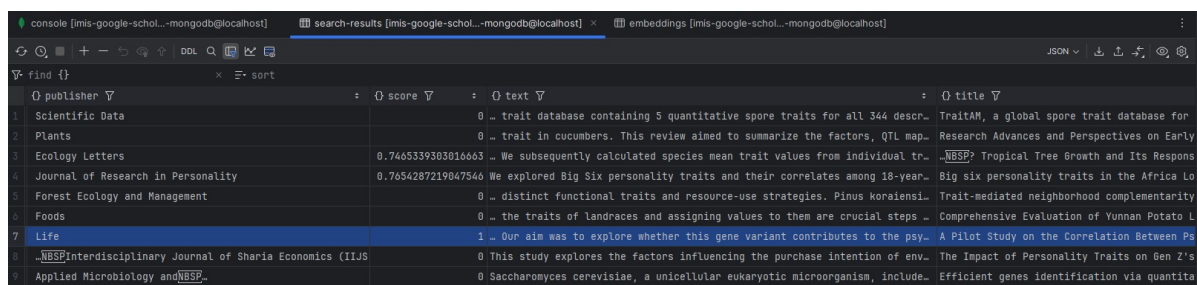
De test was succesvol. Bij wijze van proef werd de titel van één van de zoekresultaten toegevoegd aan de lijst van titels uit IMIS. Het resultaat is te zien in figuur 8.3. Deze oplossing breidt de bestaande MongoDB databank uit met embeddings. Voor artikels zonder indentificatie kan een score berekend 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-
10
11              dings.append(self.get_embedding(publication['StandardTitle']))
12
13          for i, (embedding, title) in enumerate(zip(embeddings,
14              data)):
15              doc = {
16                  "_id": i,
17                  "title": title,
18                  "embedding": embedding,
19              }
20          docs.append(doc)
21          self.db_service.set_collection("embeddings")
22          self.db_service.insert_many(docs)

```

**Codefragment 8.5:** Berekenen van de embeddings.



publisher	score	text	title
Scientific Data	0	... trait database containing 5 quantitative spore traits for all 344 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	0.7654287219847546	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
18        self.db_service.create_search_index(model=search_index_model)
```

---

**Codefragment 8.6:** 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:** Opzoeken van een titel.

# 9

## Conclusie

De IMIS collecties van het VLIZ kunnen voortaan semi-automatisch uitgebreid worden door het onderzoek van deze bachelorproef.

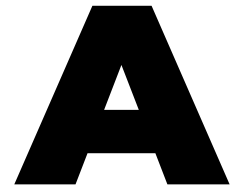
Google Scholar alerts zijn een automatische en incrementele bron van gegevens die zullen blijven stromen zolang de melding bestaat.

Web scraping technieken met LLMs verdienen bijzondere aandacht. Online modellen laten toe om een SERP te parsen en bovendien zijn ze bestand tegen interne veranderingen van de SERP. Na dit onderzoek moet er verder getest worden met verschillende lokale modellen op systemen met grafische kaarten. Het succes van de online modellen doet toch verwachten dat er ook lokaal een slaagkans moet zijn. Tot dan wordt er gewerkt met Beautiful Soup. Aangezien de structuur van de SERP vast ligt, is dit zeker geen slechte oplossing. De code van Beautiful Soup en de HTML structuur gaan hand in hand zodat wijzigingen van die laatste naar verwachting vrij eenvoudig op te vangen zijn in de code. Belangrijk daarbij is dat er een notificatiesysteem voorzien wordt dat aangaat wanneer er zo een wijziging zou optreden.

Het berekenen van een relevantiescore op basis van de frequentie van de zoekopdracht in het zoekresultaat geeft zeker aanleiding tot een nieuwe indicator, maar de mate waarin die effectief ook voorspelt of een publicatie interessant is voor IMIS is zeker vatbaar voor discussie. Verder onderzoek moet dus ook blijven inzetten op de "Natural Language Understanding" om de waarde van een publicatie voor IMIS beter te berekenen.

Het opsporen van de DOI van elke publicatie volgens een stapsgewijze procedure is efficiënt doch niet feilloos. Er blijft een minderheid van publicaties waarvoor de DOI niet bestaat of niet gevonden kan worden. Verder onderzoek moet vooral inzetten op het vinden van de juiste DOI op de webpagina van de publicatie wanneer daar meerdere DOIs aanwezig zijn. Dat is een eenvoudige taak voor een mens, maar dat is het niet voor een computer.

Tenslotte helpt het opzoeken van de gelijkheid tussen de titel van een publicatie en alle titels in IMIS aan de hand van “Semantic search” wanneer een manuele beslissing moet genomen worden om een publicatie toe te voegen aan IMIS. Bij een volledige match zal er geen twijfel zijn, maar bij een gedeeltelijke match, bijvoorbeeld in het geval van variaties in de titel, moet er beter onderzocht worden wat de threshold is om te beslissen om een publicatie aan IMIS toe te voegen.



# Google Scholar zoekopdracht

## A.1. Zoeken in Google Scholar (“Google Scholar Guide”, 2025)

### A.1.1. Basis zoeken

Google Scholar (GS) 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.

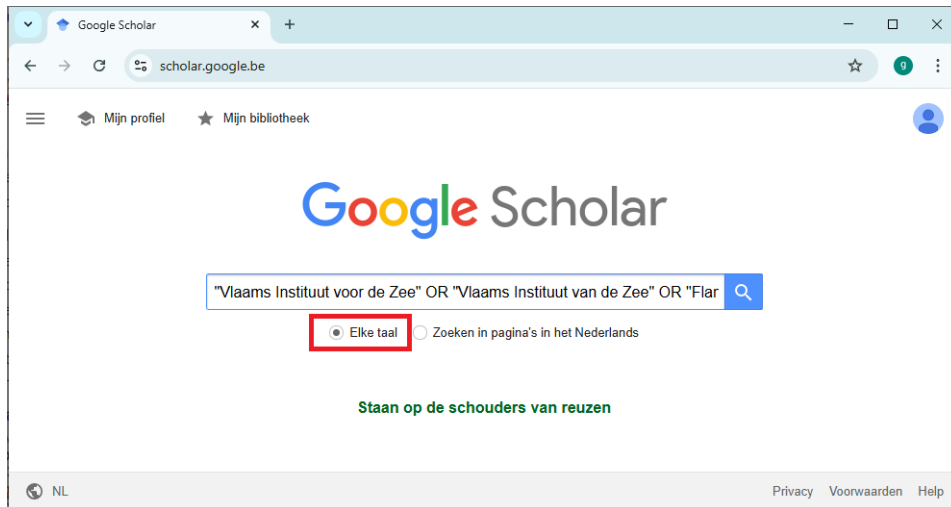
Wanneer een zoekopdracht verzonden wordt, dan is er een antwoord binnen de 3 seconden. Het resultaat kan vervolgens verder gefilterd worden:

- **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.<sup>1</sup>
- **Sorteren op datum:** Hierbij worden de resultaten gesorteerd op publicatiedatum.
- **Reviewartikelen:** Hierbij worden enkel state of the art publicaties gefilterd.<sup>2</sup>

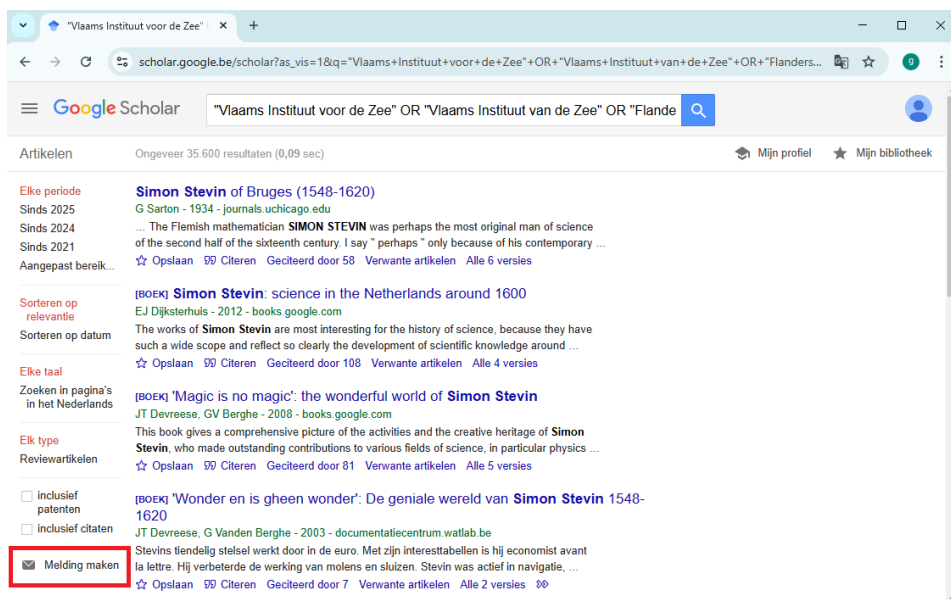
Elk resultaat kan verder uitgediept worden:

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

<sup>2</sup>Een 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 A.1:** Google Scholar user interface voor het basis zoeken van publicaties op basis van de ingevoerde zoektermen.



**Figuur A.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.

### A.1.2. Geavanceerd zoeken

GS heeft ook een meerdere geavanceerde 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.

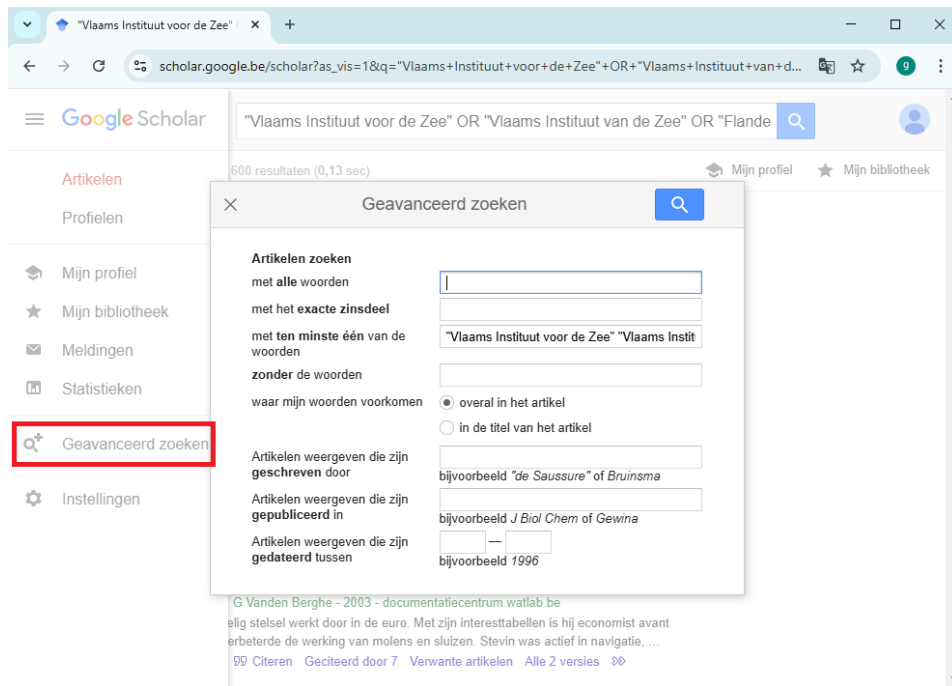
Voor alle bovenstaande filters kan ingesteld worden of er enkel in de titel of overal in de tekst mag gezocht worden.

Daarnaast zijn er 3 bijkomende geavanceerde 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 geavanceerde 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 A.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:

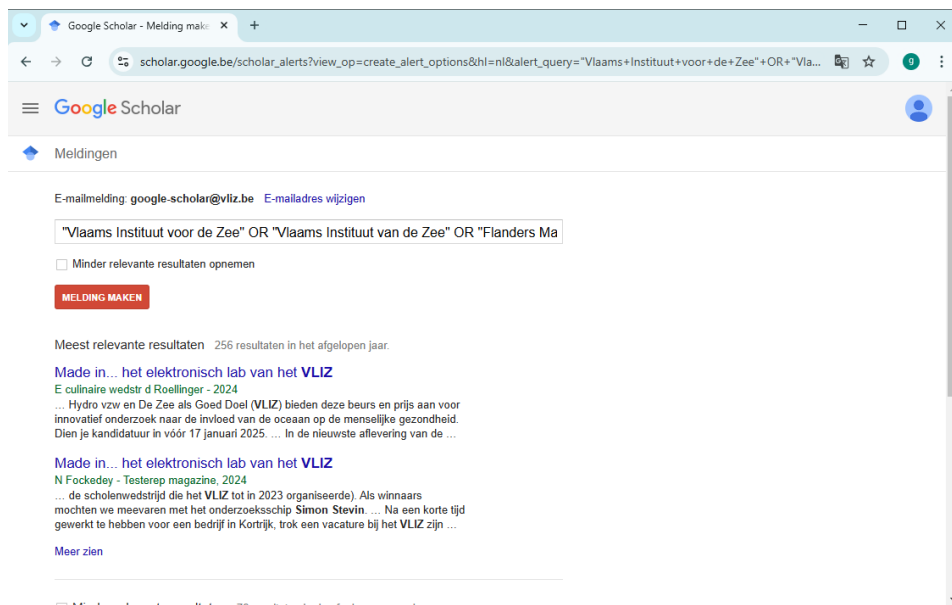
- **intitle:** De zoekresultaten bevatten de opgegeven zoekterm in de titel.
- **intext:** De zoekresultaten bevatten de opgegeven zoekterm in de tekst.
- **author:** De zoekresultaten bevatten de opgegeven auteur.
- **source:** De zoekresultaten bevatten de opgegeven uitgever.

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

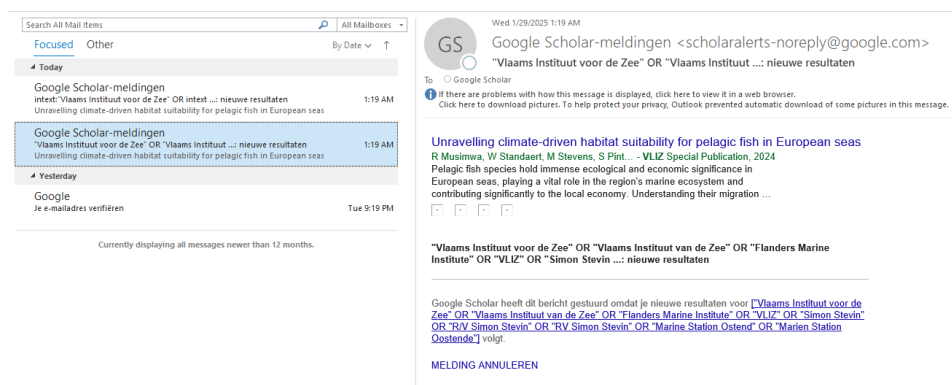
- **aanhalingstekens (" "):** De zoekresultaten bevatten de exacte tekst tussen aanhalingstekens.
- **liggend streepje (A-B):** Om aan te tonen dat 2 zoektermen sterk verbonden zijn.
- **liggend streepje (A -B):** Om de tweede zoekterm uit te sluiten van de resultaten.

## A.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

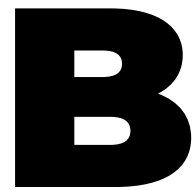


**Figuur A.4:** Google Scholar user interface voor het aanmaken van een e-mail alert voor de ingevoerde zoekopdracht.



**Figuur A.5:** Google Scholar e-mail alert met de nieuwe resultaten sinds het aanmaken van de melding en sinds de vorige melding.

klikken op 'Melding maken' zoals getoond op figuur A.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 geëactiveerd. Vanaf dan worden nieuwe publicaties die voldoen aan de filtercriteria systematisch doorgestuurd naar het e-mailadres.



## 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.

### **B.1. Inleiding**

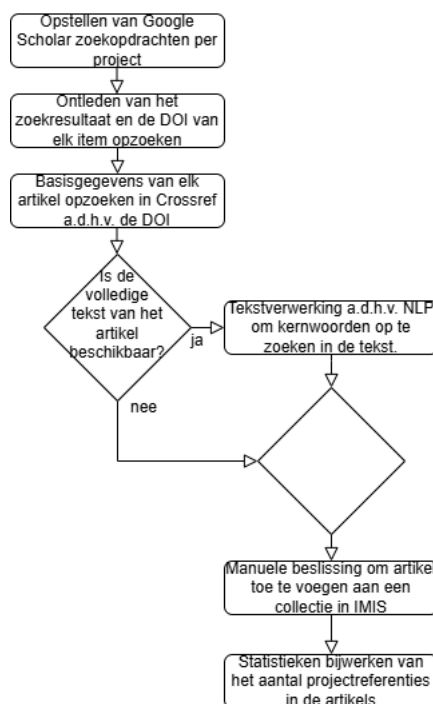
Het Vlaams Instituut voor de Zee (VLIZ) ("Vlaams Instituut voor de Zee", [2024](#)) is een pionier in zeekennis. Dit wetenschappelijk instituut gelegen in Oostende heeft onder andere een mandaat om een complete en geactualiseerde catalogus bij te houden van alle wetenschappelijke publicaties in de mariene 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 B.1:** Chronologische oplijsting 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.

## **B.2. Literatuurstudie**

IMIS (Haspeslagh & Vanden Berghe, 2002) 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).

### **B.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.

#### **B.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>
- Ahluwalia, A., & Wani, S. (2024). Leveraging Large Language Models for Web Scraping. <https://doi.org/10.48550/ARXIV.2406.08246>
- Almeida, F., & Xexéo, G. (2019). Word Embeddings: A Survey. <https://doi.org/10.48550/ARXIV.1901.09069>
- 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](https://doi.org/10.2991/978-94-6463-386-3_36)
- 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/>
- Ayuso, E., Dumfeh Brogya, M. S., Kumar Ahlawat, V., & Sain, M. (2024). From Manual to Machine: How AI is Redefining Web Scraping for Superior Efficiency: A Literature Review. *2024 International Conference on Communication, Control, and Intelligent Systems (CCIS)*, 1–9. <https://doi.org/10.1109/ccis63231.2024.10931912>
- Bast, H., Buchhold, B., & Haussmann, E. (2016). Semantic Search on Text and Knowledge Bases. *Foundations and Trends® in Information Retrieval*, 10(2–3), 119–271. <https://doi.org/10.1561/15000000032>
- Beautiful Soup*. (2025, mei 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>
- Crossref. (2024, november 1). <https://www.crossref.org/>
- Crossref API. (2025, mei 1). <https://github.com/fabiobatalha/crossrefapi>
- Crossref Commons for Python. (2025, mei 1). [https://gitlab.com/crossref/crossref\\_commons\\_py](https://gitlab.com/crossref/crossref_commons_py)
- DELPH-IN Overview. (2024, november 1). <https://github.com/delph-in/docs/wiki>
- DOIs and matching regular expressions. (2025, mei 1). <https://www.crossref.org/blog/does-and-matching-regular-expressions/>
- End-to-end Neural Coreference Resolution in spaCy. (2025, mei 1). <https://explosion.ai/blog/coref>
- 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/>
- Fanni, S. C., Febi, M., Aghakhanyan, G., & Neri, E. (2023). Natural Language Processing. In *Introduction to Artificial Intelligence* (pp. 87–99). Springer International Publishing. [https://doi.org/10.1007/978-3-031-25928-9\\_5](https://doi.org/10.1007/978-3-031-25928-9_5)
- Gemma3: Powerful AI Within Everyone's Reach. (2025, mei 1). <https://gemma3.org/>
- GloVe. (2025, mei 1). <https://nlp.stanford.edu/projects/glove/>
- Google Scholar Guide. (2025, mei 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](https://nl.wikipedia.org/wiki/Google_Scholar)
- GROBID: a machine learning software for extracting information from scholarly documents. (2024, november 1). <https://github.com/kermitt2/grobid>
- Habanero. (2025, mei 1). <https://github.com/sckott/habanero>
- Haspeslagh, J., & Vanden Berghe, E. (2002, januari 1). *IMIS: Integrated Marine Information System*. Vlaams Instituut voor de Zee. [https://www.researchgate.net/publication/33549102\\_IMIS\\_Integrated\\_Marine\\_Information\\_System](https://www.researchgate.net/publication/33549102_IMIS_Integrated_Marine_Information_System)
- Havrlant, L., & Kreinovich, V. (2017). A simple probabilistic explanation of term frequency-inverse document frequency (tf-idf) heuristic (and variations motivated by this explanation). *International Journal of General Systems*, 46(1), 27–36. <https://doi.org/10.1080/03081079.2017.1291635>
- 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](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>

- Industrial-Strength Natural LanguageProcessing*. (2025, mei 1). <https://spacy.io/>
- Introducing Meta Llama 3: The most capable openly available LLM to date*. (2025, mei 1). <https://ai.meta.com/blog/meta-llama-3/>
- Introducing PydanticAI*. (2025, mei 1). <https://pydantic.dev/>
- 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>
- Khurana, D., Koli, A., Khatter, K., & Singh, S. (2022). Natural language processing: state of the art, current trends and challenges. *Multimedia Tools and Applications*, 82(3), 3713–3744. <https://doi.org/10.1007/s11042-022-13428-4>
- 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>
- Lee, K., He, L., Lewis, M., & Zettlemoyer, L. (2017). End-to-end Neural Coreference Resolution. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. <https://doi.org/10.18653/v1/d17-1018>
- 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>
- Meier, A., & Kaufmann, M. (2019). NoSQL Databases. In *SQL amp; NoSQL Databases* (pp. 201–218). Springer Fachmedien Wiesbaden. [https://doi.org/10.1007/978-3-658-24549-8\\_7](https://doi.org/10.1007/978-3-658-24549-8_7)
- 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>
- 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](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>

- 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>
- PyMuPDF is a high-performance Python library for data extraction, analysis, conversion manipulation of PDF (and other) documents.* (2025, mei 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](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>
- 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](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, mei 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>
- 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](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](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](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](https://github.com/bart-de-paepe/google-scholar-openai/blob/main/app/src/services/parse_service.py)
- Word2vec*. (2025, mei 1). <https://code.google.com/archive/p/word2vec/>
- 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>