

Article

CWM Global Search—The Internet Search Engine for Chemists and Biologists

Alexander Kos * and Hans-Jürgen Himmler

AKos Consulting & Solutions Deutschland GmbH (AKos GmbH), Austr. 26, D-79585 Steinen, Germany;
E-Mail: contact@hjhimmler.de

* Author to whom correspondence should be addressed; E-Mail: software@akosgmbh.de;
Tel.: +49-7627-970068; Fax: +49-7627-970067.

Received: 11 October 2010; in revised form: 25 October 2010 / Accepted: 30 November 2010 /
Published: 3 December 2010

Abstract: CWM Global Search is a meta-search engine allowing chemists and biologists to search the major chemical and biological databases on the Internet, by structure, synonyms, CAS Registry Numbers and free text. A **meta-search engine** is a search tool that sends user requests to several other search engines and/or databases and aggregates the results into a single list or displays them according to their source [1]. CWM Global Search is a web application that has many of the characteristics of desktop applications (also known as Rich Internet Application, RIA), and it runs on both Windows and Macintosh platforms. The application is one of the first RIA for scientists. The application can be started using the URL <http://cwmglobalsearch.com/gswb>.

Keywords: internet search engine; meta-search engine; rich internet application; RIA; CWM global search; chemical structure search

1. Introduction

The Internet is changing how we work [2]. The Internet will also become the major information tool for chemists and biologists. Google has changed how we search and what we expect for answers. Our brain adjusts to the way we work [3]. We have to realize that the Internet is not a simple tool, but something that can control us, and this can be a frightening thought. We get used to always finding an answer, but we need to watch out that we do not forget to search not for the first best, but for the best answer; otherwise we lose on efficiency what we gain on speed.

The results of the Quick Search, structure, names and CAS Registry Numbers can be selected for a more comprehensive search using the Global Search. The results allow “cross-linking”, for example: Quick Search can be used to get a list of synonyms for ‘maslinic acid’, and then this list can be used in a Global Search. You will find the natural occurrence of this compound in an article in Wikipedia which only mentions ‘cralegolic acid’, (only one of the known synonyms for maslinic acid returned by the QuickSearch), but nowhere does the article in Wikipedia mention ‘maslinic acid’.

Figure 2. The CWM Global Search user interface Global Search—Combined CAS registry number synonym, and structure search for maslinic acid.

The screenshot displays the CWM Global Search interface. At the top, there are tabs for 'QuickSearch' and 'Global Search'. Below the tabs are buttons for 'Clear Form', 'Clear Structure', and 'Read' files (Mol File, SDF File, RXN File, RD File). The main interface is divided into several sections:

- Search Results:** A list of CAS Registry numbers and synonym names. The CAS number '4373-41-5' is selected. Synonyms include '(4aS,6aS,6bR,8aR,10R,11R,12aR,12bR,14bS)-10,11-Dihydroxy...', 'AIDS087536', 'AIDS-087536', 'Crategolic acid', 'Maslinic acid', and two variations of 'Olean-12-en-28-oic acid, 2,3-dihydroxy-, (2.alpha.,3.beta.)-'.
- Chemical Structure:** A 3D ball-and-stick model of the maslinic acid molecule is shown in the center.
- Major datasources:** A list of databases with checkboxes for selection. Selected sources include Google, PubChem, ChemSpider, eMolecules, NIST Webbook, and NCI Database.
- Specialized datasources:** A list of specialized databases with checkboxes. Selected sources include Wikipedia, ChemicalDB, ChemicalBook, ChemSynthesis, ChemExper, Buyersguide, Akos Samples, ZINC, Drugbank, and ChemBank.
- Options:** A section for search options, including 'Structure Search' (Include isomers, Include Substructures, Include similar compounds), 'Similarity coefficient' (set to 90), and 'Predefined profiles' (Search in all datasources, Search in chemistry datasources, Search in structure searchable datasources, Find physical properties, Find suppliers, Find safety information, Find spectra, Find literature, Find Open Access articles, Find drug information).
- My profiles:** A section for user-defined profiles.

At the bottom, there are buttons for 'Select/Deselect All', 'Show/Hide searchable fields', 'Save', 'Select', and 'Remove', along with a 'Start Global Search' button.

In Global Search, one can search presently over 46 different databases and Google. This includes chemistry centric databases such as ChemSpider, databases specialized in finding commercial suppliers such as eMolecules, databases with a focus on biological data such as PubChem, KEGG, ChEBI and the DrugBank, databases containing patent information such as SureChem, and literature databases such as MedLine. Some sources are gateways, like Open J-Gate, which is an electronic gateway to global journal literature in the open access domain. A complete list of the sources can be found in Appendix 1.

The user can select the data sources using predefined profiles, *i.e.*, chemistry, biology, availability, safety, *etc.*, and can create his own profiles, see Figure 2.

The search returns a collection of hyperlinks which allow direct drill down to the corresponding database to look at the data associated with the query. These drill-down pages are parsed by CWM Global Search to generate facts. These are true/false indicators that tell the user what kind of data he potentially can find in the corresponding database for his query. The result of the search is a grid with links and highlighted facts; see Figure 3. Examples of such facts are the presence of commercial

suppliers, availability of safety information or spectra, as well as the presence of known biological activities associated with the query.

Figure 3. Global Search Result Page of the combined CAS registry number synonym, and structure search for maslinic acid.

Identifier	Number of links	Status	Message	Literature
<input type="checkbox"/> NEXTBIO		NEXTBIO hit found	FULLSTRUCTURE	
			Biological activity	Clinical trials Literature
<input type="checkbox"/> NIAID		NIAID hit found	FULLSTRUCTURE	
	Physical properties		Biological test results	Literature
<input type="checkbox"/> NOVOSEEK		NOVOSEEK hit found	FULLSTRUCTURE	
				Literature
<input type="checkbox"/> CTD		Hit found	CASNUMBER	4373-41-5
			Pharmacog enomic data	Literature
<input type="checkbox"/> WIKIPEDIA		Result 1 of 1	CASNUMBER	4373-41-5
<input type="checkbox"/> PUBMED		Results: 1 to 20 of 32	CASNUMBER	4373-41-5
				Literature
<input type="checkbox"/> NOVOSEEK		Results: Pubmed(68) Free Full Text(3) U.S. Grants(0)	CASNUMBER	4373-41-5
				Literature
<input type="checkbox"/> NIAID			SYNONYMNAME	"Crategolic acid"
	Physical properties		Biological test results	Literature
<input type="checkbox"/> NIAID			SYNONYMNAME	"Maslinic acid"
	Physical properties		Biological test results	Literature
<input type="checkbox"/> Open 1 Gate		1 document(s) found	SYNONYMNAME	"Crategolic acid"

CWM Global Search supports exact structure searches, isomer searches (tautomers and stereoisomers), as well as substructure and similarity searches. At the moment, only PubChem and ChEBI support a substructure and similarity search using the drawn structure without query features. The structures can be copy and pasted from any of the major chemical drawing programs, or can be drawn directly in the JDraw [6] applet.

CWM Global Search supports single structure searches and multiple structure queries via support of SDFfiles. In addition, the program supports reaction based queries by support of RXN structures and/or RDFfiles. In case a reaction is used as a query, CWM Global Search searches the Internet for all reactants and products contained in the reaction. This is not a reaction search, but is useful for finding suppliers for starting materials, while at the same time also making sure that the product cannot be bought.

3. Additional Features

In the Quick Search and Global Search results page, a button "chemicalize" links to a page that displays calculated values for the compound; see Figure 4. Chemicalize (www.chemicalize.org) is developed by ChemAxon. Chemicalize uses ChemAxon's name to structure parsing to identify chemical structures from web pages and other text sources. It provides a large variety of predicted data related to each structure [7].

Figure 4. Chemicalize—Link to ChemAxon's calculations of molecular properties.

The screenshot displays the chemicalize.org interface for a specific molecule. The molecule is a complex polycyclic structure with multiple methyl groups and hydroxyl groups. The interface is organized into several panels:

- Molecule:** Shows the chemical structure with stereochemistry.
- Name:**
 - IUPAC name: (4aS,6aR,8aS,10R,11R,14bS)-10,11-dihydroxy-2,2,6a,6b,9,9,12a-heptamethyl-1,2,3,4,4a,5,6,6a,6b,7,8,8a,9,10,11,12,12a,12b,13,14b-icosahydricene-4a-carboxylic acid
 - Traditional name: (4aS,6aR,8aS,10R,11R,14bS)-10,11-dihydroxy-2,2,6a,6b,9,9,12a-heptamethyl-1,3,4,5,6,7,8,8a,10,11,12,12b,13,14b-tetradecahydricene-4a-carboxylic acid
- Elemental Analysis:**
 - Formula: C₃₀H₄₈O₄
 - Isotope formula: C₃₀H₄₈O₄
 - Composition: C (76.23%), H (10.24%), O (13.54%)
 - Isotope composition: C (76.23%), H (10.24%), O (13.54%)
 - Mass: 472.6997
 - Exact mass: 472.355260024
- Topology Analysis:**

Simple	Ring Counts	Path and distance
Atom count: 82		
Bond count: 86		
Cyclomatic number: 5		
Chain atom count: 12		
Chain bond count: 12		
Asymmetric atom count: 9		
Rotatable bond count: 1		
- Major Microspecies:** Shows the structure at pH=7.4, where the carboxylic acid group is deprotonated to a carboxylate anion.
- Geometry:** Includes a "Calculate Geometry" button.
- Lipinski-like filters:**
 - Lipinski's rule of five: no
 - Bioavailability: yes
 - Ghose filter: no
 - Lead likeness: no
 - Muegge filter: no
 - Weber filter: yes

At the bottom, there is a footer with the ChemAxon logo, a Creative Commons license notice, and a Creative Commons BY-NC-SA license icon.

4. What are the Differences between CWM Global Search and Other Systems such as PubChem or ChemSpider?

Unlike other systems such as PubChem or ChemSpider, CWM Global Search is NOT a database. We always search the most current snapshots of the supported data sources, thus also making sure that recently added records in the various data sources can be located.

This eliminates the problem that links to newly added PubChem records in ChemSpider and *vice versa* may not be found because of pending updates in the underlying database.

A major strength of CWM Global Search is the chemical structure search via the integrated structure editor (JDraw from Accelrys, Inc.). This structure editor supports copy/paste operations with major chemical drawing packages allowing the user to keep using his favorite structure drawing tool. Another unique feature of CWM Global Search allows you to draw a reaction; with one click you can find information for all reactants, reagents and products such as suppliers or safety information.

Comprehensive search—with a single click you can search for a structure, one or many associated CAS Registry Number(s), plus an arbitrary list of associated free text such as synonyms, brand names, and identifiers.

5. What are the Differences to SciFinder [8]?

SciFinder is the user interface for the world's largest chemistry database produced by the Chemical Abstract Service (CAS), a division of the American Chemical Society. However, it can only be accessed for a fee. While most academic institutions have access to SciFinder with an academic rate, many small and medium-sized companies simply cannot afford these fees, for them Global Search is a

very important first alternative. In all patent related cases, the most important question is whether a given structure is known or novel.

The most important aspect of the comparison with SciFinder is the fact that one cannot expect to find all answers to a given query in SciFinder. Many cases are known in which the use of Global Search produced important references that were not found in SciFinder. Thus, Global Search will provide the occasional user with quick answers, and it will give the professional information specialist the very important extra certainty that his search was as comprehensive as possible. The strength of CWM Global Search is its access to additional sources that are not considered publications, like entries in Wikipedia and databases.

6. How to Start CWM Global Search?

You can start the application using the URL <http://cwmglobalsearch.com/gswb>. If you start the application for the first time you might be asked to install the Microsoft Silverlight Plug-in. We support Internet Explorer, Firefox on Windows and Safari on Macintosh computers. Detailed information about the program can be found on the CWM Global Search homepage: www.akosgmbh.de/globalsearch. The free version is limited to Google, PubChem, ChemSpider, AKos Samples and the SureChem patent database. The free version will only show how many search results are found in the 40+ databases supported by the commercial version, and will not provide hyperlinks to the actual data.

7. Considerations

CWM Global Search relies on web services to generate the InChI names, and keys, and the availability of websites to search the underlying sources. We have no control over these web services and websites, they can be down for maintenance, moved to another location, or turned off. We periodically run searches to check for such issues, but a user should be aware that sometimes he has to re-execute a query when the web service or website is available again. According to our experience the availability fail rate of the web services and websites is very low. Since the whole application including the search engine is hosted on our server, we can upload a new version any time without involving the user, and we will do this because new interesting sources can be added monthly.

Maybe we should also discuss our business model. We try to keep the license fee as low as possible and augment this by giving sponsors room for their advertisements. A user in CWM Global Search will not interrupt his work to click on an advertisement, at least not very often. Therefore, the information that the sponsor wants to show must be in the advertisement. We run a slide show, giving each advertisement enough space to display an essential message. The picture will stay on the screen for a certain amount of time before the slide will change to a new one.

References and Notes

1. From Wikipedia, the free encyclopedia. Available online: http://en.wikipedia.org/wiki/Metasearch_engine (accessed on 1 December 2010).
2. Heuser, U.J. *Denken, wie das Netz es will*; Die Zeit: Hamburg, Deutschland, 23 September 2010.

3. Carr, N. Is Google Making Us Stupid? Available online: <http://www.theatlantic.com/magazine/archive/2008/07/is-google-making-us-stupid/6868/> (accessed on 1 December 2010).
4. Williams, A.J.; Tkachenko, V.; Lipinski, C.; Trophsa, A.; Ekins, S. Free online resources enabling crowd-sourced drug discovery. *Drug Discov. World* **2009**, *Winter*, 33–39.
5. Chemical Identifier Resolver beta 3. Available online: <http://cactus.nci.nih.gov/chemical/structure> (accessed on 1 December 2010).
6. JDraw is a Java applet structure editor from Accelrys, Inc.
7. What is chemicalize? Available online: <http://www.chemicalize.org/about.php> (accessed on 1 December 2010).
8. CAS Registry Number[®] and Synonym CAS Registry Number, SciFinder are registered trademarks of the American Chemical Society (ACS). All Rights Reserved.

Appendix

Table 1. A list of the sources searchable in Global Search. An up-to-date list of data sources can be found at: http://www.akosgmbh.de/globalsearch/databases_in_gs.htm.

Trademarks	Description	Link
<i>AKos Samples</i>	A database of approximate 6 million building blocks and screening compounds. All samples are checked for identity and purity by NMR. A network of suppliers can provide custom synthesis.	www.akosgmbh.de/AKosSamples
	BASE is one of the world's most voluminous search engines, especially for academic open access web resources. BASE is operated by Bielefeld University Library.	www.base-search.net
	BioMed Central is an STM (Science, Technology and Medicine) publisher which has pioneered the open access publishing model.	www.biomedcentral.com
	BuyersGuideChem is a directory of chemicals and chemical suppliers on the Internet.	www.buyersguidechem.de
	Chemical Entities of Biological Interest (ChEBI) is a freely available dictionary of molecular entities focusing on 'small' chemical compounds.	www.ebi.ac.uk/chebi
<u>CHEMBANK</u>	Initiative for Chemical Genetics. A freely available collection of data about small molecules (over 2000 compounds) and resources for studying their properties, especially their effects on biology.	http://chembank.broadinstitute.org
	The ChemExper Chemical Directory is mainly a supplier database for chemicals, and displays physical and chemical characteristics, structure, MSDS and more.	www.chemexper.com
Chemical Book	A supplier database mainly for the Chinese market.	http://www.chemicalbook.com

Table 1. Cont.

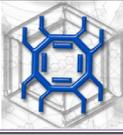
 <p>The Chemical Database The Department of Chemistry at the University of Akron</p>	<p>The Chemical Database will allow the user to retrieve information for any of 25,496 hazardous chemicals or 'generic' entries based on a keyword search.</p>	<p>http://ull.chemistry.uakron.edu/erd</p>
	<p>Chemicaland21.com aims to be a resource of individual chemical information including technical data, safety data, and related compounds.</p>	<p>http://chemicaland21.com</p>
	<p>This database allows users to search the NLM ChemIDplus database of over 370,000 chemicals.</p>	<p>http://chem.sis.nlm.nih.gov/chemidplus/</p>
	<p>ChemSpider hosts the largest and most diverse online database of chemical structures sourced from over 150 different data sources.</p>	<p>http://www.chemspider.com/</p>
	<p>ChemSynthesis is a database of compounds with their synthesis references and physical properties.</p>	<p>http://www.chemsynthesis.com/</p>
	<p>ClinicalTrials.gov is a registry of federally and privately supported clinical trials conducted in the United States and around the world.</p>	<p>http://clinicaltrials.gov/</p>
	<p>ChEBI CiteXplore combines literature search with text mining tools for biology.</p>	<p>http://www.ebi.ac.uk/citexplore</p>
	<p>Chemicals. CTD integrates a chemical subset of the Medical Subject Headings (MeSH®), the hierarchical vocabulary from the U.S. National Library of Medicine.</p>	<p>http://ctd.mdibl.org/</p>
	<p>The DrugBank database is a unique bioinformatics and cheminformatics resource that combines detailed drug (<i>i.e.</i>, chemical, pharmacological and pharmaceutical) data with comprehensive drug target (<i>i.e.</i>, sequence, structure, and pathway) information.</p>	<p>http://www.drugbank.ca/</p>
	<p>The Directory of Open Access Journals (DOAJ) lists open access journals, that is, scientific and scholarly journals that meet high quality standards by exercising peer review or editorial quality control.</p>	<p>http://www.doaj.org</p>
	<p>Envirofacts contains chemical data from several different program system databases: the <u>Aerometric Information Retrieval System</u>, the <u>Permit Compliance System</u>, and the <u>Toxics Release Inventory System</u>.</p>	<p>http://www.epa.gov/envirofw/gov/envirofw/</p>
	<p>Find Suppliers and Information for over 8 million unique chemicals!</p>	<p>http://www.emolecules.com/</p>

Table 1. Cont.

	US Environmental Protection Agency	http://www.epa.gov/
	This data source searches the European Patent Office database via the ChEBI CiteXplore search engine.	http://www.epo.org
	The FDA is responsible for protecting the public health by assuring the safety, efficacy, and security of human and veterinary drugs, biological products, medical devices, our nation's food supply, cosmetics, and products that emit radiation.	http://www.fda.gov
	Free patents online has hundreds of gigabytes of full-text data which is keyword searchable using the most powerful search engine in the industry.	http://www.freepatentsonline.com/
	Google Scholar is a freely-accessible Web search engine that indexes the full text of scholarly literature across an array of publishing formats and disciplines.	http://scholar.google.de/
	IPCS INCHEM is an invaluable tool for those concerned with chemical safety and the sound management of chemicals.	http://www.inchem.org/
	KEGG COMPOUND is a chemical structure database for metabolic compounds and other chemical substances that are relevant to biological systems.	http://www.genome.jp/kegg/compound/
	The NCI database contains 250,251 structures, which corresponds to the open part of the NCI database up until and including the latest release of the DTP cancer screen results of August 2000.	http://129.43.27.140/ncidb2/
	Look up whether a structure occurs in many different databases, both public and commercial. Currently loaded pointers to over 74 million entries from more than 100 databases, representing more than 46 million unique chemical structures.	http://cactus.nci.nih.gov/cgi-bin/lookup/search
	NextBio's integrated database contains publicly available data from a variety of sources, including GEO, caBIG, and Array Express among others.	www.nextbio.com
	National Institute of Allergy and Infectious Diseases: This database contains compounds that have been tested against HIV, HIV enzymes or opportunistic pathogens.	http://chemdb2.niaid.nih.gov
	The NIST Chemistry WebBook provides access to data compiled and distributed by NIST under the Standard Reference Data Program.	http://webbook.nist.gov/chemistry/

Table 1. Cont.

	novo seek is a biomedical search engine developed by <u>bioalma</u> for searching the published knowledge in biomedical literature.	http://www.novoseek.com
	Open J-Gate is an electronic gateway to global journal literature in open access domain.	http://www.openj-gate.com
	Database of human genetic variations on drug response.	http://www.pharmgkb.org/
	Proceedings of the National Academy of Sciences of the United States of America. PNAS Online contains the full text, figures, tables, equations, and references of all articles in PNAS dating back to 1990.	http://www.pnas.org/
	PubChem provides information on the biological activities of small molecules.	http://pubchem.ncbi.nlm.nih.gov/search/search.cgi
	PubMed is a service of the <u>U.S. National Library of Medicine</u> that includes over 18 million citations from MEDLINE and other life science journals for biomedical articles back to 1948. PubMed includes links to full text articles and other related resources.	http://www.ncbi.nlm.nih.gov/pubmed/
	<u>PubMed Central (PMC)</u> is a digital archive of life sciences journal literature that includes more than one million articles.	http://www.ncbi.nlm.nih.gov/pmc/
<i>--SIRI MSDS Index--</i>	The SIRI MSDS archive is maintained by Dan Woodard, MD dan@siri.org and Ralph Stuart, CIH. Our objective is to make critical safety information immediately and universally as accessible as possible.	http://hazard.com/msds/
	SureChem is making patent chemistry search faster, easier and more accessible.	http://www.surechem.org/
	Wikipedia's articles provide links to guide the user to related pages with additional information. There are currently more than 5000 molecules indexed in Wikipedia.	http://www.wikipedia.org/
	Database of commercially-available compounds for virtual screening.	http://zinc.docking.org/