# The LAILAPS Search Engine: A Feature Model for Relevance Ranking in Life Science Databases

Matthias Lange, Karl Spies, Christian Colmsee, Steffen Flemming, Matthias Klapperstück and Uwe Scholz

Research Group Bioinformatics and Information Technology, Leibniz Institute of Plant Genetics and Crop Plant Research (IPK), Gatersleben, Germany

{lange|spies|colmsee|flemming|klapper| scholz}@ipk-gatersleben.de

#### **Summary**

Efficient and effective information retrieval in life sciences is one of the most pressing challenge in bioinformatics. The incredible growth of life science databases to a vast network of interconnected information systems is to the same extent a big challenge and a great chance for life science research. The knowledge found in the Web, in particular in life-science databases, are a valuable major resource. In order to bring it to the scientist desktop, it is essential to have well performing search engines. Thereby, not the response time nor the number of results is important. The most crucial factor for millions of query results is the relevance ranking.

In this paper, we present a feature model for relevance ranking in life science databases and its implementation in the LAILAPS search engine. Motivated by the observation of user behavior during their inspection of search engine result, we condensed a set of 9 relevance discriminating features. These features are intuitively used by scientists, who briefly screen database entries for potential relevance. The features are both sufficient to estimate the potential relevance, and efficiently quantifiable.

The derivation of a relevance prediction function that computes the relevance from this features constitutes a regression problem. To solve this problem, we used artificial neural networks that have been trained with a reference set of relevant database entries for 19 protein queries.

Supporting a flexible text index and a simple data import format, this concepts are implemented in the LAILAPS search engine. It can easily be used both as search engine for comprehensive integrated life science databases and for small in-house project databases.

LAILAPS is publicly available for SWISSPROT data at

http://lailaps.ipk-gatersleben.de

#### 1 Introduction

"Finding information in the WWW is not much of a challenge. Just head for Google or Entrez and get the related web page or database entry." This issue can be heard frequently talking to biologists, who search information about a certain biological object [1]. However, finding reliable information about the function of a protein or seeking the protein that is involved in a certain activity in the cell cycle, is much more challenging. One has the choice of about 1,100 life science databases and billions of database records [2]. Even if one reduces the number of databases and records using database integration systems and powerful query systems, there

are still too many results for a simple query like "arginase" - an enzyme involved in the urea cycle. As shown in Figure 1, one gets 6739 hits in NCBI Entrez Protein databases, 2610 in Uniprot, and 568 in KEGG GENES (data from December 2009).



Figure 1: Search Engines in Life Science – the screenshots show examples of popular search engines, to investigate protein function: Google, Entrez, UniProt and KEGG.

Intuitively, the first choice are web search engines. Web site ranking techniques order query hits by relevance. But, trying to apply ranking methods that were developed to rank natural language text or WWW-sites to life science content and databases is questionable [3]. For example, the top ranked Google hit for "arginase" is a Wikipedia page. This is because the page is referenced by a high number of web-pages or Google assigned a manual defined priority rank. Here, the hypothesis is: A high hyperlink in-degree of a page means high popularity and high popularity means high relevance.

In order to find scientific relevant database entries, scientists need strong scientific evidence in relation to the specific research field. A dentist has other relevance criteria than a plant biologist or a patent agent. The intuitive and commonly used way at the scientist's desktop is query refinement. Criteria like who published, in which journal, for which organism, evidence scores, surrounding keywords etc. matter. Even complete search guides, e.g. for dentists were published [4].

Other ranking algorithms use Term Frequency - Inverse Document Frequency (TF-IDF) as ranking criteria. Apache-Lucene<sup>1</sup> is a popular implementation of this concept and is frequently

http://lucene.apache.org

used in bioinformatics, like LuceGene from the GMOD project [5], which is used for the EBI 'google' like search frontend EBeye. The TF-IDF approach works well, but misses the semantic context between the database entries and the query.

Another approach is the probabilistic relevancy ranking [6], whereby probabilistic values for the relevance of database fields and word combinations have to be predefined. In combination with a user feedback system, the probabilistic approach shows promising ranking performance [7].

Semantic search engines use methods from natural language processing and dictionaries to predict the semantic most similar database entries. Such conceptual search strategies, implemented in GoPubMed [8] or ProMiner [9], are frequently used algorithms in text mining projects.

The combination and abstraction of the mentioned relevance indicators motivated the development of a feature model for life science databases, a user feedback system and a machine learning ranking engine. In the next sections we present the LAILAPS system as method for relevance ranking in life science databases [10]. In particular, we present the used ranking concept, the used feature model, its implementation as neural network and ranking benchmarks.

## 2 LAILAPS Method for Relevance Ranking

Nearly every search engine, including LAILAPS, incorporates a scoring or ranking function to calculate a relevance for an entry. The central LAILAPS hypothesis is that the relevance score is context-dependent and the absolute rank position can be determined by sorting the relative scores. We apply information theory and postulate that the relevance of a database entry depends on two factors: its *content* and its *interpretation* by the search engine user. For the first factor we found that the relevance decision is based on a small number of core properties of the content. These core properties are used to deduce a feature model that expresses all important properties of a database entry as feature vector. The factor user interpretation is realized by LAILAPS in form of a feedback system and hand curated reference sets of relevant rated database entries. Both factors are used for the scoring algorithm. The algorithm uses artificial neural networks as method to estimate the relevance score of a database entry based on these factors.

#### 2.1 LAILAPS Feature Model

A relevance scoring function for life science databases is highly dependent on the underlying data. In contrast to full text data like PubMed or even traditional web sites, a database entry in a life science database is

- 1. structured and split into blocks (e.g. attributes, entities),
- 2. enriched by metadata (e.g. name of an attribute, quality information),
- 3. a compressed excerpt of a fact and
- 4. a mix of scalar, pure values and natural language text (e.g. stoichiometric biochemical reaction vs. function description).

These properties lead to the conclusion that classical ranking methods miss important properties and are suboptimal for life science database entries. Consequently, we had to define our own set of features that combines traditional and life science database specific ones. Motivated by the observation of user behavior during search engine result inspection, we introduced a set of 9 features F, presented in table 1. These features are intuitively used by scientists, who briefly screen database entries for potential relevance. The features are both sufficient to estimate the potential relevance, and efficiently quantifiable.

	feature	description			
$\overline{F_1}$	attribute	in which attribute the query			
		term was found			
$F_2$	database	in which database the found			
		entry is included			
$F_3$	frequency	the frequency of all query			
		terms in the entry and at-			
		tribute			
$F_4$	cooccurrence	express how close and in			
		which order the query terms			
		were found			
$F_5$	keyword	gives information, if good or			
		bad keyword are present near			
		to the query terms			
$F_6$	organism	to which organism the			
		database entry relates to			
$F_7$	sequence length	the length of the sequence de-			
		scribed by the database entry			
$F_8$	text position	which portion of the attribute			
		is covered by the query term			
$F_9$	synonym	gives information if the hit			
		was produced by an auto-			
		matic synonym expansion			

Table 1: Overview of the LAILAPS feature set

### 2.2 LAILAPS Relevance Ranking

The starting point of a ranking process is the query. Here we define a query as a set of AND-linked terms. Each term is a combination of alpha-numeric characters. Terms can be grouped into phrases, by quotation marks. We restrict the logical join of terms and phrases to a conjunctive form.

$$Q = t_1 \wedge t_2 \wedge \ldots \wedge t_n | t_n \in \text{term or phrase}$$
 (1)

The result of each query is a set R of database entries, which include all query terms or phrases. For these database entries, all attribute hits and related text positions are extracted:

$$find(Q) \to \{R\} \mid R = (t, d, \{(a, p)\})$$
 (2)

Where, t is the query term; d is the database entry; and  $\{(a, p)\}$  is a set of attribute-position pairs.

Based on the mentioned features, we define 9 feature functions  $find_f$ , that compute for each element r of the query result set R a scalar score  $\omega$  as element in the final feature vector  $\vec{\Omega} = (\omega_1, \dots, \omega_9)$ :

$$score_f(r) \to \omega \mid f \in F_1, \dots, F_9$$
 (3)

The final step is a ranking function, which computes a relevance score  $\tau$  from this 9-dimensional feature vector  $\vec{\Omega}$ :

$$relevance(\vec{\Omega}) \to \tau$$
 (4)

Because this function constitutes a regression problem, we chose a machine learning approach to estimate this function. For a given vector of quantified features of a particular database entry, a continous relevance score must be predicted. The focus of the relevance function is to automatically learn to recognize complex patterns from the feature vectors and make intelligent relevance predictions. Hence, we choose a machine learning approach to estimate the function  $relevance(\vec{\Omega})$ .

After we have predicted relevance scores, they can be ordered such as

$$\tau_1 < \tau_2 \Rightarrow \tau_1$$
 is less relevant than  $\tau_2$  (5)

The final relevance ranking is the order of all relevance scores:

$$(\tau_1, \dots, \tau_n) \mid \forall \tau : \ \tau_{n-1} < \tau_n \tag{6}$$

As shown in figure 2 those ranking tasks have to be executed for each query as sequential workflow.

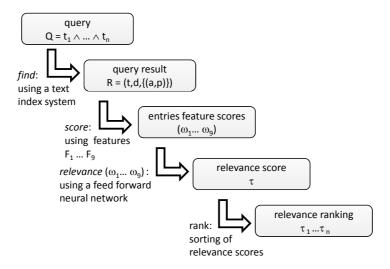


Figure 2: The LAILAPS Ranking Workflow

As mentioned before, the function *relevance* constitutes a regression problem. We have to map a vector of 9 features to a scalar, continous relevance score. In particular, supervised method perform well for text mining systems [11]. We found that artificial neural networks show best

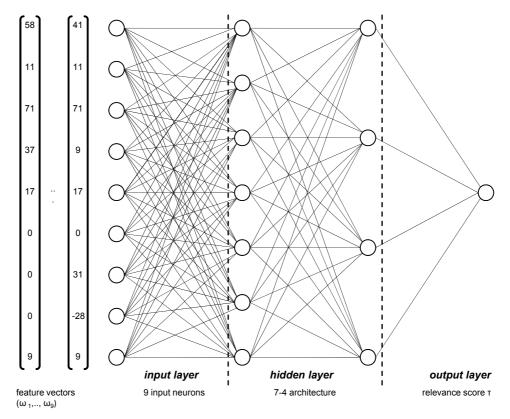


Figure 3: neural network for relevance ranking – the used neural network predict for a vector  $\vec{\omega}$  of 9 feature values the relevance of the database entry.

performance for our regression problem [12]. Using a set of reference data, we trained a feed-forward neural network with 9 neurons at the input and 7-4 neuron architecture in the hidden layer (see Figure 3). In order to train the network, we split up the training data into 80% for training and 20% for testing and used 500 training epochs. These parameters were estimated by minimizing the mean squared error  $\epsilon$  over the training set:

$$\epsilon = \frac{1}{n} \sum_{i=1}^{n} (\tau_i' - \tau_i)^2 \mid n = \text{size of training set}; \tau_i' = \text{manual score}; \tau_i = \text{predicted score}$$
 (7)

In each epoch the change of  $\epsilon$  was checked. After 500 epochs, no significant decrease was found and the final mean square error was 0.33 [12].

The crucial step for the neural network training is a set of true positive, manual evaluated relevance rankings. Our industrial and academic partners provided a set of plant metabolic queries with 1089 manually relevance ranked database records (see Table 2). This reference ranking list was separated into three confidence classes: high, medium and low. With these data, a neural network for plant metabolism was trained.

#### 2.3 LAILAPS Implementation

LAILAPS was developed as a 3-tier web application using Apache Tapestry<sup>2</sup> as web application-framework, ORACLE as database backend and a JAVA-implemented ranking logic featured by

<sup>&</sup>lt;sup>2</sup>http://tapestry.apache.org

Query Text		Category
		Split-Up
		(hi/me/lo)
industrial use case 1	20	6 / 4 / 10
"pinene synthase"	18	10 / 3 / 5
industrial use case 2	39	8 / 13 / 18
industrial use case 3		14 / 32 / 18
"gamma tocopherol methyltransferase"		21 / 9 / 8
"ent-kaurene synthase"	65	17 / 38 / 10
"chlorophyll synthase"	77	17 / 54 / 6
industrial use case 4	134	35 / 68 / 31
"cinnamyl-alcohol dehydrogenase"	214	45 / 36 / 133
industrial use case 5		3 / 4 / 10
"dihydrokaempferol 4-reductase"	65	9 / 29 / 27
"l-ascorbate peroxidase"	100	69 / 12 / 19
"morphine 6-dehydrogenase"	35	2 / 15 / 18
"zeaxanthin epoxidase"	51	21 / 2 / 28
"squalene monooxygenase"	84	24 / 30 / 30
"acetoacetyl-coa synthetase"	68	14 / 36 / 18

Table 2: Overview of the training data set.

the Java Object Oriented Neural Engine (JOONE)<sup>3</sup>. The backend database stores the loaded life science databases in an entity-attribute-value (EAV) [13] adapted database schema. This flexible concept enables the import of RFC-compatible CSV-formatted<sup>4</sup> exports from life science databases, whereas each row comprise a database record and its columns the fields.

For the imported databases, an inverse text index is computed and synonyms are loaded. In the public available system, we provide protein synonyms extracted from UNIPROT/SWISSPROT [14] and BRENDA [15]. The ranking logic

- 1. computes all matched positions per database entry,
- 2. extracts for each database entry an n-dimensional feature vector with 9 basis feature classes and
- 3. predicts a relevance probability using user specific neural networks, which maps the feature vectors to the users specific relevance score.

Since the ranking profile is computed in context of the authenticated user, a valid user login is recommended. Otherwise, a default context with a pre-trained neural network is used.

The relevance ordered query hits are assigned to the ranking profile and rendered by the middleware into a number of web pages (see Figure 4). The embedded feedback system provides a tool for the user to rate the relevance of a particular database entry. The ratings are stored in the database backend and used to accumulate user training data. The LAILAPS feedback system is

<sup>3</sup>http://www.jooneworld.com/

<sup>4</sup>http://tools.ietf.org/html/rfc4180

transparently embedded into the result browser. By opening the database detail browser, AJAX<sup>5</sup> code is injected into the original data HTML presentation, which is for example, provided by the SRS@EBI data retrieval system [16]. This code collects the user rating of the database entry and tracks user interactions. This feedback is used to enrich the original training data in the related ranking profile.

#### 3 Results and Discussion

Searching scientific databases effectively necessitates the use of contemporary software to locate desired and meaningful information according to the users scientific or project priorities. However, the combination of relevance ranking and life science data retrieval is still missing in life science information systems. LAILAPS fills this gap and provides a search engine for integrated or single instance life science databases in combination with an efficient ranking system.

We have evaluated the relevance prediction using the standard measures for precision Pr, recall Re and  $F_1$  score:

$$Pr = \frac{TP}{TP + FP}; \quad Re = \frac{TP}{TP + FN}; \quad F_1 = \frac{2 * Pr * Re}{Pr + Re}$$

As mentioned, we have used a curated reference set of plant metabolic queries. In order to decide whether a database entry has been correctly ranked or not, we do not consider its concrete ranking position. Because of combinatorial explosion, for human curators it is nearly impossible to find a correct relevance order among hundreds of database entries. Rather, the knowledge quality of a certain database entry is crucial. Consequently, the database entries were classified into three confidence classes: "HIGH", "MEDIUM" and "LOW". The "HIGH"-class comprises the top entries with proven and reliable knowledge. The class "MEDIUM" includes all those, that could be interesting but are uncertain. The class "LOW" includes all data, that has insufficient knowledge value or bad quality indicators. For each query, those classes form sub sets of a query result R such as  $R = R_H \cup R_M \cup R_L$ , whereas  $R_H \cap R_M \cap R_L = \emptyset$ . Each set forms a continuous window in the list of results. E.g., for a query result of 100 database entries, the window for  $R_H$  ranges from position 1-20,  $R_M$  from position 21-80, and  $R_L$  from 81-100.

In order to evaluate the LAILAPS rankings, we have to compare a list of relevance ordered database entries with the confidence classes of the reference set. We can't compare absolute ranking positions, because the elements in the confidence classes have no order. But we can say in which range (window) of rank positions a LAILAPS ranked entry should be, to fall into a certain confidence class. This consideration is used to define the true positives (TP), false positives (FP) and false negatives (FN):

For example, for one use case the curators sorted 20 entries into the class "HIGH". 18 of the top 20 LAILAPS ranked result are in the window  $R_H$ . In this case, the precision is  $\frac{18}{18+2} = 0.9$ .

Because all database entries of the reference set were found by LAILAPS, the recall is Re = 100%. This is because the text indexing is the basis for matching query terms. The text index uses a tokenizer, that decomposes text into words. LAILAPS use the same text decomposition

<sup>&</sup>lt;sup>5</sup>Asynchronous JavaScript and XML

error type	LAILAPS benchmark semantics
true positive (TP)	the rank position is in the same window as in the reference set
false positive (FP)	the rank position is in a different window as in the reference set
false negative (FN)	the database entry was not found by LAILAPS

Table 3: Definition of evaluation error types.

rules and the same databases as the reference retrieval systems. Furthermore, no synonym expansion of query terms was used. Thus, the hits of the reference systems could be reproduced by LAILAPS and no false negative hits exit.

The overall benchmarking result of 16 queries is shown in Figure 5. The average recall, precision and  $F_1$  values are shown in Table 4.

confidence class	precision	recall	$F_1$
"HIGH"	62%	100%	76
"MEDIUM" ∪ "HIGH"	81%	100%	90

Table 4: Evaluation of LAILAPS relevance ranking results.

In average we achieve a better precision than existing search engines [3]. Training and benchmark data for the non-industrial use cases are available by request to the authors.

The training data were collected in the application scenario of plant research and queries for protein functions. In order to bring LAILAPS to a broad community, we provide a public installation of the LAILAPS search engine. This will help to improve performance and to include more user domain specific ranking profiles. Because of limited database and project resources, the public, non-commercial version is restricted to SWISSPROT data. A comprehensive set of databases is available for registered users on request to the authors.

LAILAPS combines a clean, powerful and easy to use human computer interface with a machine learning based, context sensitive ranking system. It comprises a search engine with a self trained neural network ranking system, which brings a new quality and determinism into the scientific knowledge exploration.

#### 4 Conclusion

In this paper, we presented a feature model for relevance ranking in life science databases. The presented concept is to learn ranking pattern from the ranking behaviour of scientist, who make use of data retrieval systems. We combined user ranking profiles, a reference set of relevant database entries for 19 protein queries as training set to a ranking model. Queries are formulated as simple keyword lists and will be expanded by synonyms. The model is used to extract per database entry a feature vector. Using supervised machine learning approach, we were able to predict a user specific relevance score per feature vector. We found a feed forward neural network as best performing method for the regression problem of ranking in life science databases.

Supporting a flexible text index and a simple data import format, this concepts are implemented in the LAILAPS search engine. It can easily be used both as search engine for comprehensive

integrated life science databases and for small in-house project databases. Using expert knowledge as training data for a predefined neural network or using users own relevance training sets, a reliable relevance ranking of database hits has been implemented. LAILAPS is public available for SWISSPROT data at http://lailaps.ipk-gatersleben.de

## Acknowledgements

Thanks to Thomas Münch for the web hosting service and the valuable technical support.

#### References

- [1] A. Divoli, M. Hearst, and Wooldridge M. A. Evidence for Showing Gene/Protein Name Suggestions in Bioscience Literature Search Interfaces. In *Pacific Symposium on Biocomputing*, volume 13, pages 568–579, 2008.
- [2] M.1 Y. Galperin and G. R. Cochrane. Nucleic Acids Research annual Database Issue and the NAR online Molecular Biology Database Collection in 2009. *Nucleic Acids Research*, 37(suppl\_1):D1–4, 2009.
- [3] Matthew Richardson, Amit Prakash, and Eric Brill. Beyond pagerank: machine learning for static ranking. In *WWW '06: Proceedings of the 15th international conference on World Wide Web*, pages 707–715, New York, NY, USA, 2006. ACM.
- [4] J. Day. The quest for information: A guide to searching the internet. *Journal of Contemporary Dental Practice*, 2(4):033–043, 2001.
- [5] B. O'Connor, A. Day, S. Cain, O. Arnaiz, L. Sperling, and L. Stein. Gmodweb: a web framework for the generic model organism database. *Genome Biology*, 9(6):R102, 2008.
- [6] Nicholas C. Ide, Russell F. Loane, and Dina Demner-Fushman. Essie: A Concept-based Search Engine for Structured Biomedical Text. *Journal of the American Medical Informatics Association*, 14(3):253–263, 2007.
- [7] E. Agichtein, E. Brill, and S. Dumais. Improving web search ranking by incorporating user behavior information. In *SIGIR '06: Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 19–26, New York, NY, USA, 2006. ACM.
- [8] A. Doms and M. Schroeder. GoPubMed: exploring PubMed with the Gene Ontology. *Nucleic Acids Research*, 33(suppl\_2):W783–786, 2005.
- [9] Daniel Hanisch, Katrin Fundel, Heinz-Theodor Mevissen, Ralf Zimmer, and Juliane Fluck. Prominer: rule-based protein and gene entity recognition. *BMC Bioinformatics*, 6(Suppl 1):S14, 2005.
- [10] Matthias Lange, Karl Spies, Joachim Bargsten, Gregor Haberhauer, Matthias Klapperstück, Michael Leps, Christian Weinel, Röbbe Wünschiers, Mandy Weißbach, Jens Stein, and Uwe Scholz. The LAILAPS Search Engine: Relevance Ranking in Life Science Databases. *Journal of Integrative Bioinformatics*, 7(2):e110, 2010.

- [11] Chris Burges, Tal Shaked, Erin Renshaw, Ari Lazier, Matt Deeds, Nicole Hamilton, and Greg Hullender. Learning to rank using gradient descent. In *ICML '05: Proceedings of the 22nd international conference on Machine learning*, pages 89–96, New York, NY, USA, 2005. ACM.
- [12] J. Bargsten. Nutzerprofilgestütztes Schätzen von Rankingparametern für Suchmaschinen-Anfragen an integrierte biologische Datenbanken. Diplomarbeit, Martin-Luther-Universität Halle-Wittenberg, Institut für Informatik, 2009.
- [13] P. M. Nadkarni, L. Marenco, R. Chen, E. Skoufos, G. Shepherd, and P. Miller. Organization of Heterogeneous Scientific Data Using the EAV/CR Representation. *Journal of the American Medical Informatics Association*, 6(6):478–493, 1999.
- [14] The UniProt Consortium. The Universal Protein Resource (UniProt) 2009. *Nucleic Acids Research*, 37(suppl\_1):D169–174, 2009.
- [15] I. Schomburg, A. Chang, and D. Schomburg. BRENDA, enzyme data and metabolic information. *Nucleic Acids Research*, 30(1):47–49, 2002.
- [16] E. M. Zdobnov, R. Lopez, R. Apweiler, and T. Etzold. The EBI SRS server new features. *Bioinformatics*, 18(8):1149–1150, 2002.

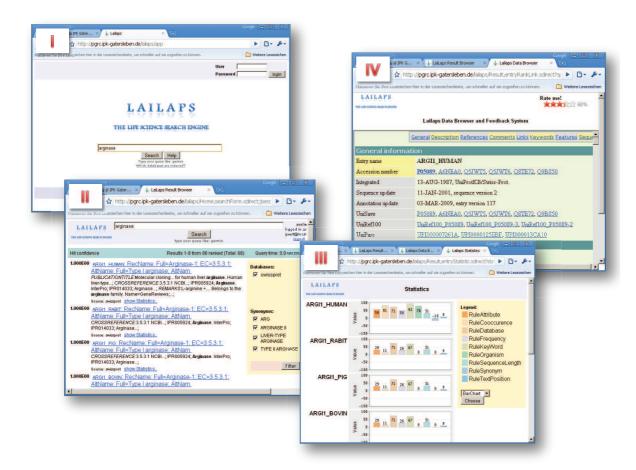


Figure 4: LAILAPS view of an example query session – The four browser windows represent the common query workflow of LAILAPS. Start screen with an optional login to load a customized ranking profile (I). The query might be specified as a single word, combination of words delimited by a whitespace or quoted phrases. The ranked query result is presented as a list of relevance sorted database accessions (II), with a short hit description, the evidence value, the hyper link to the original data source and a link to the scoring statistics in window (III). The exploration of hits is supported by a detail browser and feedback system (IV). The original data is extracted from the original source and displayed. The user have to be aware that LAILAPS is querying the original data from the SRS system, which installed at the EBI (http://srs.ebi.ac.uk). Thus, some outdated accession may result in broken links. In this page, the user might rank the relevance of the hit for later training of the user ranking profile.

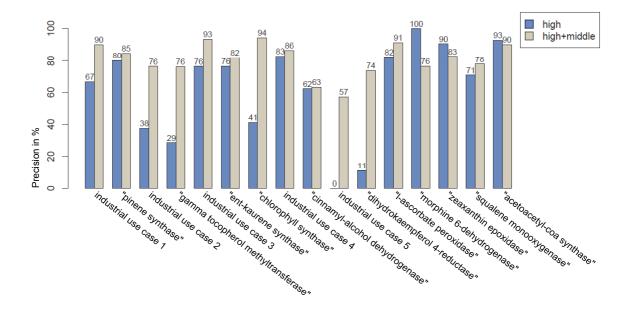


Figure 5: LAILAPS recall for use case queries