

Retrieving Biomedical Literature  
-  
An Open Source Search Engine Based on  
Open Access Resources

Hayda Almeida, Ludovic Jean-Louis, Marie-Jean Meurs



Biocuration 2016, Genève

April 2016

# Biomedical Literature Retrieval

- Scientific databases → support for research and health care
- Large amount of open access data available

**PubMed** BD

24,000,000<sup>+</sup>

PY: Since 1809

**PMC** OA

1,200,000<sup>+</sup>

Since 1973

+

=

**25,403,053 records**

- Retrieval of relevant information → critical task
- Scientific journal articles → input for many tasks (Almeida et al., 2014)



# Scientific Database Search: Challenges

## 1 Article content searched



article abstract



article full-text

### Use of full-text search:

- Better support for literature analysis tasks (Gay et al., 2005)
- Improvement in search results (Nourbakhsh et al., 2012)
- Access to more relevant information in articles (Van Auken et al., 2014)

# Scientific Database Search: Challenges

## 2 Express search in query language

- Users frequently reformulate queries (Dogan et al., 2009)
- Few users generate advanced queries (Shariff et al., 2013)
- Most searches made by inexperienced users (Yoo and Mosa, 2015)

Natural language: `alpha-amylase from Cryptococcus flavus`

Query language: `alpha-amylase OR alpha amylase AND (cryptococcus OR (cryptococcus AND flavus))`

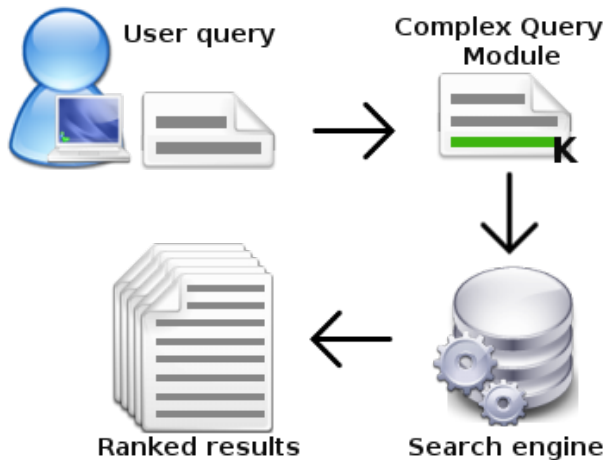
## Suggested Approach

- Search engine for biomedical open access data
- Ability to handle natural language queries
- Document Indexing Module
  - Parsing → XML (PubMed) and NXML (PMC)
  - Relevant fields → full-text, abstract, metadata
  - Indexing → map fields to index schema
- Complex Query Module
  - User input → natural language processing
  - Query types & query strategies per type
  - Query expansion → UMLS Metathesaurus annotations


## Pipeline: Document Indexing



## Pipeline: Query Search



## Document Indexing Module

- Based on Apache Solr 
- One index entry per document
- Document semantic representation: {document field, content}

Field	PubMed	BD	PMC	OA	Field	PubMed	BD	PMC	OA
<b>1</b> Article title	✓			✓	<b>8</b> Reference title	✗			✓
<b>2</b> Journal title	✓			✓	<b>9</b> Reference IDs	✗			✓
<b>3</b> Abstract	✓			✓	<b>10</b> Object captions	✗			✓
<b>4</b> Body section titles	✗			✓	<b>11</b> PMCID	✓			✓
<b>5</b> Body full content	✗			✓	<b>12</b> PMID	✓			✓
<b>6</b> Author names	✓			✓	<b>13</b> Article keywords	✗			✓
<b>7</b> Reference authors	✓			✓	<b>14</b> Publication year	✓			✓



## Complex Query Module: Query Types

- Keyword query,  $K_Q$ 
  - No stop-words among query terms
  - "AIDS versus HIV"
- Open Question query,  $O_Q$ 
  - Presents interrogative cues
  - "what is the difference between HIV and AIDS?"
- Statement query,  $S_Q$ 
  - Does not present interrogative cues
  - Has stop-words
  - "the difference between HIV and AIDS"

## Complex Query Module: Query Generation

- Each query type → different strategy
- Search Fields → where to look for query terms
- Phrase Search Fields → look for query terms appearing in sequence
- Boost → increase document relevance with a coefficient in query time

Type	Search Fields	Boost?	Phrase Search Fields	Boost?
$K_Q$	abstract, body, keywords {...}	abstract, body	title, body, authors {...}	title, authors
$O_Q$	title, abstract, body {...}	captions, body {...}	captions, abstract {...}	abstract, body
$S_Q$	body, authors, keywords {...}	title, captions {...}	title, abstract {...}	title, body

## Complex Query Module: Query Expansion

- MetaMap (Aronson and Lang, 2010) → UMLS Methathesaurus
- Avoid redundancy → annotations without any terms in user query

User query: "AIDS versus HIV"

MetaMap annotations:

```
"HIV+ [HIV Seropositivity]"
```

```
"AIDS [Acquired Immunodeficiency Syndrome]"
```

```
"HIV [HIV]"
```

Expanded query:

```
"AIDS versus HIV Acquired Immunodeficiency Syndrome"
```

## Preliminary Evaluation Data

- Large data → challenge finding manual annotations
- 19 manually annotated sets {query, target article ID}
  - Biocurators support: mycoCLAP (Strasser et al., 2015) database
  - Each enzyme entry → 1<sup>+</sup> article(s)
  - Articles retrieved from scientific literature databases
- 9 {query, PMCID}, 10 {query, PMID}

Q#	Target article ID	User query	mycoCLAP ID
Q <sub>3</sub>	PMC2780388	characterization of GH5 beta-mannanase enzyme from <i>Aspergillus niger</i>	MAN5A_ASPNG
Q <sub>4</sub>	PMC3092853	characterization of GH16 beta-glucanase from <i>Aspergillus fumigatus</i>	MLG16B_ASPFU
Q <sub>15</sub>	PMID1400249	characterization of <i>Candida albicans</i> maltase	AGL13B_CANAL
Q <sub>16</sub>	PMID12761390	beta-1,4-galactanases from <i>Humicola insolens</i> and <i>Myceliophthora thermophila</i>	GAN53A_HUMIN

## Evaluation Metrics

- Pseudo-judgement → top 20 ranked results
- Reciprocal Rank (RR)

Computed for each query

Inverse of target article ranking

$$RR = \frac{1}{\text{position}}$$

- Mean Reciprocal Rank (MRR)

Computed for all queries

RR average for the 19 {query, target article ID} sets

$$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{position}}$$

# Preliminary Results

Q#	PMC rank	bioMine rank	bioMine RR score	Q#	PubMed rank	bioMine rank	bioMine RR score
Q <sub>1</sub>	3	<b>2</b>	0.500	Q <sub>10</sub>	2	<b>1</b>	1.000
Q <sub>2</sub>	1	20	0.050	Q <sub>11</sub>	N/A	<b>7</b>	0.143
Q <sub>3</sub>	1	2	0.500	Q <sub>12</sub>	1	<b>1</b>	1.000
Q <sub>4</sub>	2	8	0.125	Q <sub>13</sub>	2	<b>1</b>	1.000
Q <sub>5</sub>	2	13	0.077	Q <sub>14</sub>	1	<b>1</b>	1.000
Q <sub>6</sub>	9	<b>1</b>	1.000	Q <sub>15</sub>	2	<b>1</b>	1.000
Q <sub>7</sub>	2	5	0.200	Q <sub>16</sub>	1	N/A	0.000
Q <sub>8</sub>	1	17	0.059	Q <sub>17</sub>	N/A	<b>1</b>	1.000
Q <sub>9</sub>	1	10	0.100	Q <sub>18</sub>	1	N/A	0.000
				Q <sub>19</sub>	1	<b>1</b>	1.000
<b>total # of queries = 19</b>				<b>MRR = 0.513</b>			

# Preliminary Results

Q#	PMC rank	bioMine rank	bioMine RR score	Q#	PubMed rank	bioMine rank	bioMine RR score
Q <sub>1</sub>	3	2	0.500	Q <sub>10</sub>	2	1	1.000
Q <sub>2</sub>	1	20	0.050	Q <sub>11</sub>	N/A	7	0.143
Q <sub>3</sub>	1	2	0.500	Q <sub>12</sub>	1	1	1.000
Q <sub>4</sub>	2	8	0.125	Q <sub>13</sub>	2	1	1.000
Q <sub>5</sub>	2	13	0.077	Q <sub>14</sub>	1	1	1.000
Q <sub>6</sub>	9	1	1.000	Q <sub>15</sub>	2	1	1.000
Q <sub>7</sub>	2	5	0.200	Q <sub>16</sub>	1	N/A	0.000
Q <sub>8</sub>	1	17	0.059	Q <sub>17</sub>	N/A	1	1.000
Q <sub>9</sub>	1	10	0.100	Q <sub>18</sub>	1	N/A	0.000
				Q <sub>19</sub>	1	1	1.000
<b>total # of queries = 19</b>				<b>MRR = 0.513</b>			

# Preliminary Results

Q#	PMC rank	bioMine rank	bioMine RR score	Q#	PubMed rank	bioMine rank	bioMine RR score
Q <sub>1</sub>	3	2	0.500	Q <sub>10</sub>	2	1	1.000
Q <sub>2</sub>	1	20	0.050	Q <sub>11</sub>	N/A	7	0.143
Q <sub>3</sub>	1	2	0.500	Q <sub>12</sub>	1	1	1.000
Q <sub>4</sub>	2	8	0.125	Q <sub>13</sub>	2	1	1.000
Q <sub>5</sub>	2	13	0.077	Q <sub>14</sub>	1	1	1.000
Q <sub>6</sub>	9	1	1.000	Q <sub>15</sub>	2	1	1.000
Q <sub>7</sub>	2	5	0.200	Q <sub>16</sub>	1	N/A	0.000
Q <sub>8</sub>	1	17	0.059	Q <sub>17</sub>	N/A	1	1.000
Q <sub>9</sub>	1	10	0.100	Q <sub>18</sub>	1	N/A	0.000
				Q <sub>19</sub>	1	1	1.000
<b>total # of queries = 19</b>				<b>MRR = 0.513</b>			



# Preliminary Results

Q#	PMC rank	bioMine rank	bioMine RR score	Q#	PubMed rank	bioMine rank	bioMine RR score
Q <sub>1</sub>	3	2	0.500	Q <sub>10</sub>	2	1	1.000
Q <sub>2</sub>	1	20	0.050	Q <sub>11</sub>	N/A	7	0.143
Q <sub>3</sub>	1	2	0.500	Q <sub>12</sub>	1	1	1.000
Q <sub>4</sub>	2	8	0.125	Q <sub>13</sub>	2	1	1.000
Q <sub>5</sub>	2	13	0.077	Q <sub>14</sub>	1	1	1.000
Q <sub>6</sub>	9	1	1.000	Q <sub>15</sub>	2	1	1.000
Q <sub>7</sub>	2	5	0.200	Q <sub>16</sub>	1	N/A	0.000
Q <sub>8</sub>	1	17	0.059	Q <sub>17</sub>	N/A	1	1.000
Q <sub>9</sub>	1	10	0.100	Q <sub>18</sub>	1	N/A	0.000
				Q <sub>19</sub>	1	1	1.000
<b>total # of queries = 19</b>				<b>MRR = 0.513</b>			

## Preliminary Results

Q#	PMC rank	bioMine rank	bioMine RR score	Q#	PubMed rank	bioMine rank	bioMine RR score
Q <sub>1</sub>	3	2	0.500	Q <sub>10</sub>	2	1	1.000
Q <sub>2</sub>	1	20	0.050	Q <sub>11</sub>	N/A	7	0.143
Q <sub>3</sub>	1	2	0.500	Q <sub>12</sub>	1	1	1.000
Q <sub>4</sub>	2	8	0.125	Q <sub>13</sub>	2	1	1.000
Q <sub>5</sub>	2	13	0.077	Q <sub>14</sub>	1	1	1.000
Q <sub>6</sub>	9	1	1.000	Q <sub>15</sub>	2	1	1.000
Q <sub>7</sub>	2	5	0.200	Q <sub>16</sub>	1	N/A	0.000
Q <sub>8</sub>	1	17	0.059	Q <sub>17</sub>	N/A	1	1.000
Q <sub>9</sub>	1	10	0.100	Q <sub>18</sub>	1	N/A	0.000
				Q <sub>19</sub>	1	1	1.000
<b>total # of queries = 19</b>				<b>MRR = 0.513</b>			

## Conclusion and Ongoing Work

- Scientific literature search in article abstracts and full-text
- Processing of natural language queries
- Target articles ranked in bioMine at first position  $\approx 50\%$  of the time
- Use of open access data
- Source code publicly available

<https://github.com/BigMiners/bioMine>

### Next steps

- Improvement of full-text document retrieval
- Development of web-based user interface

# Thank you!

## Questions?

### References

Almeida et al., *Machine Learning for Biomedical Literature Triage*, PLOS ONE, 2014.

Gay et al., *Semi-automatic Indexing of Full Text Biomedical Articles*, AMIA Annual Symposium Proceedings, 2005.

Nourbakhsh E. et al., *Medical Literature Searches: A Comparison of PubMed and Google Scholar*, Health Information & Libraries Journal, 2012.

Van Auken et al., *BC4GO: A Full-text Corpus for the BioCreative IV GO Task*, Database, 2014.

Dogan et al., *Understanding PubMed User Search Behaviour through Log Analysis*, Database, 2009.

Shariff et al., *Retrieving Clinical Evidence: A Comparison of PubMed and Google Scholar for Quick Clinical Searches*, Journal of Medical Internet Research, 2013.

Yoo and Mosa, *Analysis of PubMed User Sessions Using a Full-Day PubMed Query Log: A Comparison of Experienced and Nonexperienced PubMed Users*, Journal of Medical Internet Research, 2015.

Aronson A. and Lang F., *An Overview of MetaMap: Historical Perspective and Recent Advances*, Journal of the American Medical Informatics Association, 2010.

Strasser K. et al., *mycoCLAP, the Database for Characterized Lignocellulose-active Proteins of Fungal Origin: Resource and Text Mining Curation Support*, Database, 2015.

## Corpus Description

### Baseline Database (BD) Files

- Journal article abstract, citations, books
- Publication years since at least 1809
- 24,350,000<sup>+</sup> entries

### Open Access (OA) Subset

- Full-text journal articles
- Publication years since at least 1973
- 1,200,000<sup>+</sup> entries

Total entries indexed: **25,403,053**

**Evaluation Data:** query, PMCID

Q#	Target article ID	User query	mycoCLAP ID
Q <sub>1</sub>	PMC3068306	alpha-amylase from <i>Cryptococcus flavus</i> activity characterization	AMY13A_CRYFL
Q <sub>2</sub>	PMC3312866	<i>Aspergillus fumigatus</i> beta-glucosidase purification and characterization	BGL3C_ASPFU
Q <sub>3</sub>	PMC2780388	characterization of GH5 beta-mannanase enzyme from <i>Aspergillus niger</i>	MAN5A_ASPNG
Q <sub>4</sub>	PMC3092853	characterization of GH16 beta-glucanase from <i>Aspergillus fumigatus</i>	MLG16B_ASPFU
Q <sub>5</sub>	PMC3180650	purification and characterization of an exo-polygalacturonase from <i>Fusarium oxysporum</i>	PGX28B_FUSOX
Q <sub>6</sub>	PMC3223205	<i>Phanerochaete chrysosporium</i> GH61 purification and characterization	PMO9D_PHACH
Q <sub>7</sub>	PMC3312857	purification and characterization of an alpha-L-rhamnosidase from <i>Aspergillus nidulans</i>	RHA78E_EMENI
Q <sub>8</sub>	PMC2291056	xylanase characterization from <i>Leucoagaricus gongylophorus</i>	XYN11A_LEUGO
Q <sub>9</sub>	PMC2702311	recombinant expression and characterization of xylanase from <i>Trichoderma reesei</i>	XYN11B_TRIRE

**Evaluation Data: query, PMID**

Q#	Target article ID	User query	mycoCLAP ID
Q <sub>10</sub>	PMID20562284	bifunctional alpha-L-arabinofuranosidase /xylobiohydrolase from <i>Penicillium purpurogenum</i>	ZAX43C_PENPU
Q <sub>11</sub>	PMID10215597	enzymatic properties alpha-mannosidase <i>Aspergillus saitoi</i>	MSD47S_ASPPH
Q <sub>12</sub>	PMID20709852	characterization of <i>Magnaporthe oryzae</i> cellobiohydrolase	CBH6A_MAGOR
Q <sub>13</sub>	PMID9758835	substrate specificity of alpha-L-arabinofuranosidase from <i>Aspergillus awamori</i>	ABF51A_ASPAW
Q <sub>14</sub>	PMID7708682	cloning and characterization <i>Candida albicans</i> chitinase	CHI18B_CANAL
Q <sub>15</sub>	PMID1400249	characterization of <i>Candida albicans</i> maltase	AGL13B_CANAL
Q <sub>16</sub>	PMID12761390	beta-1,4-galactanases from <i>Humicola insolens</i> and <i>Myceliophthora thermophila</i>	GAN53A_HUMIN
Q <sub>17</sub>	PMID12427996	<i>Neotyphodium</i> sp beta-1,6-glucanase expression and characterization	BGN5A_NEOSP
Q <sub>18</sub>	PMID21653698	purification of endo-beta-1,3-galactanase from <i>Flammulina velutipes</i>	EBG16A_FLAVE
Q <sub>19</sub>	PMID9872754	<i>Aspergillus oryzae</i> beta-xylosidase optimum pH and temperature	XYL3A_ASPOR