IBM LanguageWare

# Text Mining in Life Sciences

## UIMA Framework and Knowledge Discovery at IBM
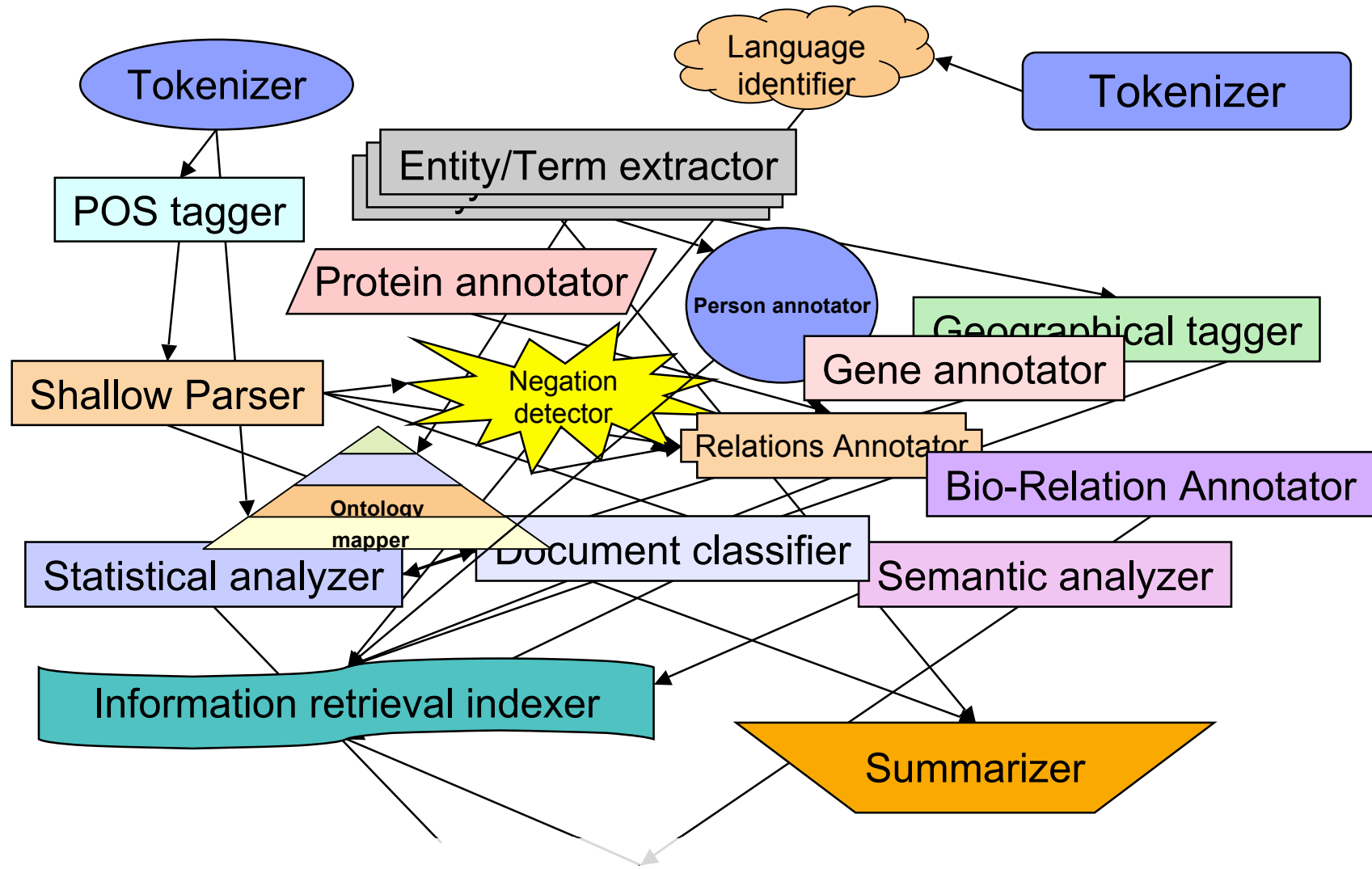
Alex Nevidomsky, IBM Dublin Software Lab.
alex_nevidomsky@ie.ibm.com

- **UIMA Concept**

- **UIMA reference implementation**

- **IBM OmniFind: UIMA-enabled platform**

- **Text analytics and domain customization**

- **UIMA applications in HCLS**

  - Clinical trials search

  - Chemical search
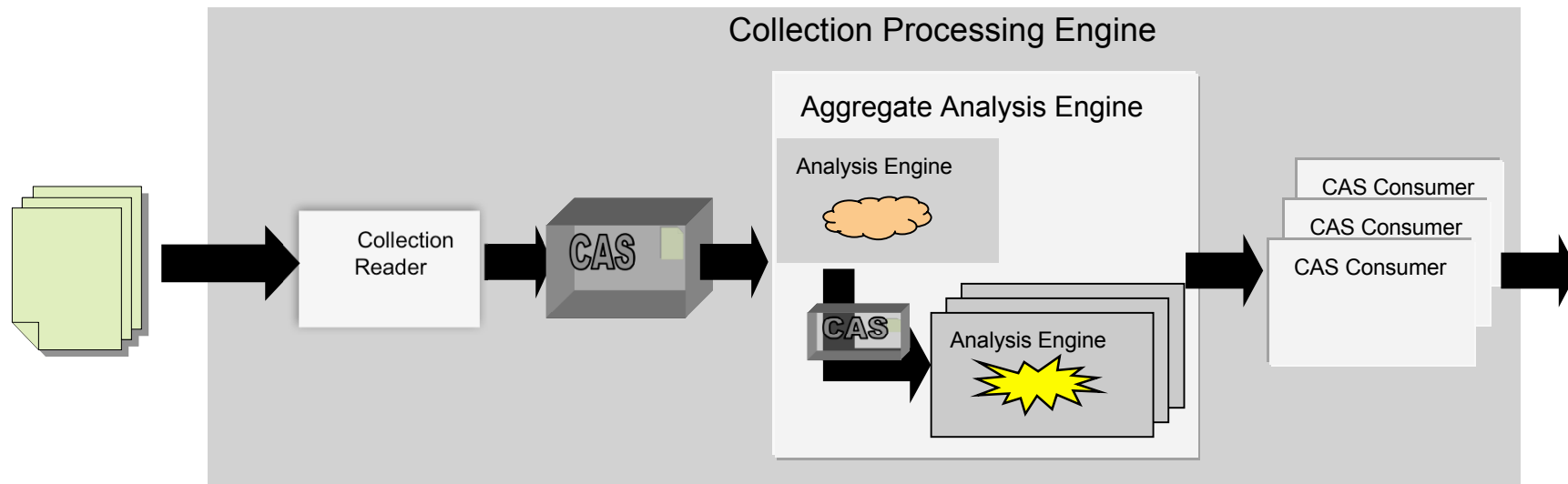
  - Knowledge discovery

## UIMA Concept

- **Analytics**

  – Specialized

  – Independently developed

  – Using different techniques

- **Combination Hypothesis**

  – Combination of different techniques (even based on mutually exclusive premises) produces better results

- **Issue of integration/interoperation**
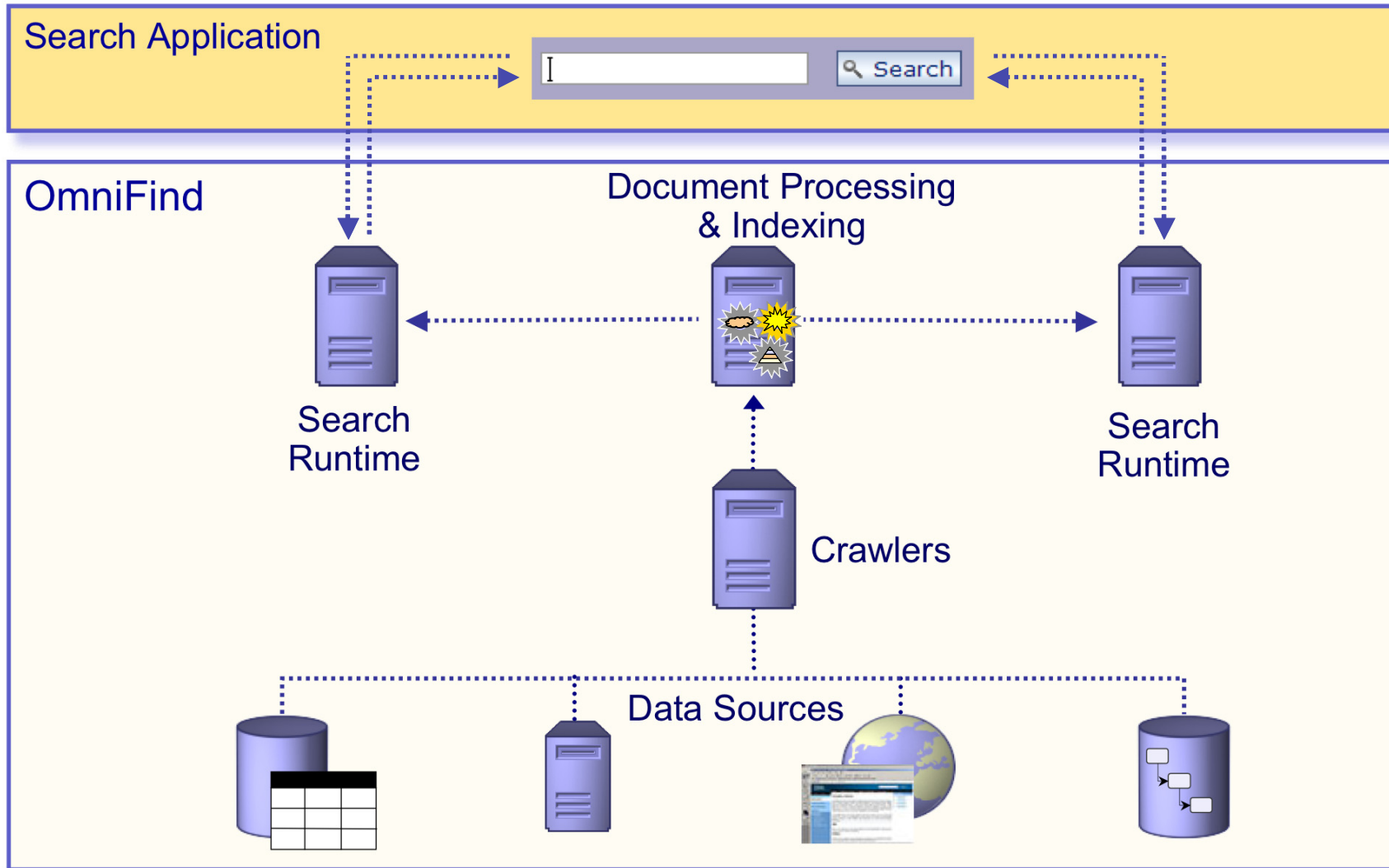
# Combination and Integration

## UIMA Concept

- **CAS (Common Annotation Structure)**

  – Container that is passed between analysis engines

  – Allows creating and consuming annotations in a standard way

# IBM Text Analytics: UIMA and LanguageWare Platform

- **UIMA: Component Software Architecture**

  – Specifies component interfaces, design patterns, data representations

- **UIMA SDK**

  – Java (C++ , Python)

  – Available from IBM AlphaWorks site (www.alphaworks.ibm.com)

- **IBM LanguageWare: multi-lingual NLP technology**

  – Provides tools for text analytics

  – Provides base UIMA text annotator

# UIMA Application in Information Retrieval: OmniFind
(WebSphere Information Integrator OmniFind Edition)

## Text Analytics in Life Sciences

- Unstructured data is major part of all available data

- Text Analytics

  – Annotation

    • Leveraging domain language

  – Relational information

    • Co-occurrence

    • Connections between nouns, through verbs

    • Sentence parsing

- Moving past Information Retrieval

# Domain Knowledge

- Domain knowledge is very important

- Domain customization

  – Any information that describes the language of a domain can prove extremely valuable in customizing the analysis

    • Vocabulary, terminology, relationships, spelling variants, abbreviations, prefixes & suffixes, rules & regular expressions

  – Semantics

- Using existing data

- IBM LanguageWare Workbench

  – Semi-automatic generation of UIMA annotators

  – Java/Eclipse based

## UIMA and Domain Customization

# Examples of UIMA Applications

## Clinical Trials Search

# Examples of UIMA Applications

## Chemical search

## The Problem: Search Patents for Chemical Structures

- **Chemicals have multiple names, trivial and official**

**Names for Valium**

ALBORAL, ALISEUM, ALUPRAM , AMIPROL ,ANSIOLIN , ANSIOLISINA , APAURIN, APOZEPAM, ASSIVAL , ATENSINE , ATILEN , BIALZEPAM , CALMOCITENE, CALMPOSE , CERCINE, CEREGULART, CONDITION, DAP, DIACEPAN, DIAPAM , DIAZEMULS , DIAZEPAM , DIAZETARD , DIENPAX, DIPAM , DIPEZONA, DOMALIUM , DUKSEN, DUXEN, E-PAM, ERIDAN, EVACALM, FAUSTAN, FREUDAL , FRUSTAN, GIHITAN, HORIZON, KIATRIUM, LA-III , LEMBROL, LEVIUM, LIBERETAS , METHYL DIAZEPINONE, MOROSAN , NEUROLYTRIL, NOAN, NSC-77518, PACITRAN, PARANTEN, PAXATE, PAXEL, PLIDAN, QUETINIL, QUIATRIL, QUIEVITA, RELAMINAL, RELANIUM, RELAX, RENBORIN, RO 5-2807, S.A.R.L., SAROMET, SEDAPAM, SEDIPAM, SEDUKSEN, SEDUXEN, SERENACK, SERENAMIN, SERENZIN, SETONIL, SIBAZON, SONACON, STESOLID, STESOLIN, TENSOPAM, TRANIMUL, TRANQDYN, TRANQUASE, TRANQUIRIT, TRANQUO-TABLINEN , UMBRIUM, UNISEDIL, USEMPAX AP, VALEO, VALITRAN, VALRELEASE, VATRAN, VELIUM, VIVAL, VIVOL, WY-3467

- Straightforward text search does not work

- Synonym expansion is insufficient

- Searching by structure is required

- **There is no explicit list of all organic chemicals**

## Chemical Annotator

- **Accurately recognizes organic chemicals in text**

    – Uses a small number of common chemical morphemes (fragments)

    • Names, prefixes, suffixes, endings, etc.

    – Uses pattern rather than a dictionary approach

    • Rules for accepting and combining fragments

    – Not restricted to "known chemicals"

o-**vinylbenzyl glyc**id**yl ether**

## Chemical Search Application

- **To create chemical search application**
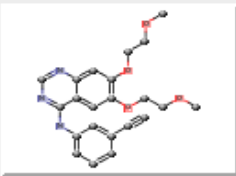
  - Taken >4M US patents

  - Extracted structured information

    - Title, authors, assignee, date, abstract, brief description, full description, claims…

  - Extracted chemical names

  - Converted chemical names into SMILES and InChI strings

  - Indexed

# Chemical Search Application

Query time: 6405.0ms

Graph Results   Tree   View Patents   Claim Analysis

Click box to draw a compound:



Powered by ChemAxon Marvin

Or enter a SMILES:

Or enter a InChI:

Search

**Similarity: 1.0**
**Patents: 98**
Name: 4-[(3-ethynyl-phenyl)amino]-6,7-bis-(2-methoxy-ethoxy)-quinazoline,
Synonyms: 15
⊟ Click to View Synonyms

    4-[(3-ethynylphenyl)amino]-6,7-bis-(2-methoxyethoxy)quinazoline
    n-(3-ethynylphenyl)-6,7-bis (2-methoxyethoxy)-4-quinazolinamine
    n-(3-ethynylphenyl)-6,7-bis (2-methoxyethoxy)-4-quinazolinamine,
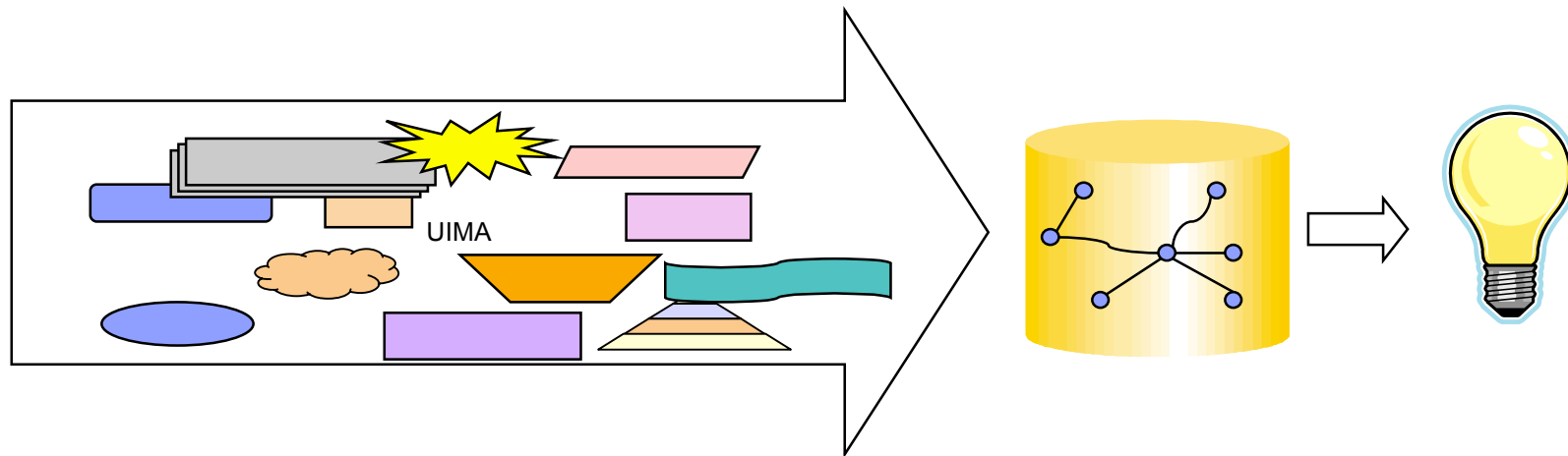    n-(3-ethynylphenyl)-6,7-bis(2-methoxyethoxy)-4-quinazolinamine
    4-[(3-ethynylphenyl)amino]-6,7-bis-(2-methoxy-ethoxy)-quinazoline,
    n-(3-ethynylphenyl)-6,7-bis(2-methoxyethoxy)-4-quinazolinamine,
    n-(3-ethynylphenyl)-6,7bis(2-methoxyethoxy)-4-quinazolinamine
    [6,7-bis(2-methoxy-ethoxy)-quinazolin-4-yl]-(3-ethynylphenyl)amine,
    [6,7-bis(2-methoxyethoxy)quinazolin-4-yl]-(3-ethynylphenyl)-amine
    [6-,7-bis-(2-methoxyethoxy)-quinazolin-4-yl]-(3-ethynylphenyl)amine
    n-(3-ethynylphenyl)-6,7-bis(2-methoxyethoxy)4-quinazolinamine
    4-[(3-ethynyl-phenyl)amino]-6,7-bis-(2-methoxy-ethoxy)-quinazoline,
    4-[(3-ethynyl-phenyl)amino]-6,7-bis-(2-methoxy-ethoxy)-quinazoline
    4-(3-ethynylphenylamino)-6,7-bis(2-methoxyethoxy)quinazoline,
    4-[(3-ethynyl-phenyl)amino]-6,7-bis-(2-methoxy-ethoxy)-quinazoline,

**Similarity: 0.963**
**Patents: 34**
Name: [6,7-bis(2-methoxyethoxy)quinazolin-4-yl]-(3-ethynylphenyl)-amine hydrochloride
Synonyms: 18
⊞ Click to View Synonyms

**Similarity: 0.962**
**Patents: 3**
Name: n-(3-ethylphenyl)-6,7-bis-(2-methoxyethoxy)-4-quinazolinamine,
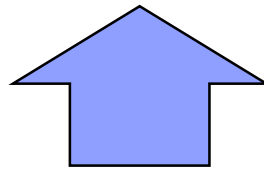Synonyms: 1
⊞ Click to View Synonyms

# Examples of UIMA Applications
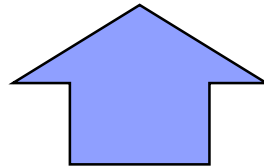
## Knowledge discovery

## Text Mining and Knowledge Discovery

Discover Knowledge

Organize Documents

Search Documents
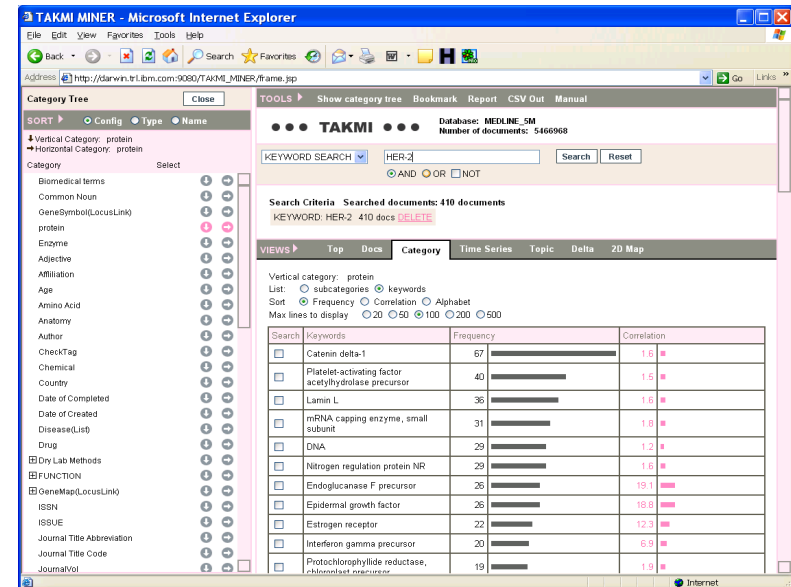
# Text Mining and Knowledge Discovery

- **TAKMI, BIW**

## Stack of UIMA components

- NLP
- Domain-specific lexical and semantic data

## Collection-level analysis

## Interactive (selection-level) analysis

- Relationship mining
- Trend analysis
- Visualization

## Conclusions

- **UIMA framework facilitates integration and interoperation of text analysis modules**

  – Essential for large scale development

    • Allows clear requirements specification

    • Easy to debug and verify

  – Useful for research cooperation

    • Simplifies integration

    • Simplifies deployment

  – Creates development community

- **UIMA simplifies development and deployment of Knowledge Discovery applications**

  – Industrial KD applications benefit most from componentized process

- **The presented works are by**

  - IBM T.J. Watson Research Center

  - IBM Almaden Research Center

  - IBM Tokyo Research Laboratory

  - IBM Software Group

- **Clinical Trials Portal developed in cooperation with IFPMA**