
The International Conference for Science & Business Information 2009

Chemical Depictions – The Grand Challenge in Patents

Dr. Marc Zimmermann

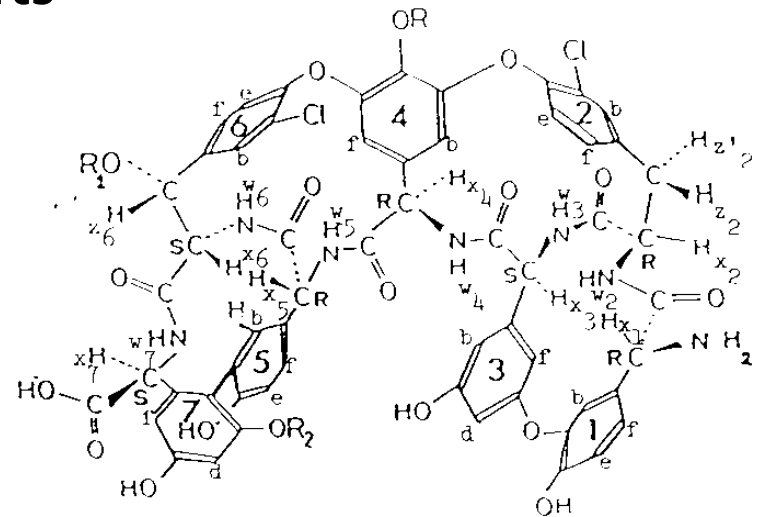
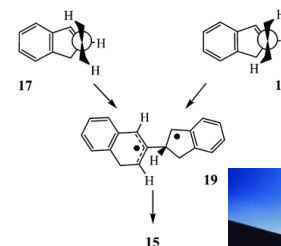


Image Extraction – an old Problem..?

- Short Overview on Chemical Structure Reconstruction
- CSR in Patents
 - Page Segmentation + Image Classification
 - Error Detection + Prediction
 - Markush Structures
 - Content Production Pipeline
- Conclusions



Archaeology



Chemistry / Pharmacology

Overview: Fraunhofer SCAI.Bio

Information extraction in the life sciences

- Biology
- Chemistry
- Medicine
- Pharmacology

Working on

- Textual sources (PubMed, Journals)
- Image sources (Journals, Thesis)

Techniques

- Dictionary based approaches (ProMiner)
- Machine learning (Linda)
- Expert systems (chemoCR)

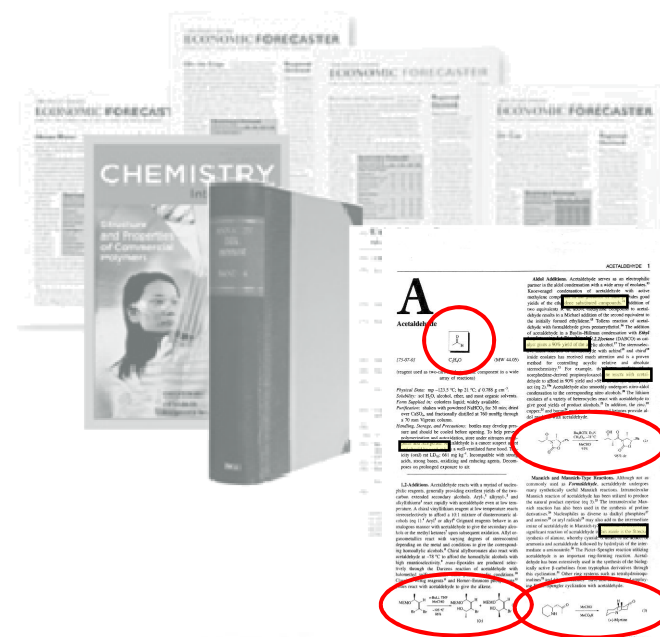


Image Reconstruction is a famous Problem

- **Kekulé: OCR-optical chemical (structure) recognition**, R. McDaniel and Jason R. Balmuth. *J. Chem. Inf. Comput. Sci.*, 32(4):373–378, 1992.
- **Chemical Literature Data Extraction: The CLiDE Project**, P. Ibison, M. Jacquot, F. Kam, A. G. Neville, R.W. Simpson, C. Tonnelier, T. Venczel and A.P. Johnson, *J. Chem. Inf. Comput. Sci.*, vol. 33(3): 338-344, 1993.
- **Apparatus and method for optical recognition of chemical graphics**, R. Casey, S. Boyer, B. Oudot and K. Zilles, US Patent 5157736, 1992. **OROCs**
- **Optical recognition of chemical graphics**, S. Boyer, Document Analysis and Recognition, Proceedings of the Second International Conference on Document Analysis, 627–631, 20-22 Oct 1993.

But what happened to them?



New Initiatives

- **Optical structure recognition software to recover chemical information: OSRA**, V.I. Filippov and C.M. Nicklaus, J. Chem. Inf. Comput. Sci., 39(3):740-743, 2009.
- **CLiDE Pro: The latest generation of CLiDE**, T.A. Valko and P.A. Johnson, J. Chem. Inf. Comput. Sci., 49(4):780-787, 2009.
- **Automated extraction of chemical structure information from digital raster images**, J. Park, R.G. Rosania, A.K. Shedden, M. Nguyen, N. Lyu and K. Saitou, Chemistry Central Journal, 3:1-16, 2009.
ChemReader
- **Machine learning and expert knowledge based evaluation of the chemoCR reconstruction result**, O. Domanova, Master Thesis, Bonn-Aachen International Center for Information Technology, University Bonn, 2009.

„Advanced Patent Search“

- Keywords (full text, title, abstract, claims)
 - Publication / Application Number
 - Applicant(s) / Inventor(s)
 - Classification (ECLA, IPC)
 - Date (filing, issue)
-
- WIPO
 - USPTO
 - EPO
 - GOOGLE Patent Search

WORLD INTELLECTUAL PROPERTY ORGANIZATION

IP SERVICES

ABOUT WIPO IP SERVICES PROGRAM ACTIVITIES RESOURCES NEWS & EVENTS

Home > IP Services > PATENTSCOPE > Patent Search

PATENTSCOPE®
Search International Patent Applications

This facility allows you to search 1,636,382 international patent applications and to view the latest [information and documents](#) available to the International Bureau.

Structured Search options results

» Keywords Front Page

AND Publication Number

AND Application Number

AND Publication Date

AND English Title

AND English Abstract

AND Applicant Name

AND Int. Class

AND Inventor Name

AND National Phase Country

AND Description

AND Claims

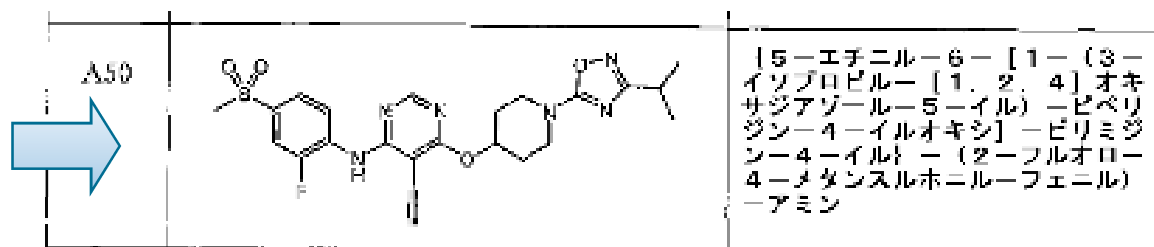
Search

Advanced Searching from the Perspective of a Chemist

What is really interesting is – find me documents containing...

- similar structures
- structures containing a polycyclic aromatic ring system
- structures fulfilling my pharmacophore
- core structure in Asian patents
- a synthesis protocol
- biological target information
- ...

graphical represented
structure information



Idea: SCAIView Patent Retrieval

SCAIView P7b build 5381 - Mozilla Firefox

http://lab-7.scai.fhg.de:8080/scaiview/

SCAIVIEW

Submit Search

Documents Entity Analysis

cancer

Result for 'cancer', NER run 'Linda' for entity **to mapping: IUP3379**, Page 0 with 1 documents per page, totals to 1 and took 10 ms.

Toggle Abstracts Select All Entity Classes Deselect All Entity Classes


Use Document Base: TRETest

- Drug Names
- Human Protein / Gene
- Non Normalized SNPs
- SNP2PubMed
- Gene2PubMed
- SNPs AltGene
- SNPs PDGene
- SNPs SZGene
- MeSH Disease
- Corpora
- ChEBI
- IUPAC
- IPC

- + Human Genes / Proteins
- + DrugBank
- Corpora
- Normalized SNPs
- SNP2PubMed
- Gene2PubMed
- SNPs PDGene
- SNPs AltGene
- SNPs SZGene
- HuGeNet Genetic Associations
- MeSH Disease
- ChEBI
- IUPAC
- International Patent Class

1. 3'-azido-2',3'-dideoxyuridine administration to treat HIV and related test protocol

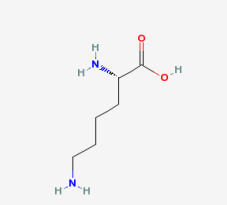
ID: US-6194391-B1



Statistics

It has been discovered that 3'-azido-2',3'-dideoxyuridine (CS-87) induces a transient mutation in HIV-1 at the 70th codon (K to R, i.e., lysine to arginine) in the reverse transcriptase region of the virus. Based on this discovery, a method is provided for treating HIV-1 that includes administering CS-87 or its pharmaceutically acceptable salt or prodrug to a human in need of a drug that induces a mutation in HIV-1 at a location other than the 70th codon of the reverse transcriptase region. This invention can be practiced by referring to the published mutation pattern for a new drug.

We claim: 1. A method for treating a human infected with HIV comprising administering CS-87 (3'-azido-2',3'-dideoxyuridine) or its pharmaceutically acceptable salt or prodrug in alternation or combination with a second anti-HIV drug that induces a mutation in HIV at a location other than the 70th codon of the reverse transcriptase region of HIV. 2. The method of claim 1, wherein the second anti-HIV drug is a protease inhibitor. 3. The method of claim 1, wherein the protease inhibitor is indinavir ((1S, 2R) 5(S)-2,3,5-trideoxy-N-(2,3-dihydro-2-hydroxy-1H-inden-1-yl)-5-(2-((1,1-dimethylethylamino)carbonyl)-4-(3-pyridinylmethyl)-1-piperazinyl)-2-(phenylmethyl)-D-erythro-pantoamide sulfate), nelfinavir, saquinavir, ritonavir, ABT-378 (N-(4(S)-(2-(2,6-dimethylphenoxy)-acetylamino)-3(S)-hydroxy-5-phenyl-1(S)-benzylperityl)-3-methyl-2-(S)-(2-oxo(1,3-diazaperhydroxyl)butanamine)), or amprenavir. 4. The method of claim 1, wherein the second anti-HIV drug is a nucleoside analogue or a nucleotide analogue. 5. The method of claim 1, wherein the nucleoside analogue is FTC (bis-2',3'-dideoxy-3'-thia-5-fluorocytidine), 3TC (beta-L-3'-thia-2',3'-dideoxycytidine), DDI (2',3'-dideoxyinosine), dd4FC (2',3'-dideoxy-2',3'-didehydro-5-fluorocytidine), DAPD (beta-D-ribofuranosyl-2,6-diaminopurine), FddA (2',3'-dideoxy-2'-beta-fluoro-adenosine), or abacavir ((1S, 4R)-4-(2-amino-6-(cyclopropylamino)-9H-purin-9-yl)-2-cyclopentene-1-methanol). 6. The method of claim 1, wherein the nucleotide analogue is adefovir (bis(pivaloyloxymethyl)-9-(2-phosphonylmethoxyethyl)adenine) or PMPA ((R)-9-(2-phosphonyl-methoxypropyl)adenine). 4.7. The method of claim 1, wherein the second anti-HIV drug is a non-nucleoside reverse transcriptase inhibitor. 8. The method of claim 1, wherein the non-nucleoside reverse transcriptase inhibitor is efavirenz, nevirapine, MKC-442 (1-ethyl-5-methyl-3-isoxazol-5-yl-5-benzyluracil), AG-1549 (Cagavirine: 5-(3,5-dichlorophenylthio)-4-isopropyl-1-(4-pyridyl)methyl-1H-imidazo[2,1-b]imidazole-2-ylnitrile carbamate), or rilpivirine. 9. The method of claim 1, wherein the second anti-HIV drug is an HIV-integrase inhibitor or a chemokine inhibitor or a fusion inhibitor. 10. A method for assessing the sensitivity of HIV-1 to CS-87 in a patient to whom CS-87 has been administered, comprising isolating a sample of HIV-1 from the patient and identifying whether a mutation has occurred at codon 70 in the reverse transcriptase region of the virus. 11. A method for treating a patient infected with a strain of HIV that is resistant to treatment with an anti-HIV agent selected from the group consisting of 3TC (beta-L-3'-thia-2',3'-dideoxycytidine), DDI (2',3'-dideoxyinosine), DDC (2',3'-dideoxycytidine), [(+)- and racemic] FTC (bis-2',3'-dideoxy-3'-thia-5-fluorocytidine), dd4FC (2',3'-dideoxy-2',3'-didehydro-5-fluorocytidine), DAPD (beta-D-ribofuranosyl-2,6-diaminopurine), FddA (2',3'-dideoxy-2'-beta-fluoro-adenosine), or abacavir. 12. The method of claim 11, wherein the anti-HIV agent is Lamivudine: Lamivudine is a nucleoside reverse transcriptase inhibitor (NRTI) with activity against Human Immunodeficiency Virus Type 1 (HIV-1) and hepatitis B (HBV), with confidence 5 0



DrugBank
swissprot
swissprot

Fertig

Goal: Multi-modal Information Extraction

US 2005/0182053A1

(19) United States
 (12) Patent Application Publication (10) Pub. No.: US 2005/0182053 A1
 Chen et al. (43) Pub. Date: Aug. 18, 2005

(54) SUBSTITUTED 3-AMINO-THIENO[2,3-B]PYRIDINE-2-CARBOXYLIC ACID AMIDE COMPOUNDS AND PROCESSES FOR PREPARING AND THEIR USES

(75) Inventors: Zhidong Chen, New Milford, CT (US); Pier Francesco Cirillo, Woodbury, CT (US); Darren DiSalvo, New Milford, CT (US); Weimin Liu, Sandy Hook, CT (US); Daniel Richard Marshall, Sandy Hook, CT (US); Lifan Wu, New Milford, CT (US); Erick Richard Roush Young, Danbury, CT (US)

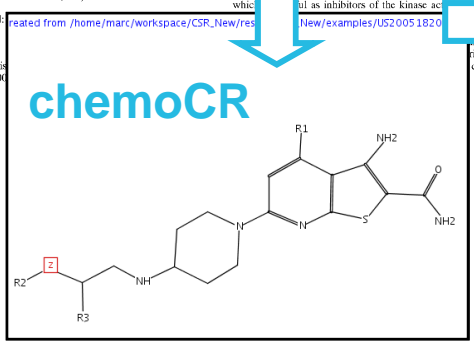
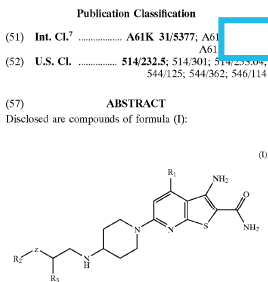
Correspondence Address: MICHAEL P. MORRIS BOEHRINGER INGELHEIM CORPORATION 900 RIDGEBURY ROAD P O BOX 368 RIDGEBURY, CT 06877-0368 (US)

(73) Assignee: Boehringer Ingelheim Pharmaceuticals, Inc., Ridgefield, CT

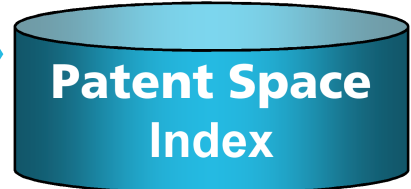
(21) Appl. No.: 11/002,828

(22) Filed: [redacted]

(60) Previous publications: [redacted]



ProMiner (NER)



3-AMINO-THIENO[2,3-BIPYRIDINE-2-CARBOXYLIC ACID AMIDE PROCESSES FOR PREPARING AND THEIR USES

RELATED APPLICATIONS

[0001] This application claims priority to U.S. patent application No. 601527,522 filed on Dec. 5, 2003. This application is also related to U.S. patent application Nos. 101453,175 and 101730,172.

TECHNICAL FIELD OF THE INVENTION

[0002] This invention relates to substituted 3-aminothieno[2,3-b]pyridine-2-carboxylic acid amide compounds useful as inhibitors of the kinase activity of the IKK (IKK) complex. The compounds are therefore useful for treatment of IKK-mediated diseases including autoimmune diseases, inflammatory diseases and cancer. The invention also relates to processes for preparing such compounds and pharmaceutical compositions comprising them.

BACKGROUND OF THE INVENTION

[0003] NF-KB or nuclear factor KB is a transcription factor that induces the expression of a large number of pro-inflammatory and anti-apoptotic genes. These include cytokines such as IL-1, IL-2, TNF- α and IL-6, chemokines including IL-8 and RANTES, as well as other pro-inflammatory molecules including COX-2 and cell adhesion molecules such as ICAM-1, VCAM-1, and E-selectin. The NF-KB family includes homo- and heterodimeric transcription factors composed of members of the Rel family (see for example P. A. Baeuerle and D. Baltimore, Cell, 1996,87, 13). Under resting conditions, NF-KB is present in the cytosol of cells as a complex with IKB. The IKB family of proteins

- proteins
- protein families
- protein complex
- compound
- process
- drug class
- disease
- pathways

Chemical Page Segmentation and Image Classification

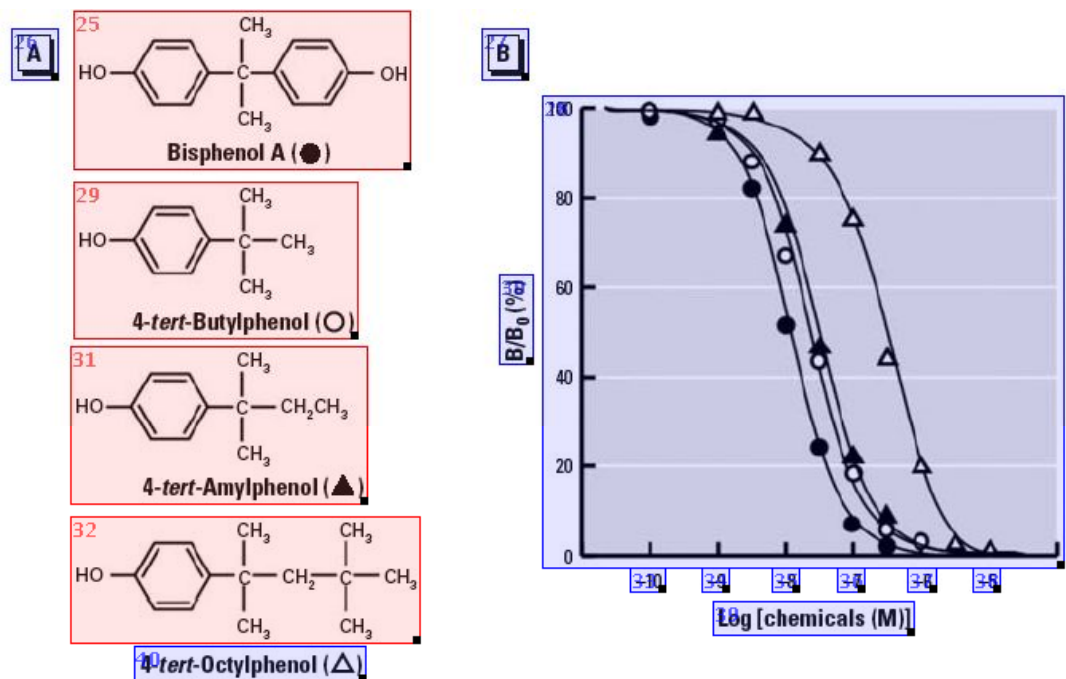
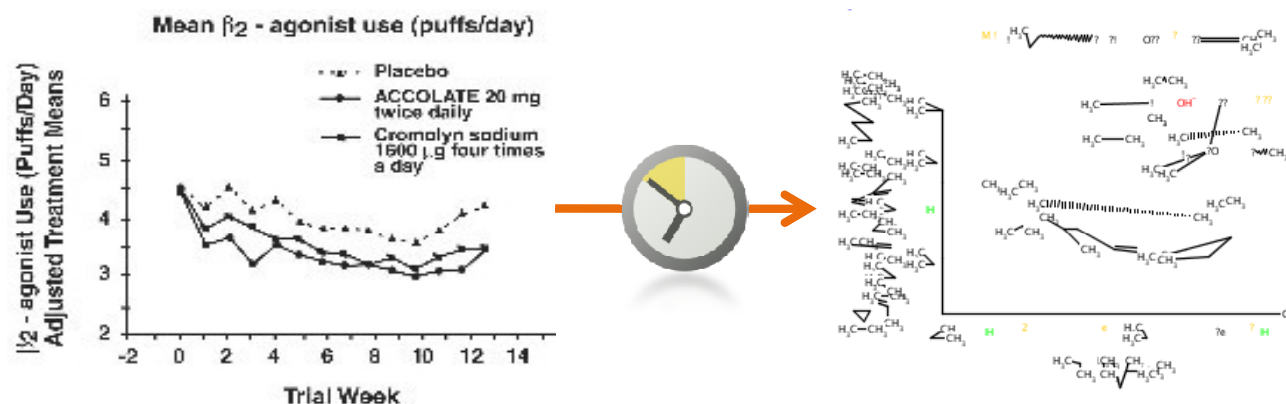


Figure 4. Chemical structure of BPA and its derivatives lacking the phenol group and their dose-response curves in the radioligand receptor binding assay for ERR- γ . (A) Chemical structure of BPA and its derivatives with the alkyl group at the position of phenol group: 4-tert-butylphenol (a methyl group); 4-tert-amylphenol (an ethyl group); and 4-tert-octylphenol (a tert-butyl methyl group). (B) Binding activities of BPA, 4-tert-butylphenol, 4-tert-amylphenol, and 4-tert-octylphenol examined by the competitive binding assay using [³H]BPA and GST-ERR- γ -LBD; representative curves indicate the IC₅₀ value closest to the mean IC₅₀ from at least five independent assays for each compound.

Motivation for Image Classification

- Example for an image without chemical content and it's reconstruction with chemoCR



- Effective reconstruction of big image stacks
- Full page scans: classification of page segments

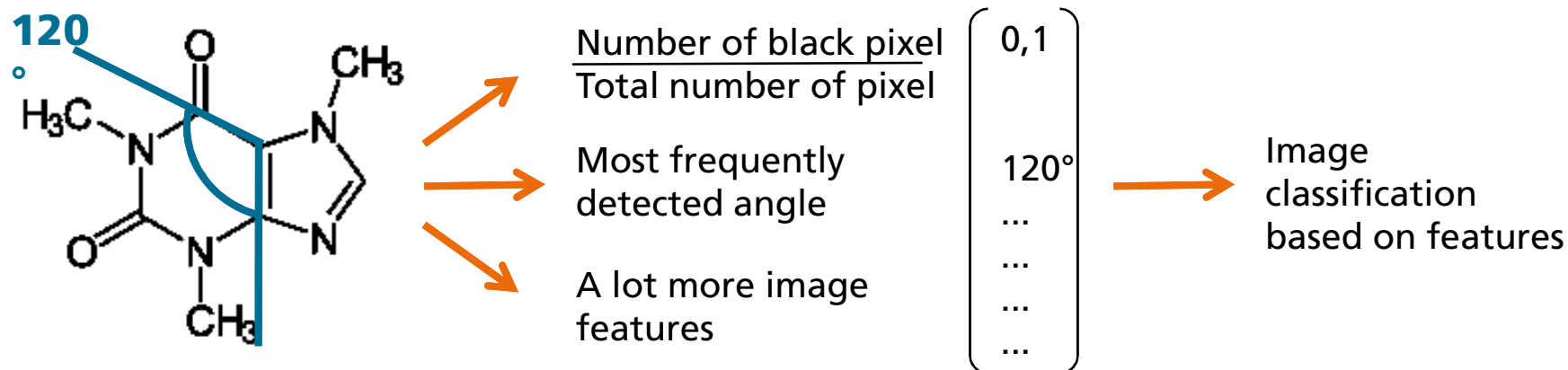
Concept of Image Classification

■ Problem

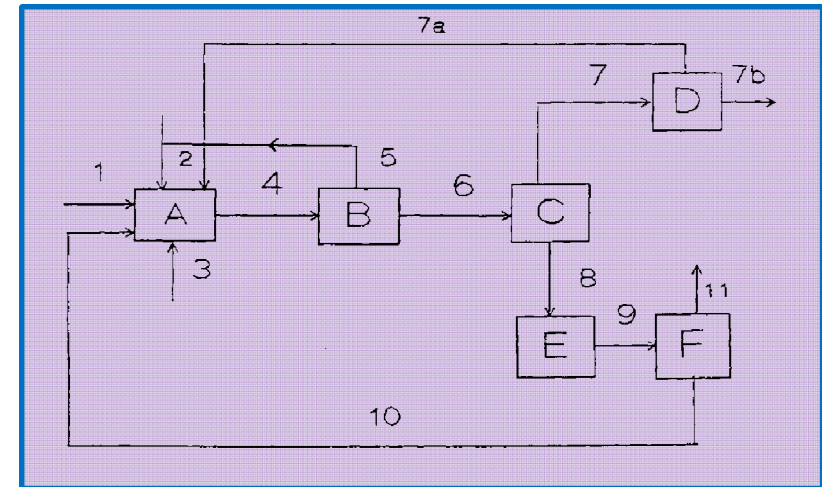
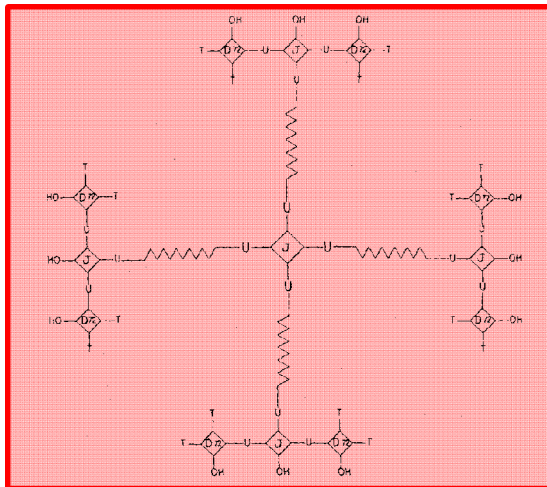
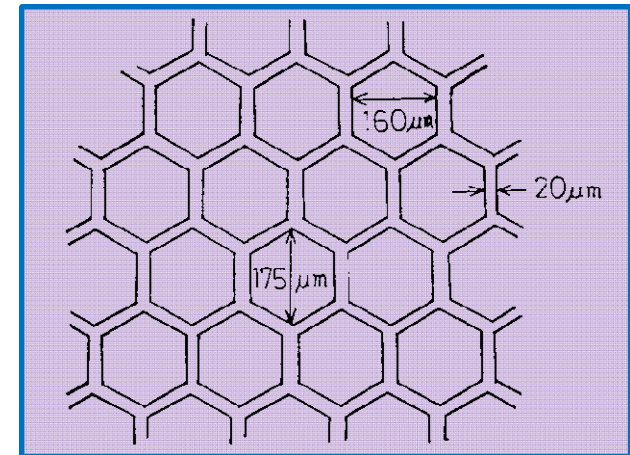
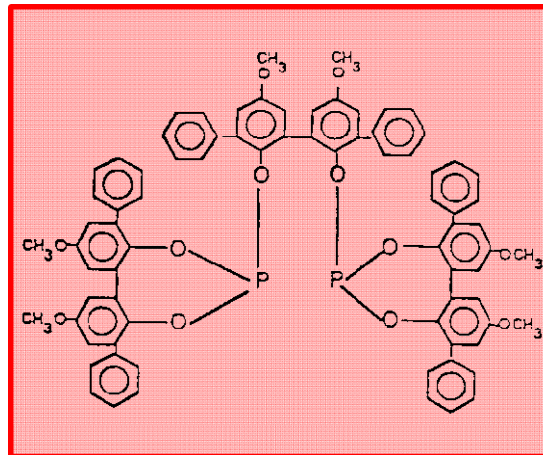
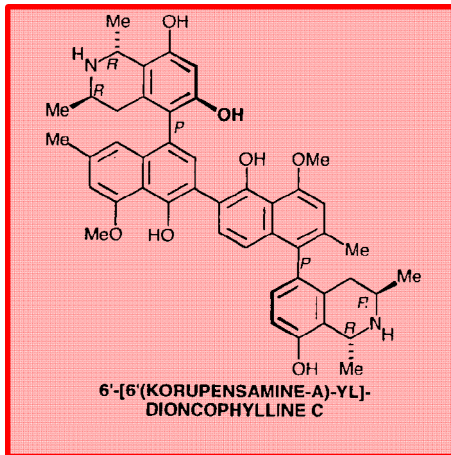
- No well defined rule, whether a image is a structural formula or not.

■ Approach

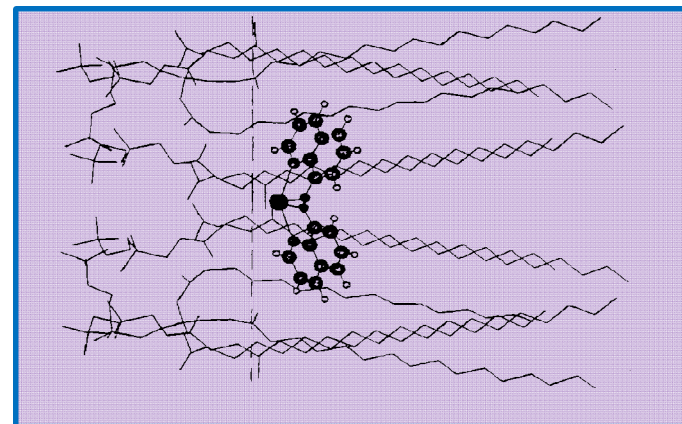
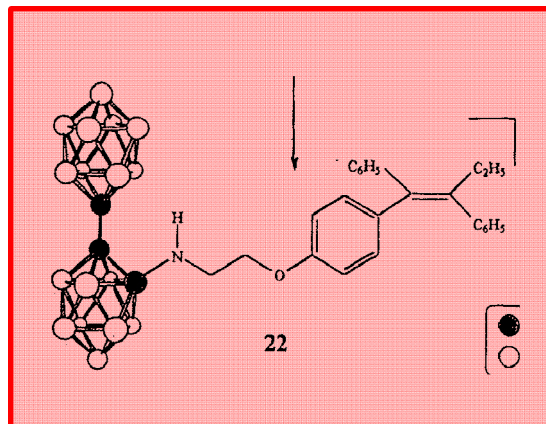
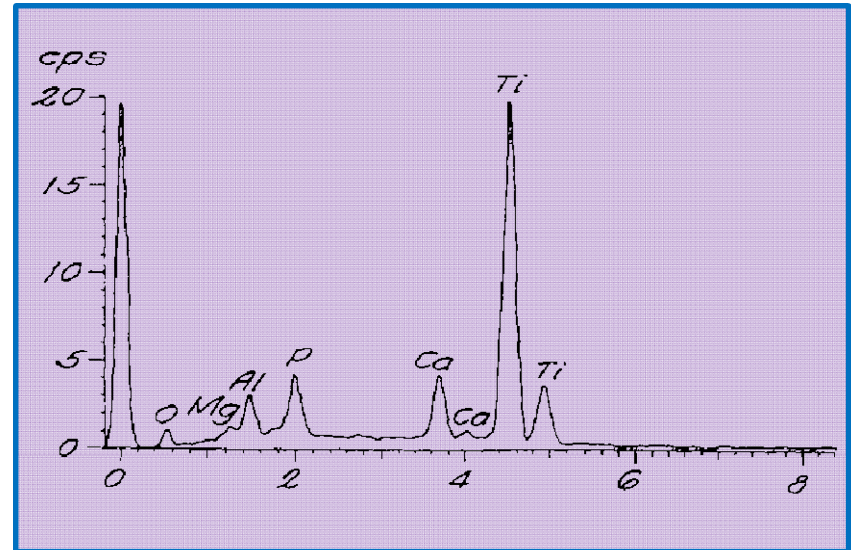
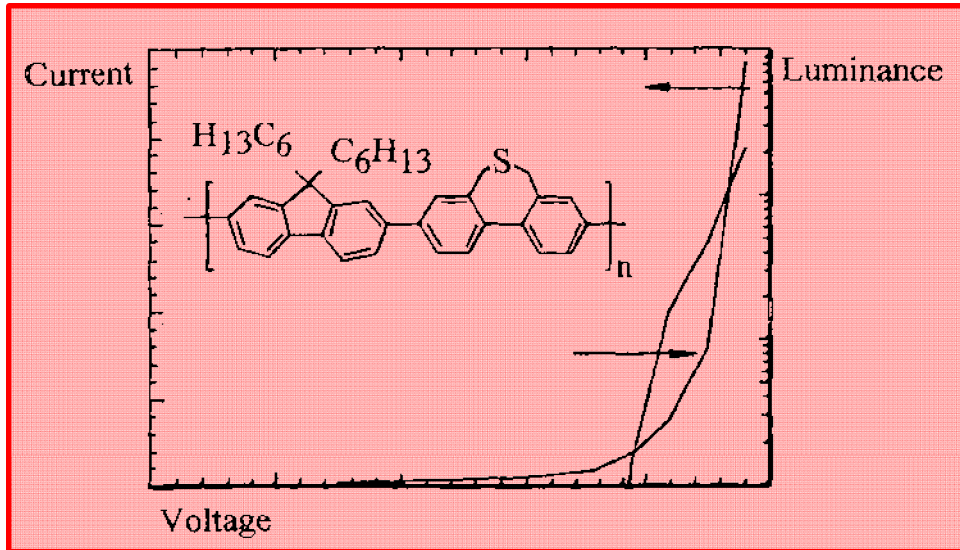
- Calculate comparable **image features** and use machine learning



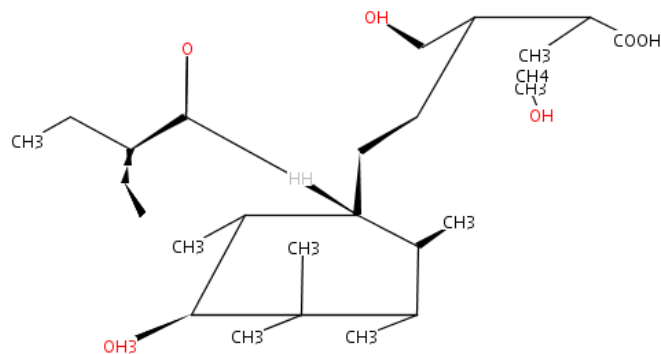
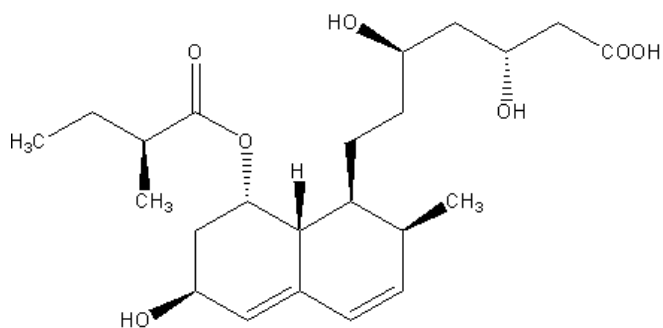
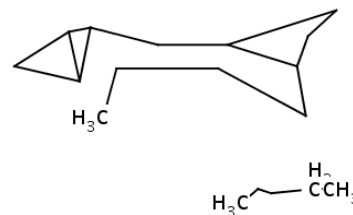
Some Examples from Patents (I)



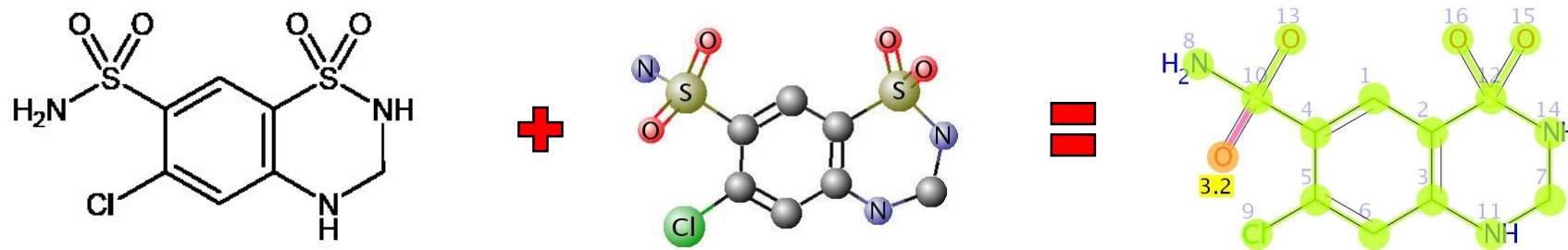
Some Examples from Patents (II)



Beautiful Artwork But Wrong Molecule



Quality Measure: Graph Matching



- SimilarityMCD (Minimal Chemical Distance)

- Module from InfoChem

- Graph-matching on

- Reconstruction result of chemoCR

- The reference molecule

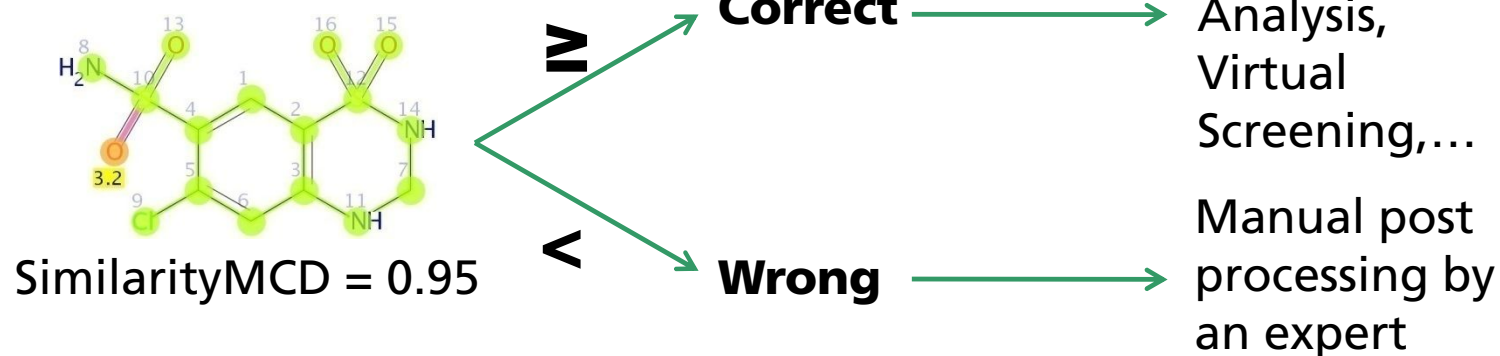
- Results in

- Numerical value, $[0,1]$

OK
bad

Automatic Assessment of the Reconstruction Quality

- SimilarityMCD in [0,1]
 - 1.0 – fully correct reconstruction
- Threshold



- Aim
 - Minimize false positive rate

Challenge

■ Problem

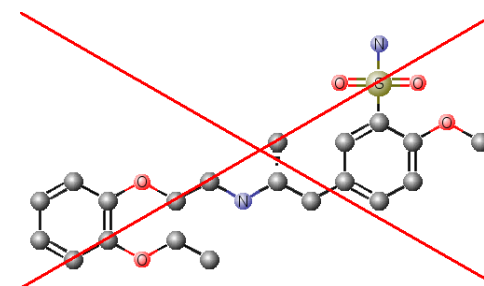
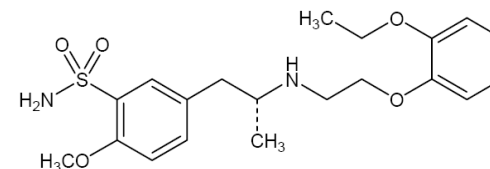
- Predict the quality of the reconstruction result without a reference molecule

■ Solution

- Machine learning

■ Expected results

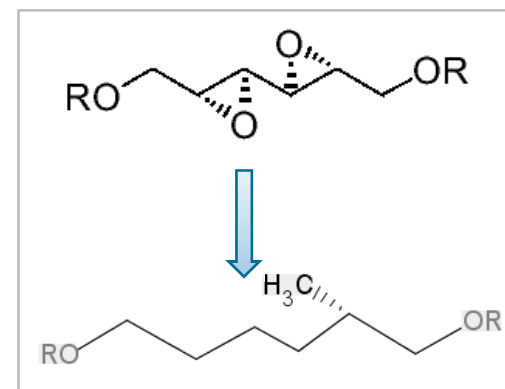
- Quality of new reconstructions estimated by trained models



Chemical Error Classification System

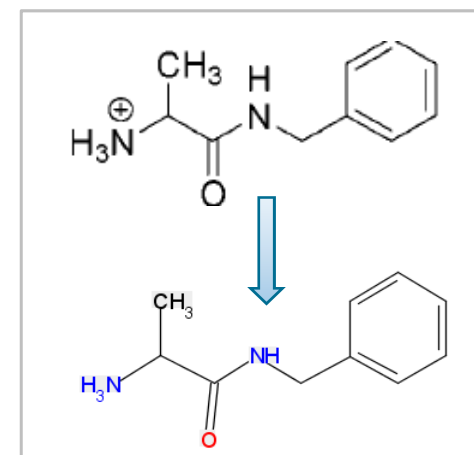
■ MISSED

- BOND_MISSED
 - COMPLETE_BOND_MISSED
 - ORDER_BOND_MISSED
 - **CHIRAL_BOND_MISSED**



■ SYMBOL_MISSED

- ATOM_SYMBOL_MISSED
- ISOTOPE_SYMBOL_MISSED
- **CHARGE_SYMBOL_MISSED**
- RADICAL_SYMBOL_MISSED



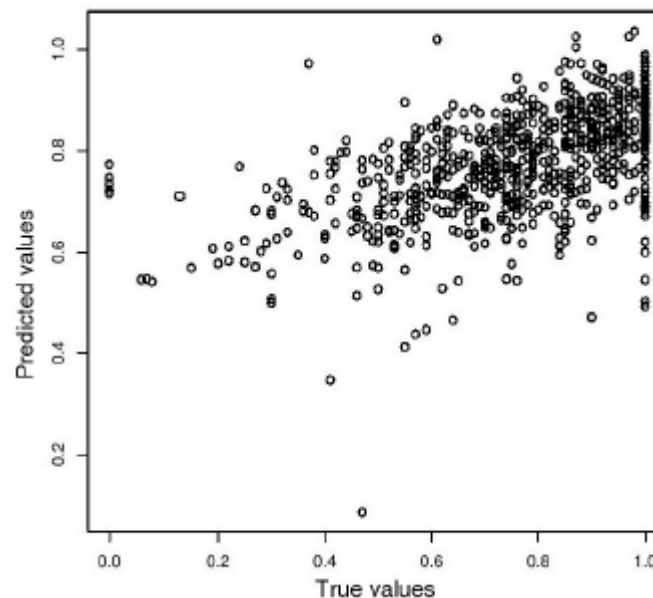
Chemical Error Classification System

- Assigned based on comparison of
 - Result of the reconstruction
 - Reference molecule
- Hierarchical system
- 4 super classes
 - AS_STRUCTURE
 - ERROR
 - MISSED
 - ADDED

1. _AS_STRUCTURE
 - a) GRAPHIC_AS_STRUCTURE
 - b) TEXT_AS_STRUCTURE
 - i. CHAR_TEXT_AS_STRUCTURE
 - ii. NUMBER_TEXT_AS_STRUCTURE
 - c) ARROW_AS_STRUCTURE
 - i. ELECTRON_ARROW_AS_STRUCTURE
 - ii. REACTION_ARROW_AS_STRUCTURE
 - d) BRACKET_AS_STRUCTURE
 - i. POLYMER_BRACKET_AS_STRUCTURE
 - ii. REACTION_BRACKET_AS_STRUCTURE
 - e) NEWMAN_AS_STRUCTURE
 - f) CHARGE_AS_STRUCTURE
2. _ERROR
 - a) KEKULE_RING_ERROR
 - b) BOND_ERROR
 - i. COORD_BOND_ERROR
 - ii. H_BOND_ERROR
 - iii. WAVY_BOND_ERROR
 - iv. CROSS_BOND_ERROR

Some preliminary Results

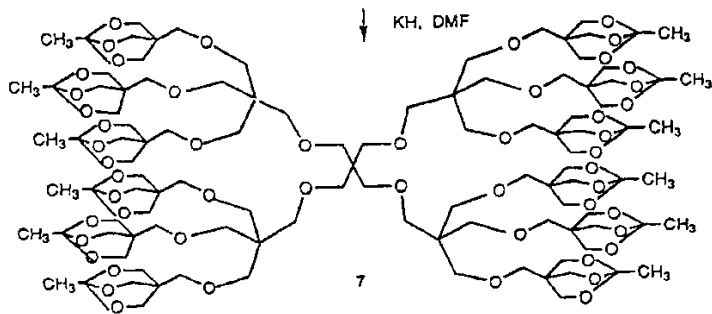
Performance measure	Training data	Test data
MSE *)	2.466×10^{-2}	2.976×10^{-2}
Correlation coefficient	5.816×10^{-1}	5.393×10^{-1}



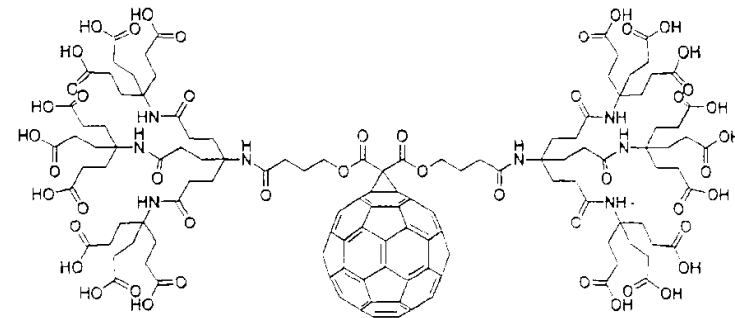
*) mean squared error

Category: Complex Molecules (I)

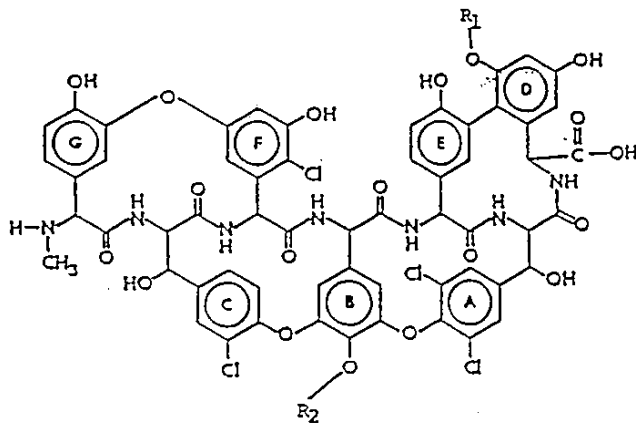
EP 0751952 B1: Displacement Chromatography of Proteins



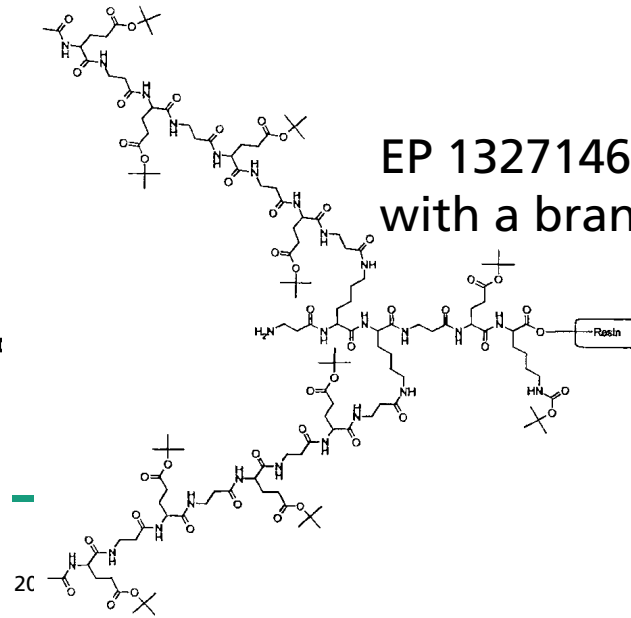
EP 1056475 B1: Dendrimere Fullerene derivate



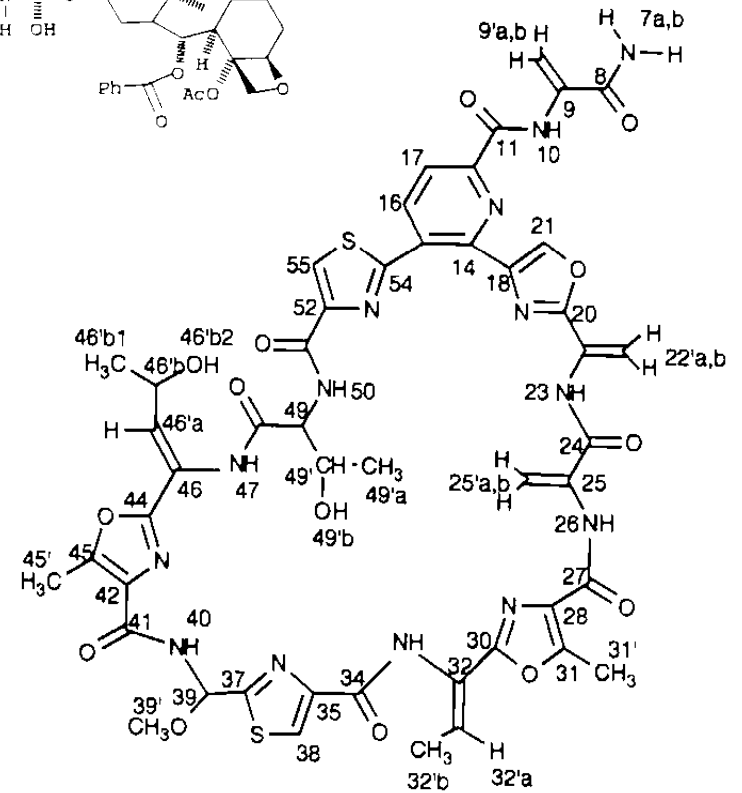
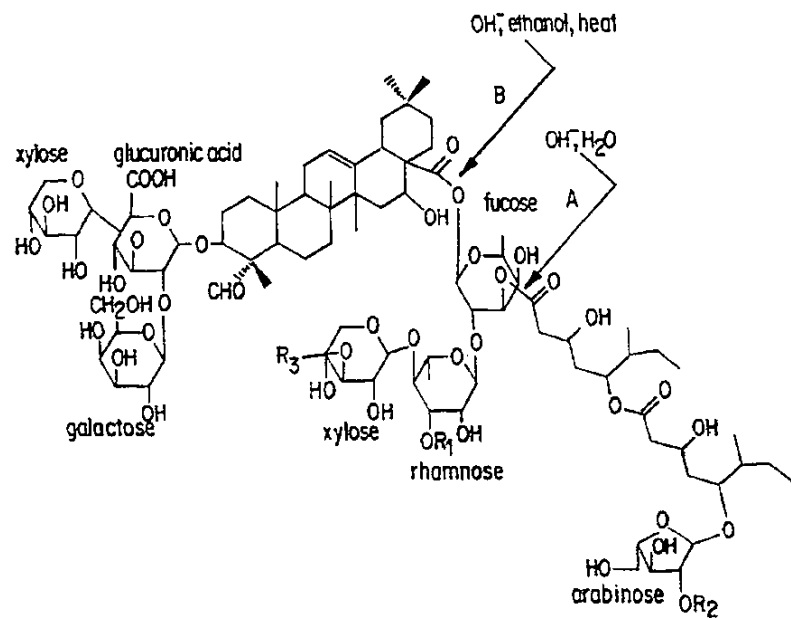
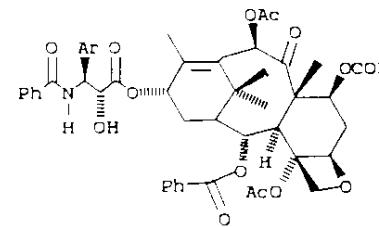
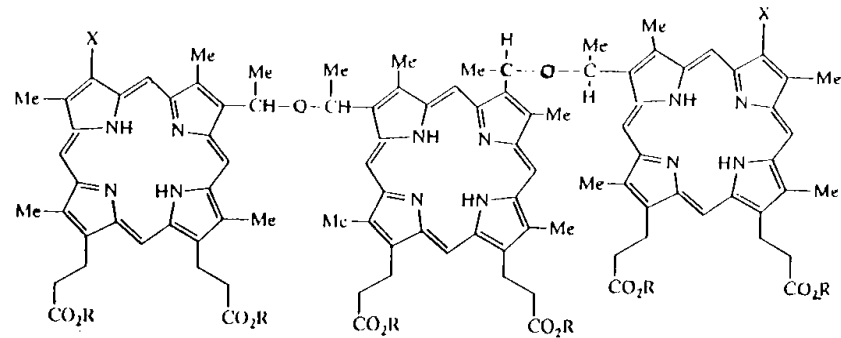
EP 0211490 B1: Antibiotics of the Vancomycin-Class



EP 1327146 B1: Compounds with a branched Linker



Category: Complex Molecules (II)

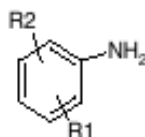


Why is Markush famous?

- Dr. Eugene A. Markush (1887-1968)
- Founder of Pharma Chemical Corporation Bayonne, NJ in 1917
- Did not(!) invent generic structures
- Patent for method of preparing a generic group of pyrazalone dyes for wool or silk



The process for the manufacture of dyes which comprised coupling with a halogen-substituted pyrazolone, a diazotized unsulphonated material prepared from a compound having the formula



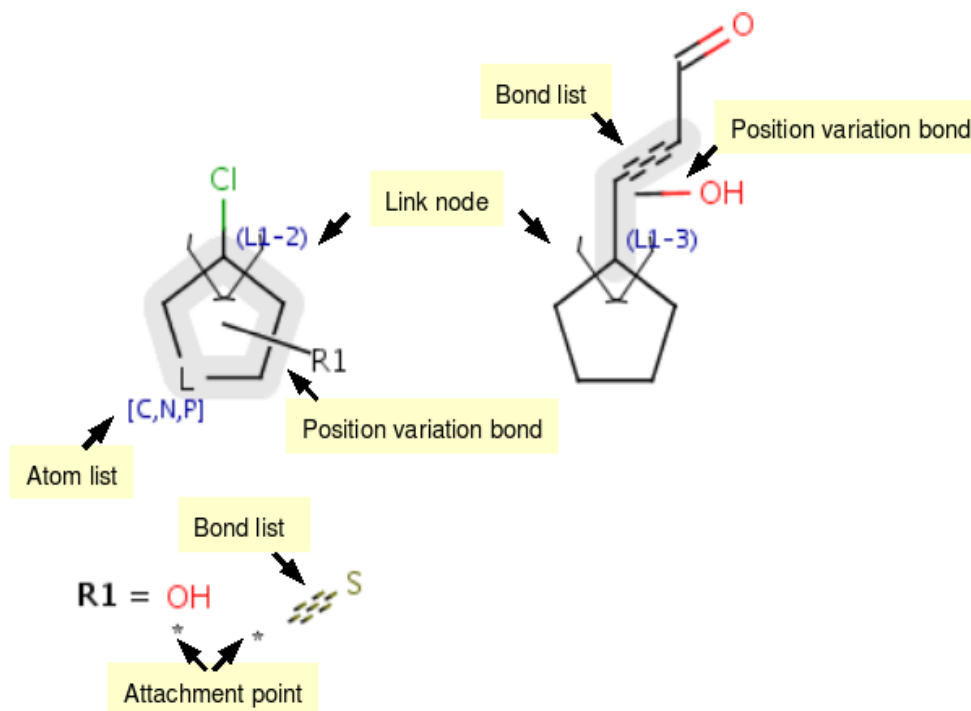
wherein R1 and R2 are independently selected from the

First claim of patent US 1506316

- The claim was challenged as being too unspecific
- Granted in 1924 by the US commissioner of Patents

Markush Structures in Principle

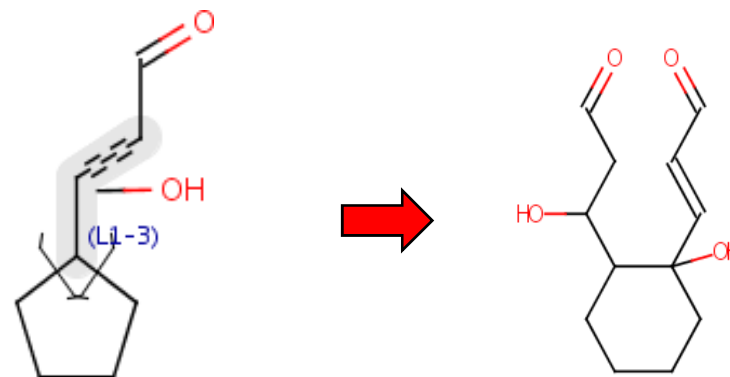
- No general definition
- Core structure
 - Multiple elements
 - Nested elements
- Additional information
 - R-group definition



Markush Structure Analysis

■ Core Structure

- Nested elements Order?
- Link node most powerful
- Creates just valid structures



■ Additional Information

- Structure list
- Formal string
- Structured text
- Natural language

- Structured text → E is selected from the group consisting of CR^9R^9 , CR^9 , O, N and NH;
- Formal strings → G is CH or N;
J is CR_9R_9 , CR_9 , C=O, C=S, or N;
- Natural language → the dashed line represents an optional double bond, provided that J is not C=O, C=S or CR_9R_9 when a double bond is present;
Y is $-NR_3$, O or S;
- Structured text → n_2 and n_3 are each independently 0-2;
 n_4 is 0-3;
-

Textual Annotations

A, B and D are each independently selected to be CR_{4b} or N, provided that at least one A, B or D is N;

E is selected from the group consisting of CR^9R^9 , CR^9 , O, N and NH;

G is CH or N;

J is CR_9R_9 , CR_9 , C=O, C=S, or N;

the dashed line represents an optional double bond, provided that J is not C=O, C=S or CR_9R_9 when a double bond is present;

Y is $-NR_3$, O or S;

n_2 and n_3 are each independently 0-2;

n_4 is 0-3;

R_1 is aryl or heteroaryl, each of which may be optionally substituted with one or more substituents selected from R_4 ;

R-group
Chem. Element
Relation

Wanted: an Editor

The screenshot displays the chemoCR software interface, which is used for chemical structure reconstruction. The main window is titled "chemoCR" and contains a menu bar (File, Workflow, Rule System, View, Help) and a toolbar. The central area is divided into two main panels:

- Left Panel (Loaded Molecule (0)):** This panel shows a chemical structure labeled "Formula I". The structure consists of a fused ring system with a benzene ring (atoms A, B, D) and a five-membered ring (atoms E, J). A nitrogen atom (N) is part of the five-membered ring and is bonded to a group Y (with substituent R₁) and a group G. The group G is further substituted with a ring system containing atoms n₂ and n₃, and a group R₂. Below the structure, the text reads: "I. A compound of Formula I" and "and enantiomers, diastereomers and pharmaceutically acceptable salts thereof wherein:". This is followed by a list of definitions for A, B, D, E, G, J, Y, n₂, n₃, n₄, R₁, and R₂.
- Right Panel (Core):** This panel shows a smaller version of the chemical structure from the left panel. Below the structure is a search bar with "benzene" entered and a dropdown menu showing a benzene ring structure. Below the search bar is a table of "R-Group Definitions" and a table of "R Logic".

R-Group Definitions:

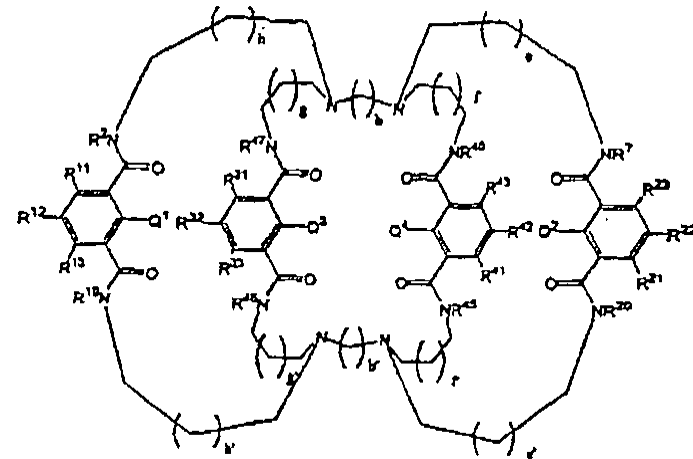
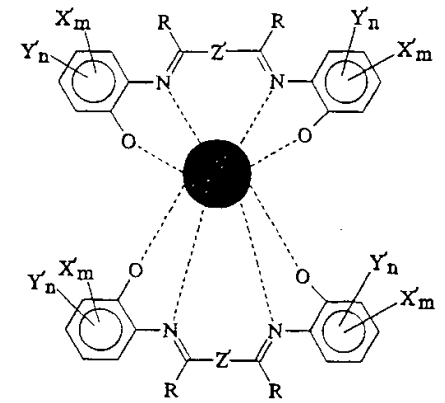
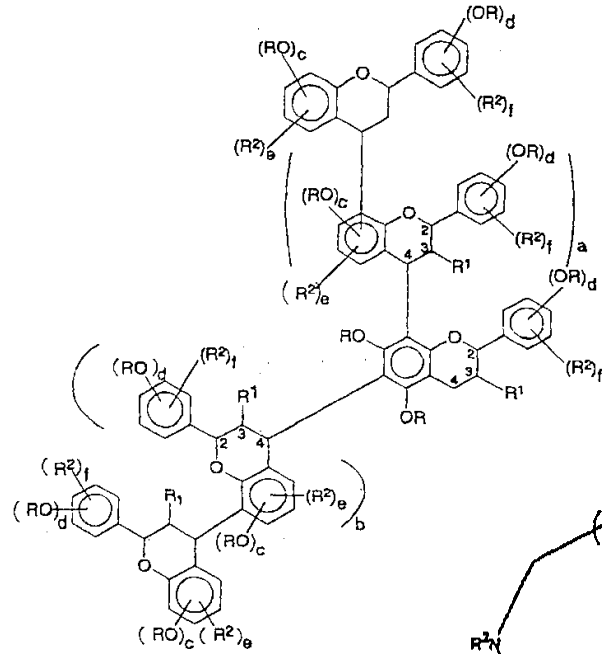
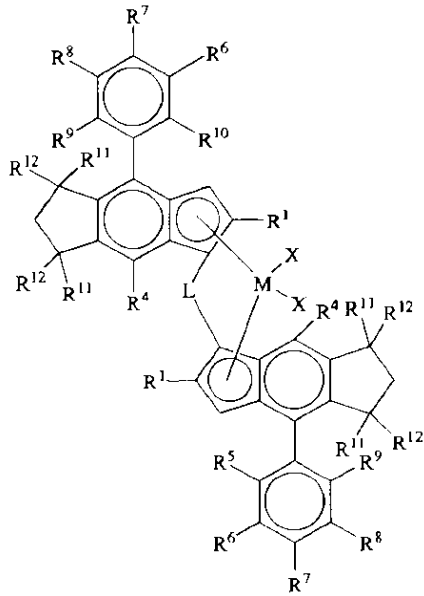
RID	Definition
B	B = CR _{4b} v N
D	D = CR _{4b} v N
E	E = CR ⁹ R ⁹ v CR ⁹ v O v N v NH
G	G = CH v N
J	J = CR _{9R_9} v CR ₉ v C=O v C=S v N

R Logic:

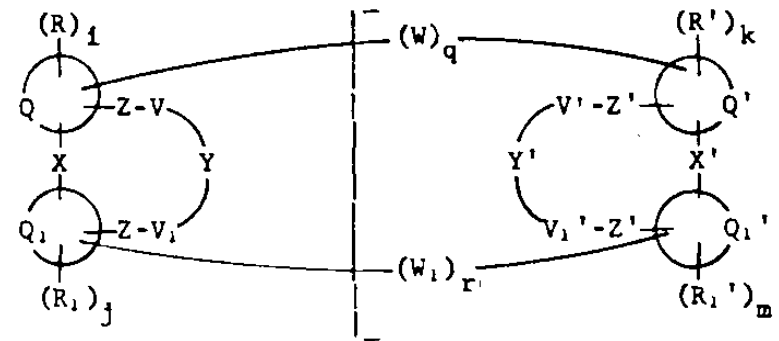
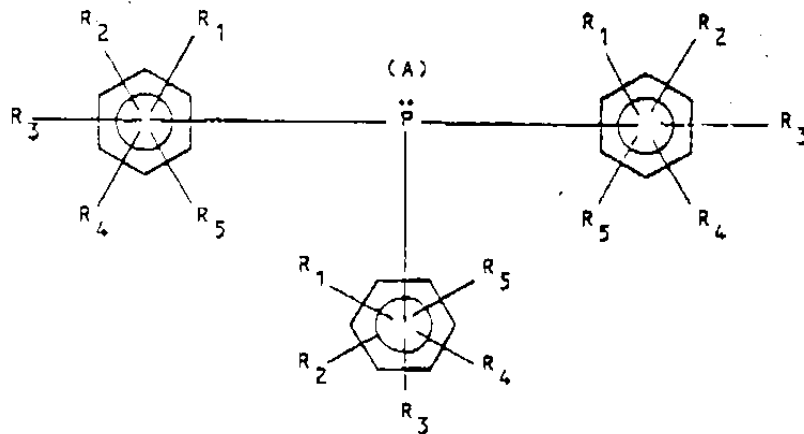
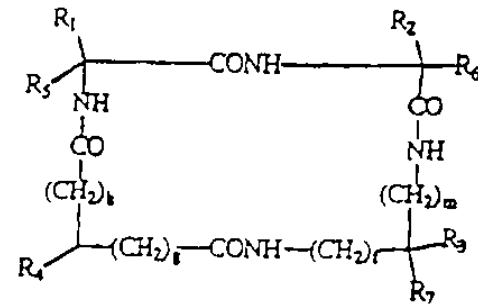
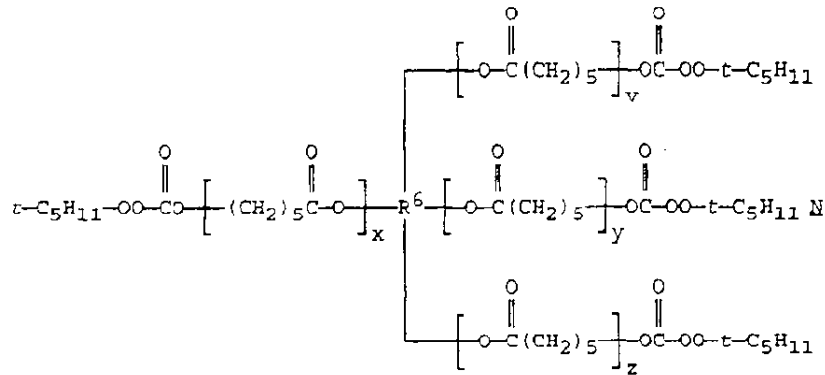
ID	Logic
1	(A v B v D)+ = N
2	IF DOUBLE THEN J NOT C=O v C=S v CR _{9R_9}
3	R_1 = ARYL v HETEROARYL ??? R_4
4	R_2: CYCLOALKYL = CYCLOALKYL v R_6*
5	R_2: ARYL = ARYL v R_6*

At the bottom of the right panel, there are buttons for "Add Definition", "Add Logic", and a context menu with options "Edit", "Assign to R-Group", and "Remove".

Some Examples from Patents (I)



Some Examples from Patents (II)



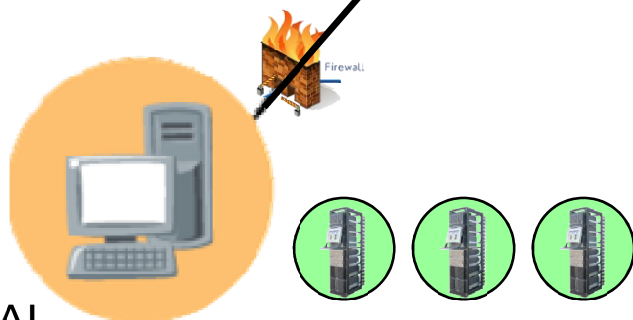
The Grand Patent Challenge: Content Production

- Is it computationally feasible?
- How do we use a protocol DB?
- What is the best job size?
- What's in there?



JUGGLE: Computing Element

UNICORE
D-Grid



SCAI

- Storage Element: Input (200Gb) and Output
- License Server
- ProtocolDB

Job Definition in the Grid

What is the optimal job size?

- Load balancing
- Context
- Dependencies
- Queuing mechanism
- Job management
- Failures
- Overhead

(19) **United States**
(12) **Patent Application Publication** (10) Pub. No.: US 2005/0182053 A1
(43) Pub. Date: Aug. 18, 2005

(54) **SUBSTITUTED 3-AMINO-THIENO[2,3-B]PYRIDINE-2-CARBOXYLIC ACID AMIDE COMPOUNDS AND PROCESSES FOR PREPARING AND THEIR USES**

(51) Int. Cl.⁷ A61K 31/5377; A61K 31/496; A61K 31/4743
(52) U.S. Cl. 514/232.5; 514/301; 514/253.04; 544/125; 544/362; 546/114

(75) Inventors: **Zhidong Chen**, New Milford, CT (US); **Pier Francesco Cirillo**, Woodbury, CT (US); **Darren DiSalvo**, New Milford, CT (US); **Weimin Liu**, Sandy Hook, CT (US); **Daniel Richard Marshall**, Sandy Hook, CT (US); **Lifen Wu**, New Milford, CT (US); **Erick Richard Roush Young**, Danbury, CT (US)

(57) **ABSTRACT**
Disclosed are compounds of formula (I):

(1)

Correspondence Address:
MICHAEL P. MORRIS
BOEHRINGER INGELHEIM CORPORATION
900 RIDGEBURY ROAD
P O BOX 368
RIDGEBELD, CT 06877-0368 (US)

(74) Assignee: **Boehringer Ingelheim Pharmaceuticals, Inc.**, Ridgefield, CT

(21) Appl. No.: **11/002,828**
(22) Filed: **Dec. 2, 2004**
Related U.S. Application Data
(60) Provisional application No. 60/527,522, filed on Dec. 5, 2003.

wherein the variables R₁, R₂, R₃ and Z are described herein, which are useful as inhibitors of the kinase activity of the IκB kinase (IKK) complex. The compounds are therefore useful in the treatment of IKK mediated diseases including autoimmune diseases inflammatory diseases and cancer. Also disclosed are pharmaceutical compositions comprising these compounds and processes for preparing these compounds.

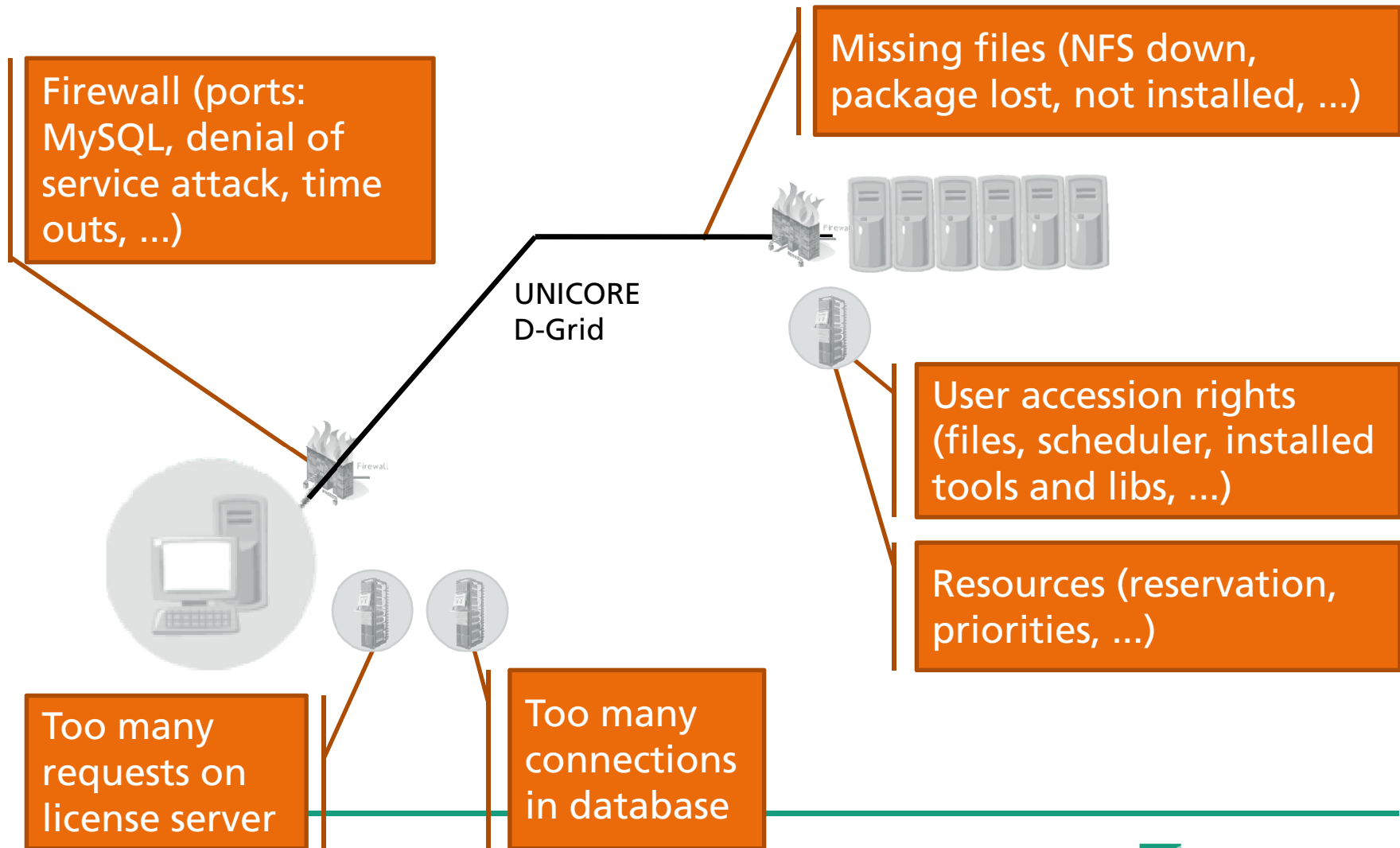
Different Scheduling Strategies

- Let us assume:
 - We process 25,000 documents
 - A document has on average 20 pages
 - A page contains on average 10 molecules

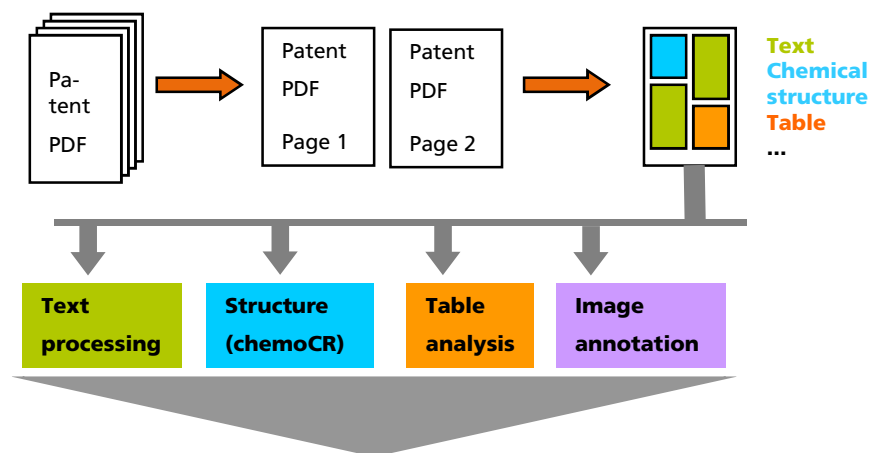
	Per Document	Per Page	Per Chemistry
Number of Jobs	25,000	500,000	5,000,000
Run Time (sec.)	2400	120	12

Where do we achieve the best crunching factor?

Technical Issues and Pitfalls



Outlook: Next Generation?



Web based Retrieval System
Linking Text and Structural Searches

Results from the Image & Text Extraction

Images	Total	AVG.
Documents	23,945	
Run Time (h)	45,360	0.1
Pages	692,576	28.9
Structures	531,505	22.2
Max page	946	
Data (Gb)	100.9	
Full Text	Total Docs	AVG.
Genes/ Proteins	824,415	30.27
IUPAC	128,997	2.39
Drug	170,857	45.66

SQL Shell

```
SQL Statement:
select url, predictedMCD, predictedLabelsSVM_full_radial, predictedLabelsSVM_polyn_f5, predictedLabelsSVM_radial_f5, predictedLabelsKNN_f5, ipidid.value as image,
i.value as molecule from infochem_structure i join (
select i.id, image_id, predictedMCD, predictedLabelsSVM_full_radial, predictedLabelsSVM_polyn_f5, predictedLabelsSVM_radial_f5, predictedLabelsKNN_f5, value, url
from item i join (
select reconstruction_id, image_id, predictedMCD, predictedLabelsSVM_full_radial, predictedLabelsSVM_polyn_f5, predictedLabelsSVM_radial_f5,
predictedLabelsKNN_f5, value, url from item i join (
select item_id, ip.image_id, predictedMCD, predictedLabelsSVM_full_radial, predictedLabelsSVM_polyn_f5, predictedLabelsSVM_radial_f5, predictedLabelsKNN_f5,
value, url from image_data id join (
select image_id, predictedMCD, predictedLabelsSVM_full_radial, predictedLabelsSVM_polyn_f5, predictedLabelsSVM_radial_f5, predictedLabelsKNN_f5, value, url from
prediction score in join)
```

History

Description:

Next SQL Statement Previous SQL Statement

Result:

#	sel...	url	predictedMCD	predictedLabel...	predictedLabel...	predictedLabel...	predictedLabel...	image	molecule
424	<input checked="" type="checkbox"/>	file:/home/thors...	0.87629207380...	-1	-1	-1	bad		ICI FHG SCAI 2008
425	<input checked="" type="checkbox"/>	file:/home/thors...	0.9786556288...	-1	-1	-1	good		fragment(s) of created from WO2 ICI FHG SCAI 2008
426	<input checked="" type="checkbox"/>	file:/home/thors...	0.8760117354...	-1	-1	-1	good		frameH3CCH3created from W ICI FHG SCAI 2008
427	<input checked="" type="checkbox"/>	file:/home/thors...	0.9844572242...	-1	-1	-1	good		fragment(s) of created from WO2005000000-01-0007 ICI FHG SCAI 2008
428	<input checked="" type="checkbox"/>	file:/home/thors...	0.8304874700...	-1	-1	-1	good		fragment(s) of created from W ICI FHG SCAI 2008
429	<input checked="" type="checkbox"/>	file:/home/thors...	0.9020889198...	-1	-1	-1	bad		ICI FHG SCAI 2008

Send Clear Invert Export Close row height: 107

SCAI Chem Patents Focus on Documents

- Collection of 56 000 patents also available as text and image

IPC Class/At least one Entity mention	Number of Patents
A61 Medical Science	34,704
C12 Enzyme Technology	4,851
A61 + C12	35,381
Gene/Protein	43,103
ATC Class	52,552
MeSH Disease	38,474
MedDRA	55,748
Drug Name	44,001
IUPAC	52,720

Example of IUPAC Annotation

tetrahydrodiazepine and diazepane. The compound according to claim 1, wherein ring B is a ring selected from pyrrole, imidazole, pyrazole, thiazole, oxazole, pyridine, pyrimidine, dihydrotriazine, pyrazine and dihydropyrimidine. The compound according to claim 1, wherein the compound of formula (I) is 5,6,7,8-tetrahydropyrimido[4,5-d]pyrimidine-2-carbonitrile, pyrazolo[1,5-a]pyrimidine-5-carbonitrile, 2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, pyrazolo[1,5-a][1,3,5]triazine-2-carbonitrile, 1H-pyrimido[4,5-e][1,3,4]oxadiazine-7-carbonitrile, 1H-pyrazolo[3,4-

cyclic group containing
om or a nitrogen atom; n
same or different, in
a salt thereof, a solvate
a C5-7 monocyclic
ven-membered
cted from an oxygen
m 1, wherein ring A is a
rahydropyrimidine,

dihydropyridazine, pyridazine, dihydropyrimidine, dihydropyrazine, dihydrotriazine, pyrazole, dihydropyrazole, pyrrole, imidazole, triazole, thiophene, furan, dihydrofuran, oxadiazine, tetrahydrodiazepine and diazepane. The compound according to claim 1, wherein ring B is a ring selected from pyrrole, imidazole, pyrazole, thiazole, oxazole, pyridine, pyrimidine, dihydrotriazine, pyrazine and dihydropyrimidine. The compound according to claim 1, wherein the compound of formula (I) is 5,6,7,8-tetrahydropyrimido[4,5-d]pyrimidine-2-carbonitrile, pyrazolo[1,5-a]pyrimidine-5-carbonitrile, 2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, pyrazolo[1,5-a][1,3,5]triazine-2-carbonitrile, 1H-pyrimido[4,5-e][1,3,4]oxadiazine-7-carbonitrile, 1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, imidazo[1,2-a]pyrimidine-2-carbonitrile, 1,3-benzothiazol-2-carbonitrile, 5,6,7,8,9-tetrahydro-5H-pyrimido[4,5-e][1,4]diazepine-2-carbonitrile, 5H-pyrrolo[3,2-d]pyrimidine-2-carbonitrile, pyrido[2,3-d]pyrimidine-2-carbonitrile, 5,7-dihydrofuro[3,4-d]pyrimidine-2-carbonitrile, 6,7-dihydro-5H-cyclopenta[d]pyrimidine-2-carbonitrile, 9H-purine-2-carbonitrile, or 1,3-benzoxazole-2-carbonitrile. The compound according to claim 1, which is a compound of formula (I-A) or formula (I-B) wherein R is a hydrogen atom or a substituent; n1 is 0, or an integer of from 1 to 3; and R is wherein T, T, T and T each is independently a bond or a spacer having from 1 to 10 atoms of the principle chain; ring 1, ring1 and ring2 each is

group which may have a substituent(s). 1C 3 2 1 2C C C The compound according to claim 1, which is selected from the group consisting of: (1) 8-(2,2-dimethylpropyl)-7-oxo-5,6,7,8-tetrahydropyrimido[4,5-d]pyrimidine-2-carbonitrile, (2) 1-(2,2-dimethylpropyl)-2-(4-methoxybenzyl)-3-oxo-2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (3) 1-(2,2-dimethylpropyl)-3-[(4-methoxybenzyl)oxy]-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (4) 4-(2,2-dimethylpropoxy)-1,3-benzothiazol-2-carbonitrile, (5) 3-(4-biphenylmethoxy)-1-(2,2-dimethylpropyl)-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (6) 2-(4-biphenylmethyl)-1-(2,2-

substituent, and other symbols have
to claim 6, wherein R is a
E each is independently a bond
(=O)-, -C(=O)NR-, -NRC(=O)-,
a hydrogen atom or a
containing a nitrogen atom which
herein -T- is -E-E-E- in which all
clic carbocyclic group which may
ich may have a substituent(s). 1B
same meanings as described in
ring 1 and/or ring 2 is (1) C3-10
ed monocyclic or bicyclic

heterocyclic group which may have a substituent(s). 1C 3 2 1 2C C C The compound according to claim 1, which is selected from the group consisting of: (1) 8-(2,2-dimethylpropyl)-7-oxo-5,6,7,8-tetrahydropyrimido[4,5-d]pyrimidine-2-carbonitrile, (2) 1-(2,2-dimethylpropyl)-2-(4-methoxybenzyl)-3-oxo-2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (3) 1-(2,2-dimethylpropyl)-3-[(4-methoxybenzyl)oxy]-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (4) 4-(2,2-dimethylpropoxy)-1,3-benzothiazol-2-carbonitrile, (5) 3-(4-biphenylmethoxy)-1-(2,2-dimethylpropyl)-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (6) 2-(4-biphenylmethyl)-1-(2,2-dimethylpropyl)-3-oxo-2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (7) 1-(2,2-dimethylpropyl)-3-oxo-2-(2-thienylmethyl)-2,3-dihydro-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (8) 1-(2,2-dimethylpropyl)-3-[2-(4-morpholinyl)ethoxy]-1H-pyrazolo[3,4-d]pyrimidine-6-carbonitrile, (9) 9-(2,2-dimethylpropyl)-6-(4-methoxybenzyl)-7-oxo-6,7,8,9-tetrahydro-5H-pyrimido[4,5-e][1,4]diazepine-2-carbonitrile, (10) 8-(2,2-dimethylpropyl)-6-(4-methoxybenzyl)-7-oxo-5,6,7,8-tetrahydropyrimido[4,5-d]pyrimidine-2-carbonitrile, (11) 4-[(2,2-dimethylpropyl)amino]pyrido[2,3-d]pyrimidine-2-carbonitrile, (12) 1-(2,2-dimethylpropyl)-3-[4-[(4-methyl-1-piperazinyl)methyl]phenyl]-1H-pyrimido[4,5-e][1,3,4]oxadiazine-7-carbonitrile, (13) 2-[6-cyano-1-(2,2-dimethylpropyl)-3-oxo-1,3-dihydro-2H-pyrazolo[3,4-d]pyrimidin-2-yl]-N-[2-(dimethylamino)ethyl]acetamide, (14) 4-[(2,2-dimethylpropyl)amino]-5,7-dihydrofuro[3,4-d]pyrimidine-2-carbonitrile, and (15) 4-[(2,2-dimethylpropyl)amino]-6,7-dihydro-5H-cyclopenta[d]pyrimidine-

But: Huge chemical Name Regions in the same Patent

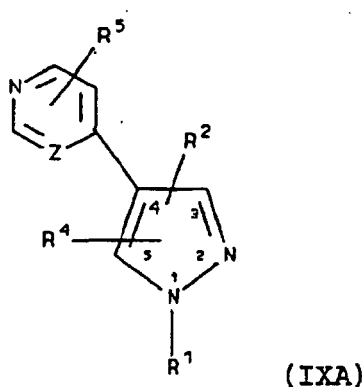
and which may be partially saturated or fully saturated, . Three- to fifteen-membered mono-, polycyclic aromatic heterocyclic group, and polycyclic heterocyclic group having spiro bond, and polycyclic bridged heterocyclic group which contain 1 to 5 hetero atoms selected from an oxygen atom(s), a nitrogen atom(s) and a sulfur atom(s) which may be oxidized, and which may be partially saturated or fully saturated include, for example, pyrrole, imidazole, triazole, tetrazole, pyrazole, pyridine, pyrazine, pyrimidine, pyridazine, triazine, azepine, diazepine, furan, pyran, oxepine, thiophene, thiopyran, thiepine, oxazole, isoxazole, thiazole, isothiazole, furazan, oxadiazole, oxazine, oxadiazine, oxazepine, oxadiazepine, thiadiazole, thiazine, thiadiazine, thiazepine, thiadiazepine, indole, isoindole, indolizine, benzofuran, isobenzofuran, benzothiophene, isobenzothiophene, dithianaphthalene, indazole, quinoline, isoquinoline, quinolizine, purine, phthalazine, pteridine, naphthyridine, quinoxaline, quinazoline, cinnoline,

primary saturated or fully p, and polycyclic h contain 1 to 5 hetero ch may be oxidized, and zole, triazole, tetrazole, pyran, oxepine, azole, oxazine, diazepine, indole, dithianaphthalene, ne, quinoxaline, ne, benzimidazole, azepine, benzotriazole, carbazole, iazine, phenoxazine, dine, pyrazoloisoquinoline, zoline, imidazolidine, popyrazole (pyrazolidine), perazine, ydropyridazine, ne, tetrahydroazepine, drodiazepine), oxirane,

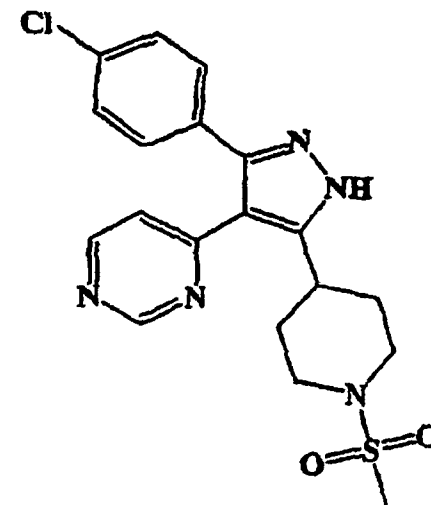
oxetane, dihydrofuran, tetrahydrofuran, dihydropyran, tetrahydropyran, dihydrooxepine, tetrahydrooxepine, perhydrooxepine, thirane, thietane, dihydrothiophene, tetrahydrothiophene, dihydrothiopyran, tetrahydrothiopyran, dihydrothiepine, tetrahydrothiepine, perhydrothiepine, dihydrothiazole, tetrahydrothiazole (oxazolidine), dihydroisoxazole, tetrahydroisoxazole (isoxazolidine), dihydrothiazole, tetrahydrothiazole (thiazolidine), dihydroisothiazole, tetrahydroisothiazole (isothiazolidine), dihydrofurazan, tetrahydrofurazan, dihydrooxadiazole, tetrahydrooxadiazole (oxadiazolidine), dihydrooxazine, tetrahydrooxazine, dihydrooxadiazine, tetrahydrooxadiazine, dihydrooxazepine, tetrahydrooxazepine, perhydrooxazepine, dihydrooxadiazepine, tetrahydrooxadiazepine, perhydrooxadiazepine, dihydrothiadiazole, tetrahydrothiadiazole (thiadiazolidine), dihydrothiazine, tetrahydrothiazine, dihydrothiadiazine, tetrahydrothiadiazine, dihydrothiazepine, tetrahydrothiazepine, perhydrothiazepine, dihydrothiadiazepine, tetrahydrothiadiazepine, perhydrothiadiazepine, morpholine, thiomorpholine, oxathiane, indoline, isoindoline, dihydrobenzofuran, perhydrobenzofuran, dihydroisobenzofuran, perhydroisobenzofuran, dihydrobenzothiophene, perhydrobenzothiophene, dihydroisobenzothiophene, perhydroisobenzothiophene, dihydroindazole, perhydroindazole, dihydroquinoline, tetrahydroquinoline, octahydroquinoline, perhydroquinoline, dihydroisoquinoline, tetrahydroisoquinoline, octahydroisoquinoline, perhydroisoquinoline, dihydrophthalazine, tetrahydrophthalazine, perhydrophthalazine, dihydronaphthyridine, tetrahydronaphthyridine, perhydronaphthyridine, dihydroquinoxaline, tetrahydroquinoxaline, perhydroquinoxaline, dihydroquinazoline, tetrahydroquinazoline, perhydroquinazoline, tetrahydropyrrolopyridine, dihydrocinnoline, tetrahydrocinnoline, perhydrocinnoline, benzoxathiane, dihydrobenzoxazine, dihydrobenzothiazine, pyrazinomorpholine, dihydrobenzoxazole, perhydrobenzoxazole, dihydrobenzothiazole, perhydrobenzothiazole, 4,5,6,7-tetrahydrothieno[3.2-c]pyridine, dihydrobenzimidazole, perhydrobenzimidazole, dihydrobenzazepine, tetrahydrobenzazepine, dihydrobenzodiazepine, tetrahydrobenzodiazepine, benzodioxepane, dihydrobenzoxazepine, tetrahydrobenzoxazepine, dihydrocarbazole, tetrahydrocarbazole, perhydrocarbazole, dihydroacridine, tetrahydroacridine, perhydroacridine, dihydrodibenzofuran, dihydrodibenzothiophene, tetrahydrodibenzofuran, tetrahydrodibenzothiophene, perhydrodibenzofuran, perhydrodibenzothiophene, tetrahydrodibenzothiophene, tetrahydro-β-carboline, dihydroazepinoindole,

Category: largest Document (I)

- 946 pages
- EP 1144403 B1
- SUBSTITUTED PYRAZOLES AS P38 KINASE INHIBITORS
- G.D. Searle LLC

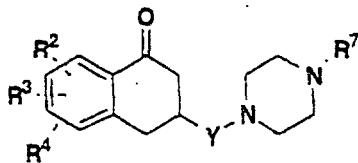



<p>European Patent Office Office européen des brevets</p>		<p>(11) EP 1 144 403 B1</p>
<p>EUROPEAN PATENT SPECIFICATION</p>		
(43) Date of publication as mentioned in the gazette of the patent: 06.10.2004 Bulletin 2004/41	(51) Int. Cl. ⁷ C07D 401/04, A61K 31/44-15, A61K 31/447, A61K 31/445, A61K 31/46, A61K 31/50, A61K 31/505, A61K 31/52, C07D 405/14, C07D 401/14, C07D 403/14, C07D 403/04	
(41) Application no. 99061746.2	(86) International application number: PCT/US01/9503507	
(21) Date of filing: 17.11.1999	(87) International publication number: WO 2002/031962 (02.06.2000 Gazette 2002/23)	
<p>(54) SUBSTITUTED PYRAZOLES AS P38 KINASE INHIBITORS SUBSTITUIERTE PYRAZOLE ALS P38 KINASE-INHIBITOREN PYRAZOLES SUBSTITUEES UTILES COMME INHIBITEURS DE KINASE P38</p>		
<p>(84) Designated Contracting State: AT BE CH CY DE DK ES FR GB GR IE IT LI LU MC NL PT SE</p>		
(30) Priority: 20.11.1998 US 984993	<ul style="list-style-type: none"> Herman, Gunar Skolnik, I. 0977 (US) Herrmann, Susan J. Kilwardt, MO 9332 (US) Haggerty, Michele St. Charles, MO 63039 (US) Huang, He Chicago, IL 60618 (US) Khanlou, Sh. K. Libertyville, IL 60068 (US) Kozlowski, Francis J. Prospect Heights, IL 60070 (US) Liao, Shuyuan Northbrook, IL 60062 (US) Mao, Suzanne Crestedfield, MO 63017 (US) Parra, Richard A. Excelsior, IL 60524 (US) Perry, Thomas D. Red Bank, IL 62258 (US) Rio, Shaileshwar N. St. Louis, MO 63141 (US) Sarkiss, Sharm P. Crestedfield, MO 63017 (US) South, Michael S. Crestedfield, MO 63017 (US) Shelley, Michael A. Libertyville, IL 60068 (US) Talley, John Jeffrey St. Louis, MO 63199 (US) Vargiac, Michael L. Bethesda, MD 20821 (US) 	
(43) Date of publication of application: 17.10.2001 Bulletin 2001/42	(73) Inventor: G.D. Searle LLC St. Louis, Missouri 63141 (US)	
(72) Inventors: Arendtsmeeren, Ashok Baltimore, IL 60620 (US) Clegg, Michael Stockton, IL 60217 (US) Collins, Paul Grapeville, IL 60626 (US) Cruz, Joseph Z. Germantown, IL 60426 (US) Dennis, Robert Bethesda, MD 20821 (US) Flynn, Daniel Thousand Oaks, CA 91320 (US) Gale, Jeffrey Stockton, IL 60217 (US) Granato, Matthew J. Crestedfield, MO 63017 (US) Hosar, Catherine C. Crestedfield, MO 63017 (US)		
<p>Note: This document is intended for publication in the Journal of the European Patent Office. Any person who is not a party to the European Patent Office procedure for the grant of a European patent should not make any statement in writing or by any other means. It shall not be deemed to have been filed with the competent authority in any country.</p>		
EP 1 144 403 B1	Printed on 10/16/04	(CONT. next page)



Category: largest Document (II)

- 715 pages
- EP 0904262 B1
- CONFORMATIONALLY RESTRICTED AROMATIC INHIBITORS OF MIRCOSOMAL TRIGLYCERIDE TRANSFER PROTEIN AND METHOD
- BRYSTOL-MYERS SQUIBB




 Europäische Patentamt
 European Patent Office
 Office européen des brevets

(11) EP 0 904 262 B1

(12) EUROPEAN PATENT SPECIFICATION

(43) Date of publication and mention of the grant of the patent: 21.04.2004 Bulletin 2004/17

(51) Int. Cl.⁷ C07C 217/04, C07D 471/10, C07D 233/79, C07D 401/08, C07D 403/08, A61K 31/24, A61K 31/645, A61K 31/415, A61K 31/44, A61K 31/42

(21) Application number: 9730385.9

(22) Date of filing: 13.01.1997

(86) International application number: PCT/US97/00587

(87) International publication number: WO 1997/02024 (24.01.1997 Gazette 1997/03)

(54) CONFORMATIONALLY RESTRICTED AROMATIC INHIBITORS OF MICROSOMAL TRIGLYCERIDE TRANSFER PROTEIN AND METHOD.
 KONFORMATIONSEINGESCHRÄNKTE AROMATISCHE INHIBITOREN DES MICROSOMAL FUNKTIONSTRANSFERPROTEINS (A7) VERFAHREN.
 COMPOSES AROMATIQUES A CONFORMATION RESTREINTE SERVANT D'INHIBITEURS DE LA PROTEINE DE TRANSFERT DE TRIGLYCERIDE VICHOSOMIQUE ET PROCEDURE ASSOCIEE

(84) Designated Contracting States: AT BE CH DE DK ES FR GB GR IE IT LI LU MC NL PT SE

(30) Priority: 16.01.1996 US 93366 P
 06.06.1996 US 17224 P
 05.11.1996 US 59370 P

(45) Date of publication of application: 31.03.1998 Bulletin 1998/12

(73) Proprietor: BRYSTOL-MYERS SQUIBB COMPANY
 Princeton, NJ 08543-4000 (US)

(72) Inventor:
 - BILLEN, Scott A.
 Hopewell, NJ 08525 (US)
 - DREEDER, John H.
 Eastampton, NJ 08060 (US)
 - LAWRENCE, P. Michael
 Tuxedo, PA 15070 (US)
 - MAGNOLI, David R.
 Hamilton, NJ 05080 (US)
 - POSS, Michael A.
 Lawrenceville, NJ 08648 (US)

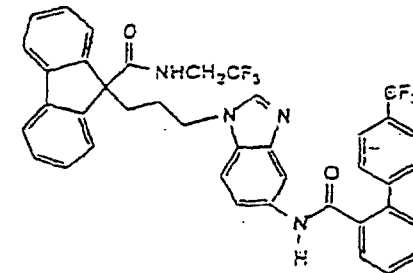
(74) Representatives: VOSBOM & PARTNER
 Shoups Plaza 4
 8105 Mountain View
 - ROELL, Jeffrey A.
 Newcomer, PA 19063 (US)
 - BUSCHGANG, William A.
 Summiton, NJ 08558 (US)
 - BUSCH, Richard G.
 Franklin Park, NJ 08823 (US)
 - TRUC, James A.
 Lawrenceville, NJ 08648 (US)

(56) References cited:
 EP-A- 842 857 WO-A-96/0205
 WO-A-86/0446 25.03.86 131
 GB-A- 773 758 GB-A- 1 147 832
 US-A- 5 758 596 US-A- 5 848 197
 US-A- 4 197 212 US-A- 4 277 465
 US-A- 4 282 716 US-A- 4 864 328
 US-A- 5 172 498 US-A- 5 232 208

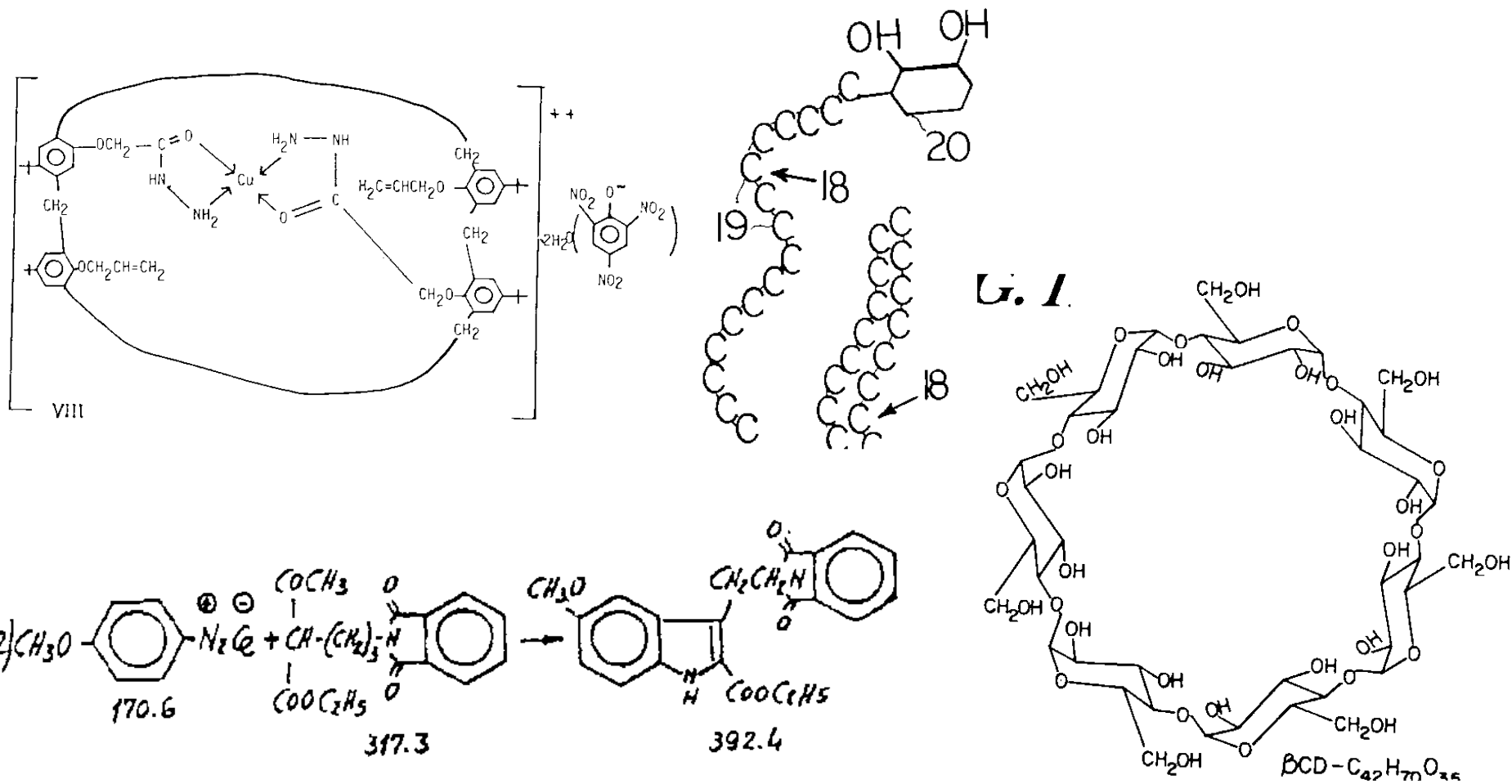
Remarks:
 The file contains technical information submitted after the publication date and not included in this specification.

Note: After nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition has been dealt with (Art. 93(1) European Patent Convention).

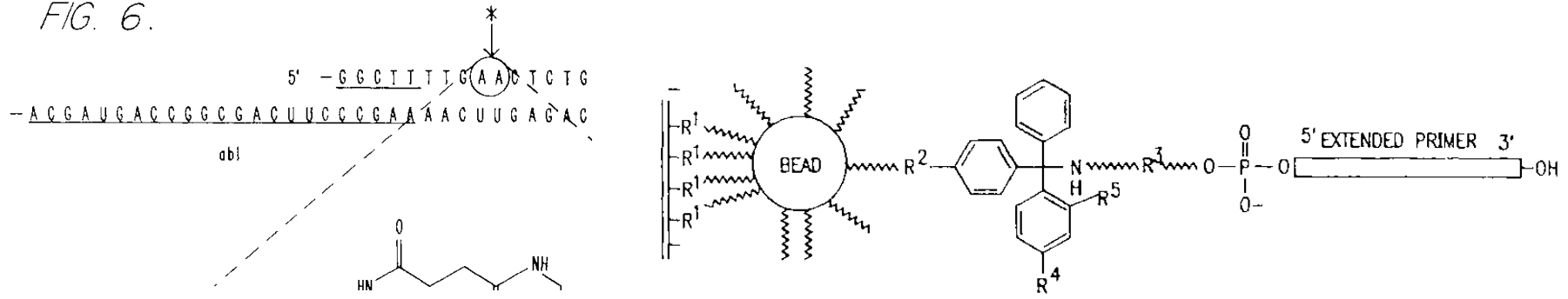
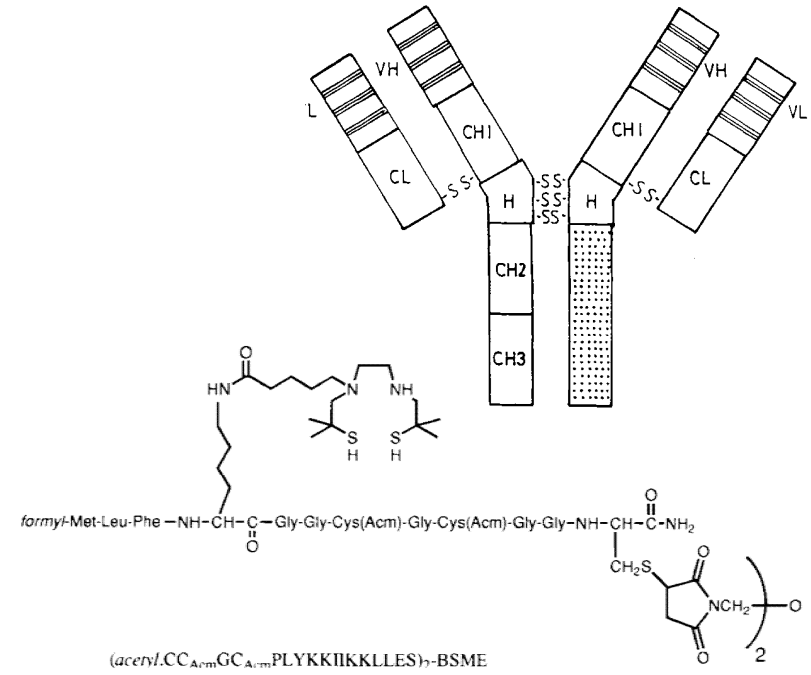
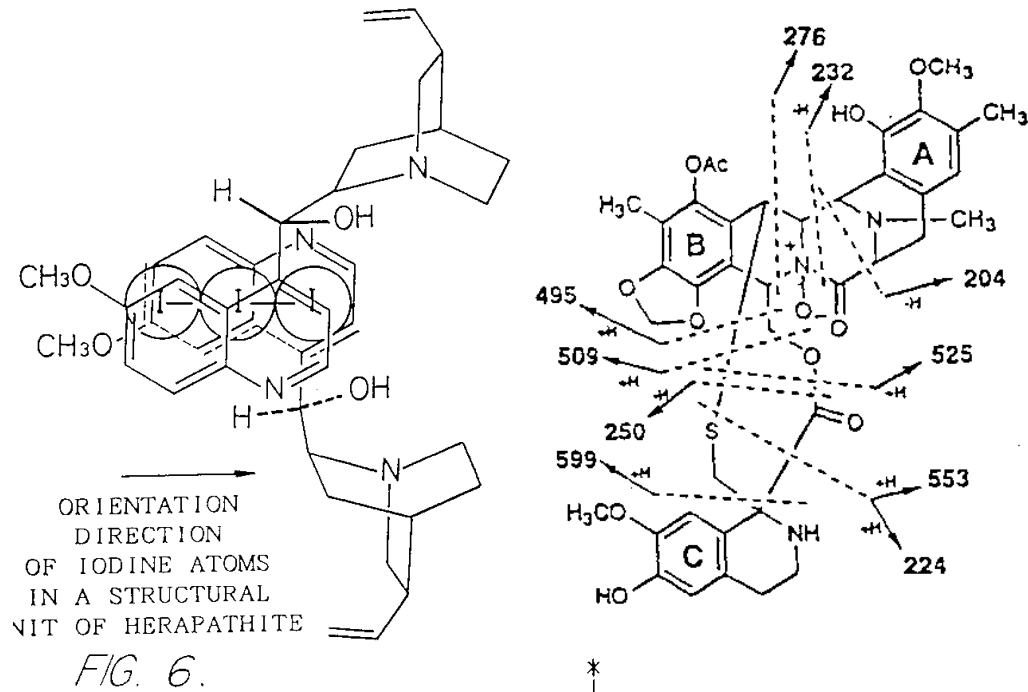
EP 0 904 262 B1



Category: That's Organic Chemistry

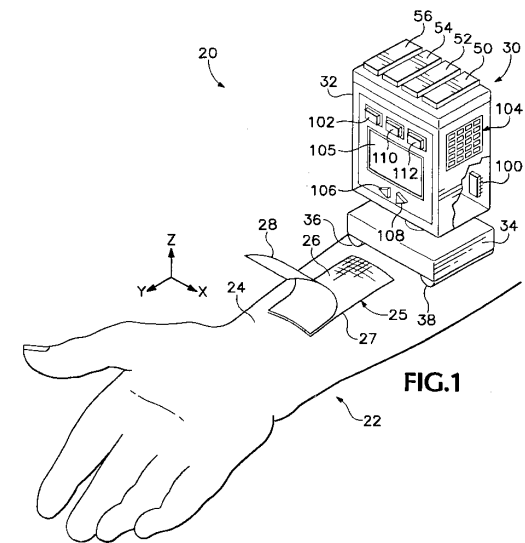


Category: Sort of Chemistry



Conclusions

- Image reconstruction is an ongoing effort
- Patents offer a lot of challenges
- In need of critical assessments and gold standards
 - Image reconstruction
 - Retrieval tasks
 - Hardware infrastructures
- In need of strategies
 - Deal with reconstruction errors
 - Extended file formats & search algorithms
 - Result visualizations



chemOCR

chemical compound
reconstruction



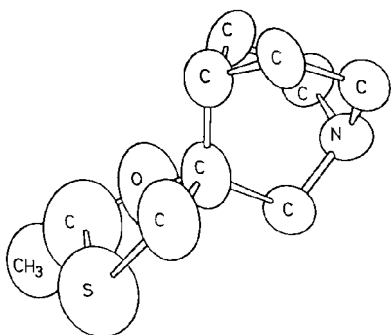
<http://www.scai.fraunhofer.de/chemocr>



Hans Kraut
Dr. Josef Eiblmaier
Dr. Peter Löw



Dr. Jörn Kohlhammer



Fachhochschule
Bonn-Rhein-Sieg



Tanja Fey, Carina Haupt
Sebastian Ginzel
Prof. Dr. Ralf Thiele

Le Thuy Bui Thi
Prof. Dr. Noltemeier



Ludwig-
Maximilians-
Universität
München

LMU

Peter Kral
Karsten Borgwardt
Prof. Dr. Hans-Peter Kriegel

ITAM

INSTITUTO TECNOLÓGICO AUTÓNOMO DE MÉXICO

Miguel Alvarez
Santiago Akle Serrano
Prof. Dr. Maria Elena Algorri

b-it

Bonn-Aachen International Center for
Information Technology

Yuan Wang, Olga Domanova
Wei Wang, Albert Ou
Prof. Dr. M. Hofmann-Apitius



RheinAhrCampus - Remagen

Nils Frings, Michael Lenz
Prof. Dr. Kschischo



Fraunhofer
SCAI

Dr. Marc Zimmermann
Dr. Christoph Friedrich
Angelika Weihermüller
Thorsten Prinz