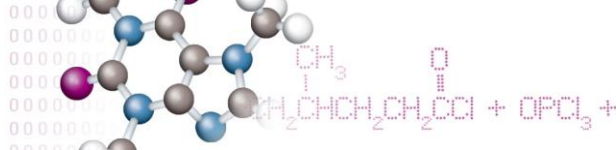


-7.5909	-1.3883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



# ChemProspector: Advanced Mining and Searching of Chemical Content in Patent Documents

Josef Eiblmaier (InfoChem), Hans Kraut (InfoChem), Larisa Insenko (InfoChem), Heinz Saller (InfoChem), Peter Loew (InfoChem)

*ICIC 24 – 27 October, Vienna*

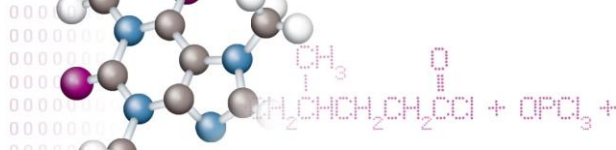
Supported by:



on the basis of a decision  
by the German Bundestag



-7.5909	-1.3883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



# Outline



## » Introduction

- › ChemProspector, a THESEUS project
- › Markush in a nutshell

## » Major Goals and Approach

## » First Results

## » Outlook



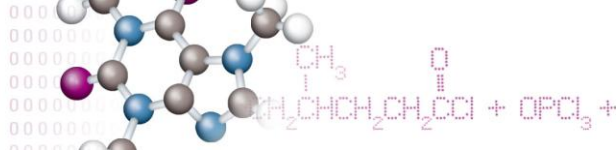
© cora / PIXELIO, www.pixelio.de

Supported by:



on the basis of a decision  
by the German Bundestag

-7.5909	-1.1833	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



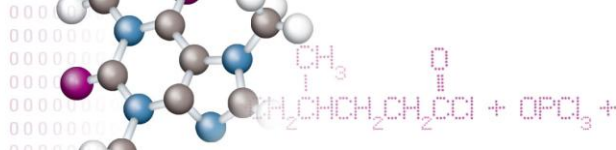
- » “New Technologies for the Internet of Services”
- » Research program initiated by the Federal Ministry of Economy and Technology (BMWi)
- » Supported with approx. 100 million Euros
- » Facilitate access to information, combine data to form new kinds of knowledge and lay the groundwork for new services on the Internet
- » Duration: five years (2007 - 2011)
- » Divided into six application scenarios combined through core technology cluster
- » Phase one: development of core technologies (2007 - 2008)
- » Phase two: THESEUS SME (2009 - 2011)

Supported by:



on the basis of a decision  
by the German Bundestag

-7.5909	-1.1833	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



# ChemProspector: Basic data

» Main emphasis:

*‘The automatic extraction of Markush Structures  
from patent documents’*

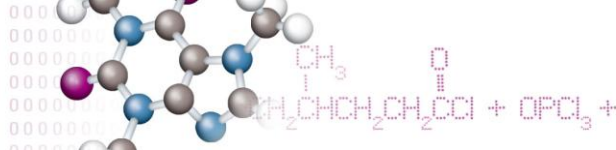
- » Research SME-project within the THESEUS research program
- » Application scenario ORDO (“Ordnung Digitaler Information”)
- » Duration: July 2009 – end of 2011





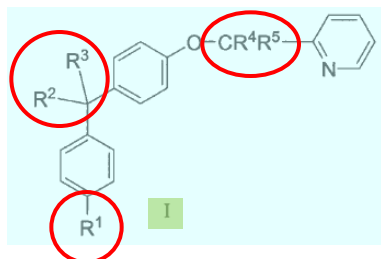


-7.5909 -1.3883 -0.4896 C  
 -9.9889 -0.1947 -0.3218 O  
 -10.2795 2.5700 0.1393 C  
 -8.5096 -1.8624 -0.6298 H  
 -7.5171 -2.3313 0.2821 H  
 -7.0177 -1.8762 -1.3207 H



DETAILED DESCRIPTION OF THE INVENTION

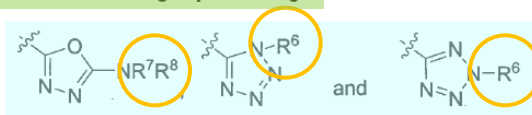
The instant invention provides a compound represented by structural formula I



10

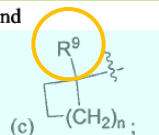
the pyridyl-N-oxide analog of formula I, and the pharmaceutically acceptable salts, esters and solvates thereof wherein:

R<sup>1</sup> is selected from the group consisting of:



15

R<sup>2</sup> is selected from the group consisting of (a) -C<sub>1-6</sub>alkyl optionally substituted with 1-3 of fluoro, (b) -C<sub>3-6</sub> cycloalkyl and



n is an integer selected from 0, 1, 2 and 3;

R<sup>3</sup> is selected from the group consisting of -H, -F, -OH, -CH<sub>3</sub> and -CF<sub>3</sub>;

20

R<sup>4</sup> is selected from the group consisting of -H and -C<sub>1-4</sub>alkyl;

R<sup>5</sup> is selected from the group consisting of -H and -CH<sub>3</sub>; and

R<sup>6</sup> is selected from the group consisting of -H, -C<sub>1-6</sub>alkyl optionally substituted with 1-3 fluoro, -C<sub>3-6</sub> cycloalkyl optionally substituted with 1-3 fluoro and -CH<sub>2</sub>-R<sup>10</sup>;

25

R<sup>7</sup> is selected from the group consisting of -H, -C<sub>1-6</sub>alkyl optionally substituted with 1-3 fluoro, -C<sub>3-6</sub> cycloalkyl optionally substituted with 1-3 fluoro, -COC<sub>1-6</sub>alkyl and -COC<sub>3-6</sub> cycloalkyl;

Supported by:



on the basis of a decision by the German Bundestag



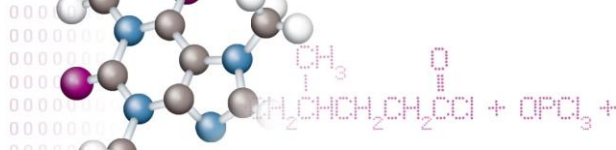




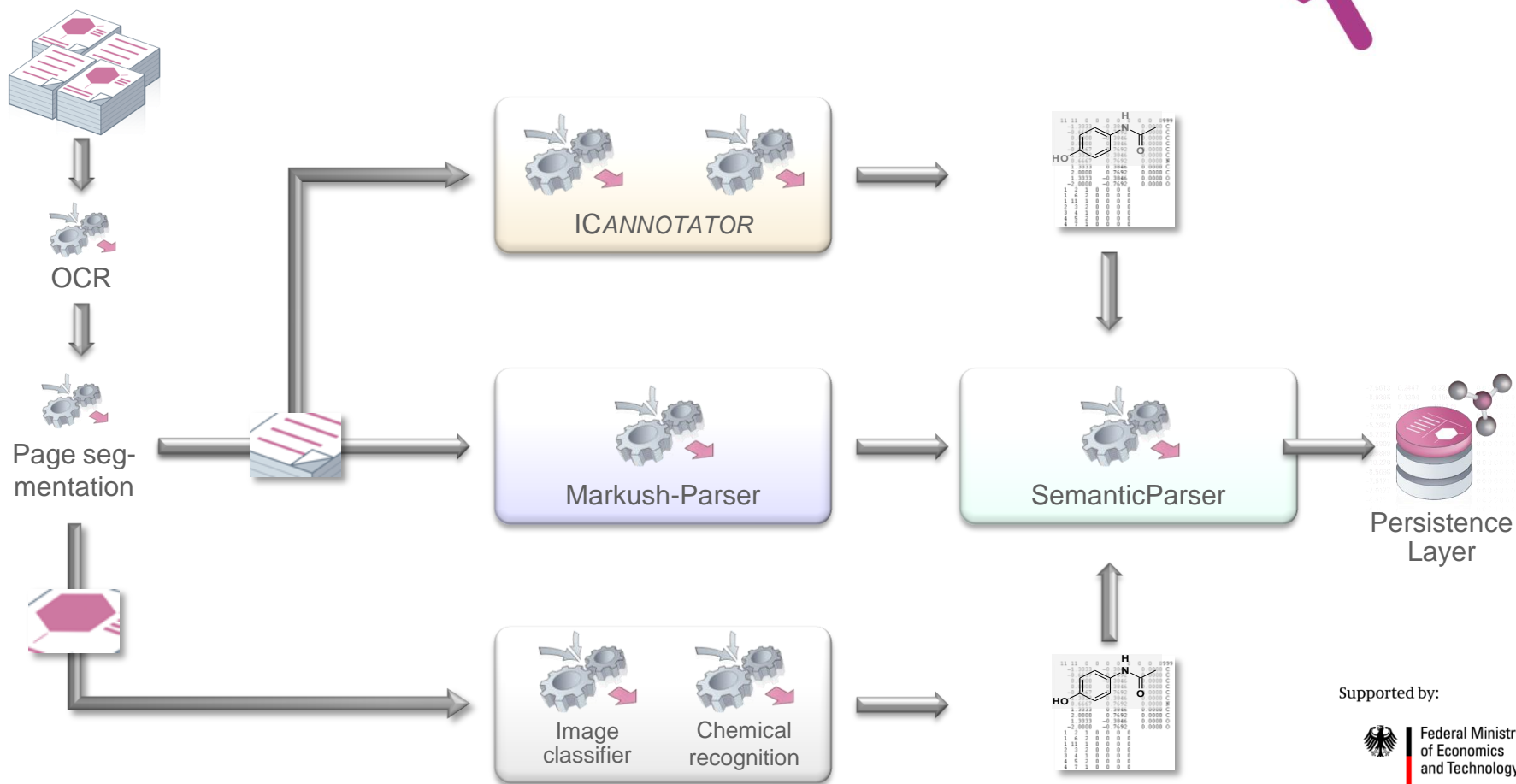




-7.5909	-1.3883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H
-1.4704	1.2424	0.1059	C



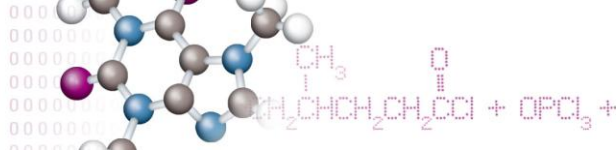
# ChemProspector: Approach



Supported by:



on the basis of a decision  
by the German Bundestag



# ICANNOTATOR



» Extracts chemical named entities

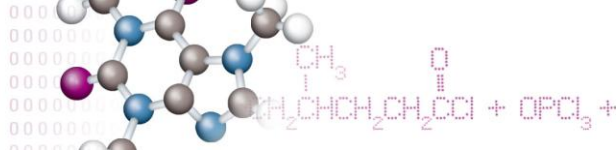
The aromatic hydrocarbon residue may preferably be an aryl group having a carbon number of 6 to 15. Examples thereof include phenyl, cresyl, xylyl, 2,6-dimethylphenyl, 2,4,6-trimethylphenyl, butylphenyl, nonylphenyl and the like.

<b>Exact chemical entities</b>	<i>methyl, ethyl, n-propyl, phenyl, chloro, nitro, amino, hydroxy, hydrogen, carbon, 1-naphthyl, 2-pyridyl, tosyl, piperidyl ...</i>
<b>Generic and homology groups, fragments</b>	<i>alkyl, alkoxy, aryl, halogenid, hydrocarbon ...</i>
<b>Combinations</b>	<i>alkylamino, 4-aryl-phenyl, ...</i>

supported by:



on the basis of a decision  
by the German Bundestag



# Markush-Parser



» Extracts Markush specific entities

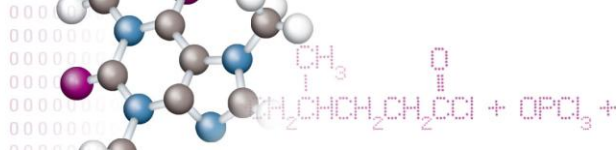
<b>Formula definitions</b>	<i>formula 1, general formula (I), derivatives represented by (3), ...</i>
<b>Variable definitions</b>	<i>R, R<sup>1</sup>, R<sup>2</sup>, R', A, X, Y, Z, Ar, ...</i>
<b>Wherein definitions</b>	<i>where, wherein, in which, ...</i>
<b>Link group</b>	<i>represents, may be, one of, is selected from, ...</i>
<b>Chain lengths</b>	<i>3-20 carbon atoms, ...</i>
<b>Topologic definitions</b>	<i>branched or unbranched, ...</i>
<b>Bond types</b>	<i>may contain double bonds, ...</i>
<b>References</b>	<i>as defined above, ...</i>
<b>Substitutions</b>	<i>optionally substituted by, ...</i>

supported by:



on the basis of a decision  
by the German Bundestag





# One simple sample



US 20030180237A1

(19) **United States**

(12) **Patent**  
Lersch et al.

(54) **COMPOSITION AND TONIC DISORDER AND/OR CL ACTIVE SUBSTANCE**

(75) Inventors: P  
C  
U

Correspondent  
SCULLY SC  
400 GARDE  
GARDEN C

(73) Assignee: G

(21) Appl. No.:

(22) Filed:

(30) Foreign  
Feb. 15, 2002

Pub

(51) Int. Cl.<sup>7</sup> .....

at least one of the compounds of the general formula (I) and/or salts thereof

(I)

in which

R<sup>1</sup> is H, alkyl, hydroxyalkyl, or a carboxyalkyl radical having 2 to 30 carbon atoms; and

Supported by:

on the basis of a decision  
by the German Bundestag





# First Results

## » Sample test set of four USPTO-Patents

Patent	Markush Structures	Markush Structures extracted	Comment
US20070123581	8	6	One structure level 2 One structure has no ligands
US20030180237	4	4	
US20040171886	16	12	Four structures level 2
US20040030160	7	4	Three structures level 2
<b>Sum</b>	<b>35</b>	<b>26</b>	

Supported by:



on the basis of a decision  
by the German Bundestag





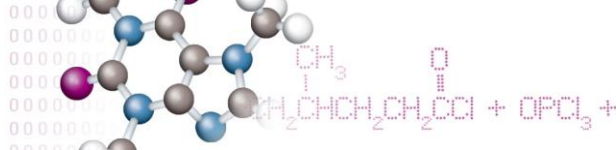








-7.5909	-1.1333	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



## Next steps/Outlook



- » Markush storage and retrieval (Extension of ICCARTRIDGE)
- » Extension of grammar rules to level 2 Markush Structures
- » Page segmentation, image recognition

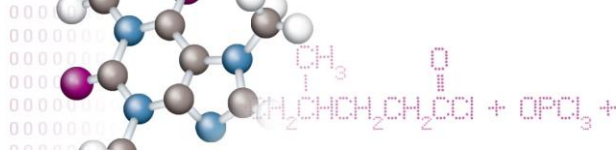


Supported by:



on the basis of a decision  
by the German Bundestag

-7.5909	-1.3883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



# Acknowledgements



» The InfoChem ChemProspector team

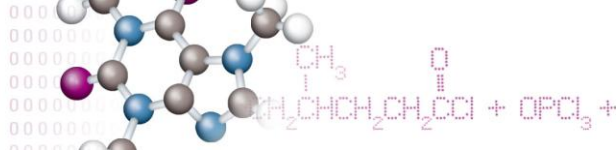


» The German Federal Ministry of Economy and Technology (BMWt)





-7.5909	-1.3883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H



© P. Storz / PIXELIO, www.pixelio.de

# Thank you!

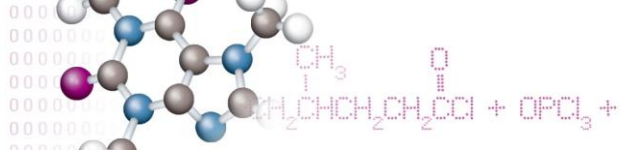
Supported by:



on the basis of a decision  
by the German Bundestag



-7.5909	-1.8883	-0.4896	C
-9.9889	-0.1947	-0.3218	O
-10.2795	2.5700	0.1393	C
-8.5096	-1.8624	-0.6298	H
-7.5171	-2.3313	0.2821	H
-7.0177	-1.8762	-1.3207	H
-1.4791	-1.8121	0.1059	C



© ediathome / PIXELIO, www.pixelio.de

# Questions?

Supported by:



on the basis of a decision by the German Bundestag