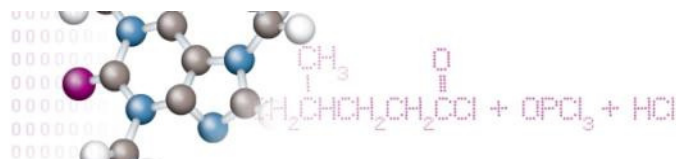




|          |         |           |
|----------|---------|-----------|
| -7.3209  | -1.0003 | -0.4030 C |
| -9.9889  | -0.1947 | -0.3218 O |
| -10.2795 | 2.5700  | 0.1393 C  |
| -8.5096  | -1.8624 | -0.6298 H |
| -7.5171  | -2.3313 | 0.2821 H  |
| -7.0177  | -1.8762 | -1.3207 H |
| -4.4781  | 1.2434  | 0.1858 C  |



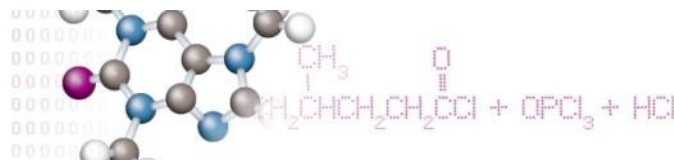
# Progresses in Automated Chemical Structure Recognition in Text and Images

V. Eigner-Pitto and H. Saller.

InfoChem GmbH, Landsberger Strasse 408, Munich, 81241, Germany



|          |         |         |
|----------|---------|---------|
| -7.2209  | -1.0003 | -0.4030 |
| -9.9889  | -0.1947 | -0.3218 |
| -10.2795 | 2.5700  | 0.1393  |
| -8.5096  | -1.8624 | -0.6298 |
| -7.5171  | -2.3313 | 0.2821  |
| -7.0177  | -1.8762 | -1.3207 |
| -4.4781  | 1.2434  | 0.1858  |



# Historical Importance of Chemisches Zentralblatt



1830 Chemisches Zentralblatt

1969

First and oldest abstracts journal in chemistry

Covers chemical literature from 1830 to 1969

Describes the „birth“ of chemistry as science (vs. alchemy)

1840



1907 Chemical Abstracts

Biggest and single abstracts source in chemistry

Currently >31 million papers and patents

Content 1840-1906 added retrospectively

1772 1817 Gmelin Handbook

1771 1881 Beilstein Handbook





|          |         |         |   |
|----------|---------|---------|---|
| -7.2209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |



# Chemisches Zentralblatt: History (1)

1st year: 544 pages/400 abstracts  
After 6 months weekly

„Chemisches Zentralblatt“

1st edition of the General Index

The editorial office moves from Leipzig to Berlin



Introduction of Patent abstracts

**Patente:** C. W. SIEMENS, Verfahren zur Herstellung von Eisen und Stahl aus Eisenerzen und Eisenoxyden in Regeneratorkamrnen. (D. P. Nr. 4135 v. 12. Sept. 1877. Ber. Chem. Ges. II. 2159. London.)  
L. H. E. CLAIR ANTHONY, Fabrication von Gas. Destillation von Knochen, Thier- und Fäcalstoffen, zuerst bei 300–500°, giebt Leuchtgas; dann Einleiten von Dampf auf die zum Glühen erhitzten Kohlen, giebt ein Gemenge von H und CO zu Heitzzwecken. (E. P. Nr. 4134 v. 6. Nov. 1877. Toulouse.)  
E. W. PARSELL, Fabrication von kaustischen Alkalien. Einwirkung des Aetzkalkes auf die Carbonlösung unter höherem Druck, wodurch die Dichte der letzteren auf 1,200 (statt 1,11) gesteigert werden darf. (E. P. Nr. 4144 v. 7. Nov. 1877. Liverpool.)

1830 1850 1856 1862 1879 1883 1897 1902 1919 1924 ...

Vergleichende Analyse der *Anagyris foetida* (eines Surrogats der *Senna*), des *Cytisus alpinus* und der *Cornilla varia*; von PESCHIER und JACQUEMIN.  
Das wesentlichste Resultat der nachfolgenden, von PESCHIER und JACQUEMIN angestellten, Untersuchungen ist, dass alle drei oben genannte Pflanzen einen gemeinschaftlichen, mit dem der Sennesblätter übereinkommenden, Stoff, der ihnen purgirende Eigenschaften verleiht, besitzen, das Cytisin oder Cathartin, in sofern beide Principe nach den Untersuchungen dieser Chemiker identisch sind (vergl. den nachfolgenden Artikel).  
I. Die *Anagyris foetida* ist ein, zur Familie der Leguminosen gehöriger, Strauch im südlichen Europa, dessen Blätter frisch einen stinkenden Geruch verbreiten, und von einigen Aerzten als Surrogat der Sennesblätter angewandt werden \*.

Title changes in „Chemisch-Pharmaceutisch Central-Blatt“

Systematische Inhalts-Uebersicht.  
I. Allgemeine und physikalische Chemie.  
Diffusion von Gasen durch poröse Körper 225. Ueber das Verhalten des Inductionsfunkentrommes in verschiedenen Gasen 463. Ueber die Durchdringlichkeit des Eisens für Gas in höherer Temperatur 654. Ueber die molekulare Bewegung der Gase 374. Chemische und Volta'sche Wärme 109. Einwirkung von Hitze auf Flüssigkeiten: Gasgehalt eine Bedingung des Siedens 209. Ueber Kältemischungen 591, 1111. Ueber die Hervorbringung von Schmelztemperaturen 936. Bestimmung hoher Temperaturen mittels eines thermoelektrischen Platin-Palladiumpaars 285, 296, mittels des Gaspyrometers 292. Wechselwirkung zwischen chemischer und mechanischer Kraft 321. Ueber Lösung 606.  
II. Anorganische Chemie.  
Allgemeines. — Ueber die sogenannte Passivität der Metalle 424. Ueber die Sechswertigkeit des Eisens und Aluminiums 337. Ueber die Atomigkeit des Sauerstoffs, Schwefels, Selen und Tellurs 382. Die Brechungsquotienten einer Lösung von salpetersaurem Wismuthoxyd 816. Ueber die Zusammensetzung der unorga-

Introduction of a systematic table of contents

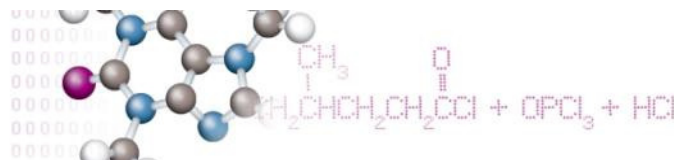
CZ splits in two parts: „technical“ and „scientific“

Holding passes to the *Deutschen Chemischen Gesellschaft*  
Introduction of Patent Register

CZ is re-unified in one journal of abstracts



|          |         |         |   |
|----------|---------|---------|---|
| -7.2209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |



# Chemisches Zentralblatt: History (2)

Introduction of Formula Index

Berlin wall does not disturb production of CZ

Editorial Office East Berlin

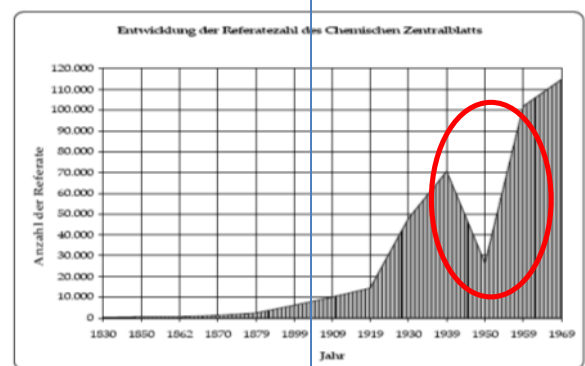
Introduction of SRD „Schnell Referat Dienst“ („quick abstract service“)



Double production of CZ in East and West Germany



... 1925      1940-45      1947-49      1950      1961      1967      1969



Difficulties in collecting information

Reunification of the journal under East and West German organizations

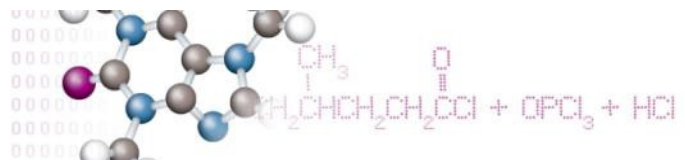


GDR office declares unable to afford production of SRD and of journal  
Chemisches Zentralblatt ceases publication





|          |         |         |        |
|----------|---------|---------|--------|
| -7.3209  | -1.0003 | -0.4030 | 0.0000 |
| -9.9889  | -0.1947 | -0.3218 | 0.0000 |
| -10.2795 | 2.5700  | 0.1393  | 0.0000 |
| -8.5096  | -1.8624 | -0.6298 | 0.0000 |
| -7.5171  | -2.3313 | 0.2821  | 0.0000 |
| -7.0177  | -1.8762 | -1.3207 | 0.0000 |
| -4.4781  | 1.2434  | 0.1858  | 0.0000 |



# Chemisches Zentralblatt: Content



- 140 years publication
- 900,000 pages
  - 700,000 contain ca. 2 million abstracts
  - 200,000 are indexes
  - 120,000 cover time period 1830-1907

**Ankündigung.**

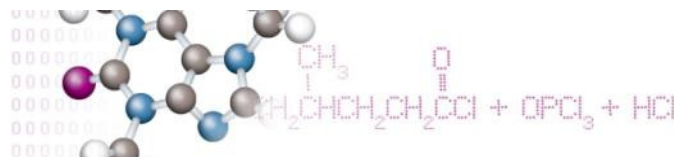
Wenn gleich die grosse Anzahl schon bestehender deutscher pharmaceutischer Zeitschriften dem Publicum schien genügen zu können, so hat doch eben diese Vervielfältigung derselben eine wechselseitige Beschränkung ihres Nutzens nothwendig mit sich führen müssen, indem sie einerseits eine grosse Zersplitterung der einheimischen Untersuchungen, andererseits eine häufige Wiederholung ausländischer Abhandlungen darin zur Folge gehabt hat. Rechnen wir hinzu, dass mehrere dieser Zeitschriften auch solche Abhandlungen und Notizen aufnehmen, die, wiewohl an sich werthvoll, doch nur von sehr entferntem Interesse für den Pharmaceuten sind, und dass Originaluntersuchungen, welche den wesentlichen Reichthum jener Zeitschriften bilden sollen, neben den, für den Pharmaceuten wirklich wichtigen, Resultaten oft ein schwer übersehbares Detail von Nebendetails enthalten müssen, deren Kenntniss für sein Fach von keinem wesentlichen Nutzen ist, so wird man leicht erkennen, dass das Bedürfniss so Vieler, mit dem möglichst geringsten Aufwande von Zeit und Mitteln den Fortschritten der Pharmacie in ihrem ganzen Umfange zu folgen, nicht hinreichend dadurch befriedigt wird; und diess lässt uns hoffen, dass ein Unternehmen, welches eine Abhülfe für dieses Bedürfniss verspricht, Beifall und Theilnahme beim pharmaceutischen und selbst dem ärztlichen Publicum, in so weit seine wissenschaftlichen Interessen mit denen des erstern zusammenfallen, finden werde.

Der Zweck dieses Unternehmens, welchem wir den obigen Titel beilegen, wird in Bezug zur Pharmacie ein ähnlicher seyn, als der Zweck  
1r Jahrgang. 1

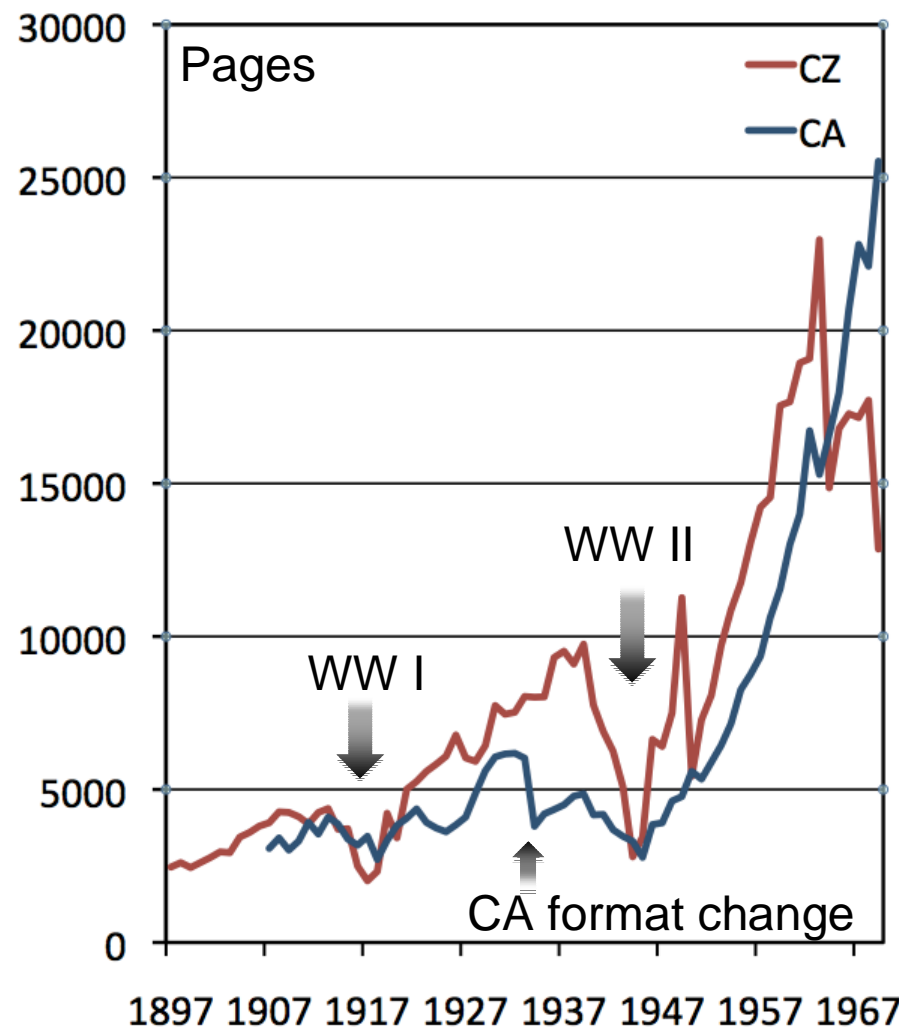
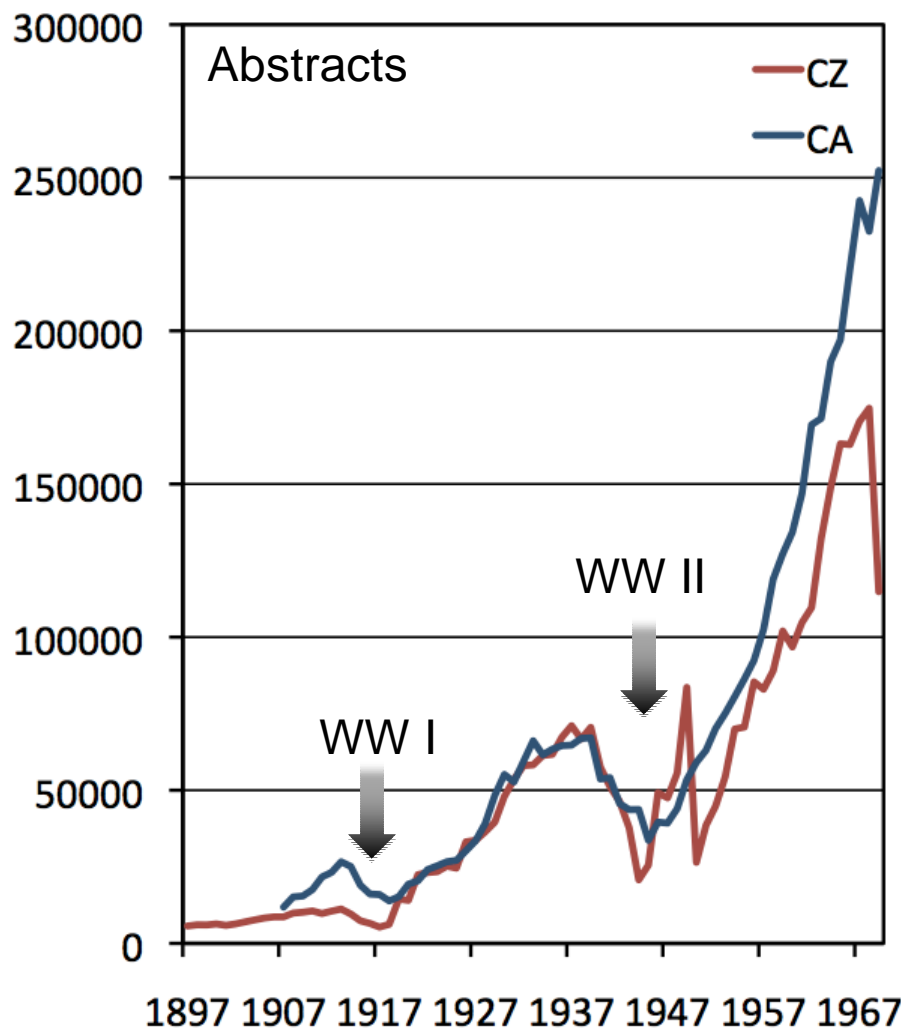




|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | C |
| -9.9889  | -0.1947 | -0.3218 | O |
| -10.2795 | 2.5700  | 0.1393  | C |
| -8.5096  | -1.8624 | -0.6298 | H |
| -7.5171  | -2.3313 | 0.2821  | H |
| -7.0177  | -1.8762 | -1.3207 | H |
| -4.4781  | 1.2434  | 0.1858  | C |

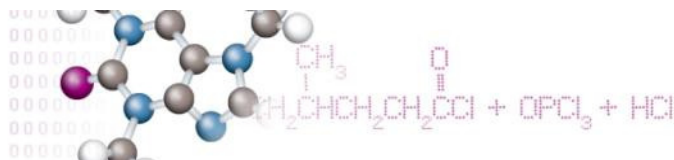


# Chemisches Zentralblatt vs. CA: Quantity





|          |         |           |
|----------|---------|-----------|
| -7.3209  | -1.0003 | -0.4030 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 |
| -4.4781  | 1.2434  | 0.1858 C  |

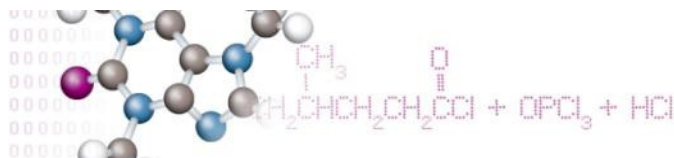


## Chemisches Zentralblatt vs. CA: Quality

- Many textbooks on chemical literature claim better quality of Chemisches Zentralblatt than CA for pre-WW II
  - H. Skolnik, *The literature matrix of chemistry*, 1982: „outstanding A/I service“
  - R.E. Maizell, *How to find chemical information*, 3rd ed. 1998, citing E.J. Crane, „[...] has value because of [...] good abstracts“
  - M. Mücke, *Die chemische Literatur*, 1982, „Zwar war CA zahlenmässig [...] dem Chemischen Zentralblatt überlegen, doch war dies gerade umgekehrt, was die Qualität der Referate betraf.“
  - R.T. Bottle, J.F. Rowland, *Information Sources in Chemistry*, 4th ed. 1993, „Before WW II, many chemists regarded CZ as superior in coverage to CA; its abstracts were longer and more informative [...]“

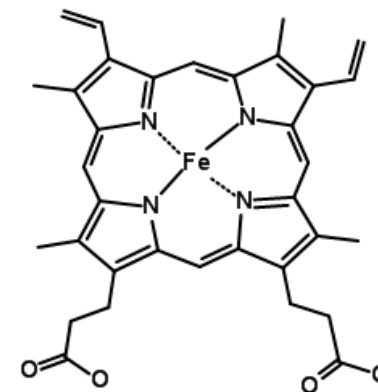


|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | O |
| -9.9889  | -0.1947 | -0.3218 | O |
| -10.2795 | 2.5700  | 0.1393  | C |
| -8.5096  | -1.8624 | -0.6298 | H |
| -7.5171  | -2.3313 | 0.2821  | H |
| -7.0177  | -1.8762 | -1.3207 | H |
| -4.4781  | 1.2434  | 0.1858  | C |



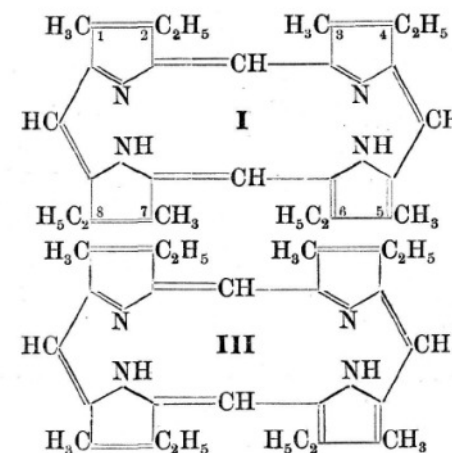
## Chemisches Zentralblatt vs. CA: Quality

Example: Hans Fischer, Georg Stangler, Synthese des Mesoporphyrings, Mesohämins und über die Konstitution des Hämins, Justus Liebigs Ann. Chem. 459(1927), 53-98.



|                    | CZ I(1928), 528 | CA 22:11339 (1928), 1363 |
|--------------------|-----------------|--------------------------|
| Length (pages)     | 7.5             | 1                        |
| Length (words)     | 3,882           | 690                      |
| Length (chars)     | 24,308          | 4,695                    |
| Compounds          | ~ 120           | ~ 70                     |
| Structure formulas | Yes             | No                       |

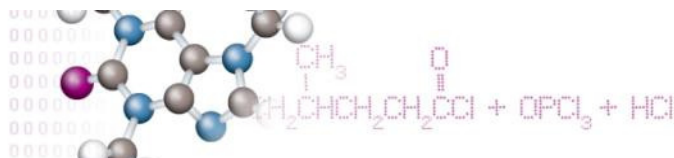
1928. I. D. ORGANIS







|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |

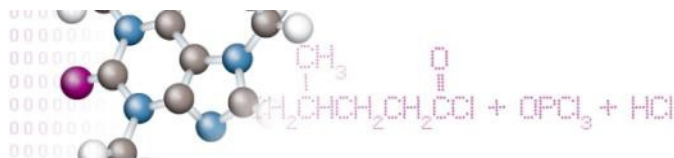


## Chemisches Zentralblatt: Conclusion

- Covers chemical literature from 1830 to 1969
- Very good abstract quality
  - Better quality (length, details) than CA for pre-WW II period 1907-1940
- Contains also important patent information
- Invaluable information in indexes (e.g. synonyms of ancient chemical names)
- Only comprehensive abstract journal on the market up to 1907
  - More comprehensive than CA for 19th century literature
  - Complements Beilstein and Gmelin handbooks for 19th century literature



-7.3209 -1.0003 -0.4030 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  
-4.4781 1.2434 0.1858 C



## Importance of Chemisches Zentralblatt: Example

Org. Lett., 2006, 8 (19), pp 4279–4281



*N*-Substituted diaza[12]annulenes are obtained by one-pot reaction of *N*-(2,4-dinitrophenyl)pyridinium chloride with amines in moderate to high yields. The <sup>1</sup>H NMR spectrum reveals that diamagnetic ring current is generated in the diaza[12]annulene ring. The *N*-substituted diaza[12]annulenes are electrochemically active in solution.

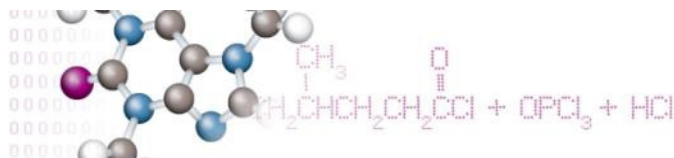
The authors have retracted this paper on November 15, 2007 (Org. Lett. 2007, 24, 5139)

Chemisches Zentralblatt., 1904, 2, 1145

**Th. Zincke, Über Dinitrophenylpyridiniumchlorid und dessen Umwandlungsprodukte.** Das Dinitrochlorbenzol-(2,4,1) verbindet sich mit Pyridin direkt zu einem gut charakterisierten quaternären Ammoniumsalz, dem *Dinitrophenylpyridiniumchlorid*, C<sub>5</sub>H<sub>5</sub>N<sup>v</sup>(Cl)C<sub>6</sub>H<sub>3</sub>(NO<sub>2</sub>)<sub>2</sub>. Diese Pyridinverb. reagiert in sehr eigenartiger Weise mit Alkali und mit verschiedenen primären und sekundären Aminen; durch Alkali wird sie in einen schön rot gefärbten Körper umgewandelt, der die Zus. der freien Ammoniumbase besitzt, während die Amine mit großer Leichtigkeit eine Spaltung unter B. von Dinitranilin herbeiführen.



|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |

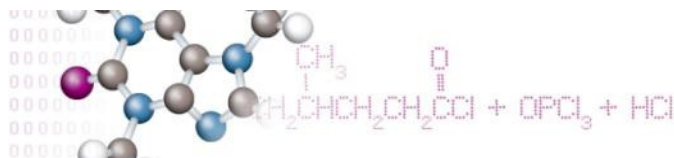


# Chemisches Zentralblatt: Digitalization

- Relevant for documentation of prior art
- Continuous and growing demand of the information
- FIZ Chemie Berlin has scanned the whole work and offers a full text searchable database for the web and the dataset for integration in Intranets
- ETH Zurich has bought the digitalized raw material (pdfs with OCRed text in the background) from FIZ and is creating a database offering full text search
- CAS has performed automatic translation (German → English) of the 1897-1907 volumes and included them in CAplus



|          |         |           |
|----------|---------|-----------|
| -7.3209  | -1.0003 | -0.4030 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 |
| -4.4781  | 1.2434  | 0.1858 C  |



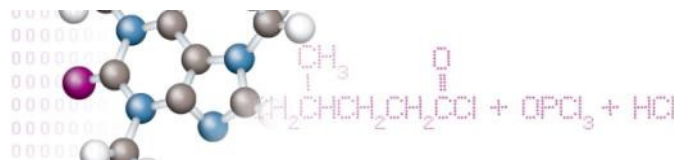
## InfoChem Motivation

- Text search in Chemisches Zentralblatt:
  - Abstracts in German language
  - High number of old German chemical names
- Chemists think in structures!!!
- Language independent structure search would help ALL scientists to access this historical source and to use the relevant information of this art
- Required technology for structure search projects
- Optimized German-English dictionaries
- 30 million SPRESI names

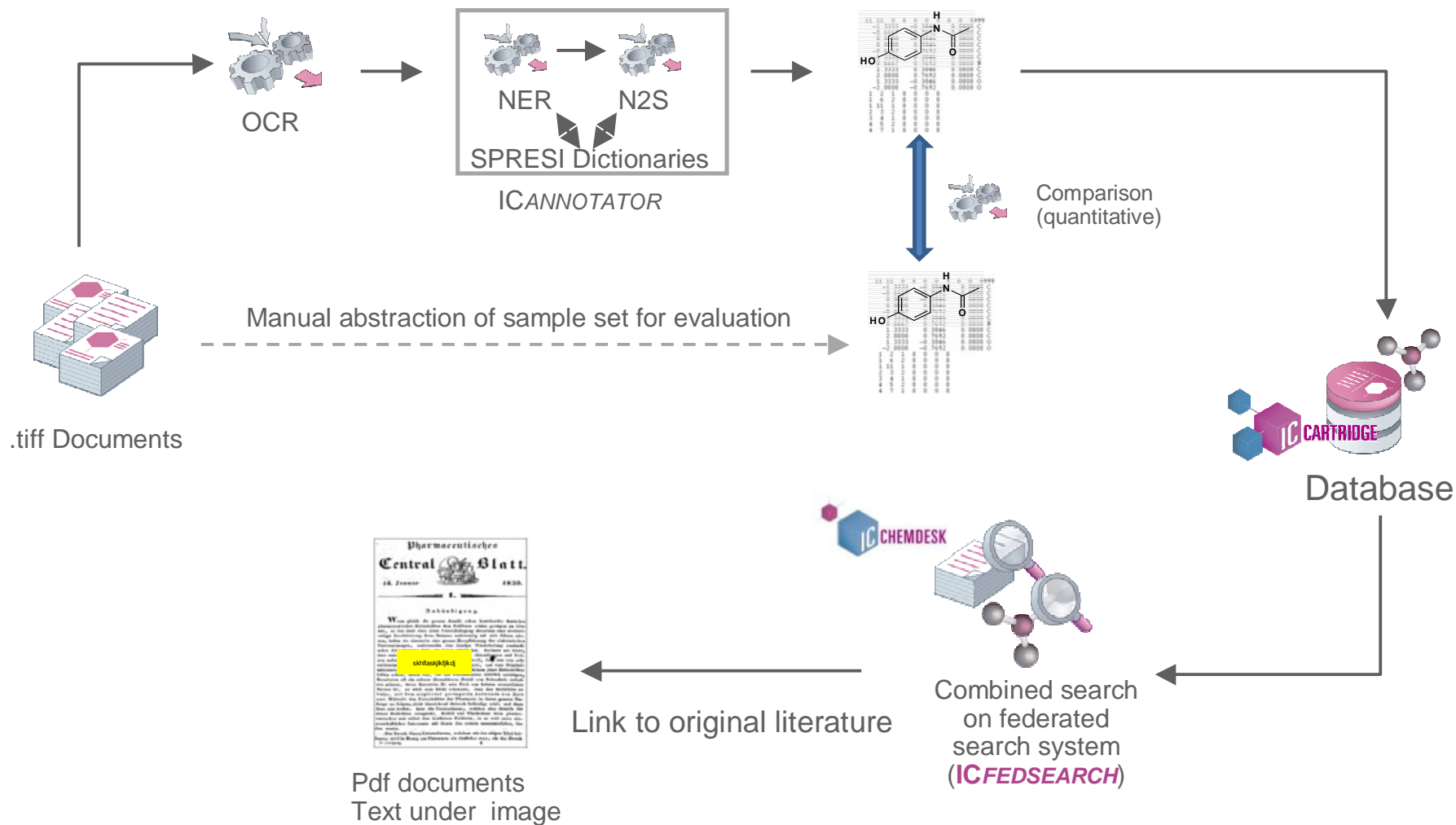




|          |         |         |   |
|----------|---------|---------|---|
| -7.2209  | -1.0663 | -0.4030 | C |
| -9.9889  | -0.1947 | -0.3218 | O |
| -10.2795 | 2.5700  | 0.1393  | C |
| -8.5096  | -1.8624 | -0.6298 | H |
| -7.5171  | -2.3313 | 0.2821  | H |
| -7.0177  | -1.8762 | -1.3207 | H |
| -4.4781  | 1.2434  | 0.1858  | C |



# Overview of Approach and Applied Technology



# Challenges OCR (1)

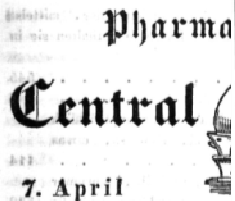
1830

1870

1910

1930

1969



**Pharma**

**Central**

7. April

**Dorrenner**

Da die Masse des aufzu Verhalten gegen den vorlufigen zu gewinnen anfangt, so Mittheilungen an das pharmazeutische Zeit, wie im gegenwartigen in fortlaufender Nummer mit der Zur Raumersparniß fur G Wissenschaften betreffen, ist zuglich die Personal-, Local-, Handels- an uns wir nochmals einladen, wie dieses schon beim vorigen Hefen in großem Interesse

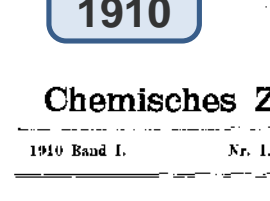


**Chemisches Z**

1910 Band I. Nr. 1.

**Allgemeine und physikalische Chemie**

**Inhalt.** Wochenbericht. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen.

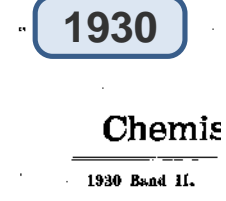


**Chemisches Z**

1910 Band I. Nr. 1.

**Allgemeine und physikalische Chemie**

**Inhalt.** Wochenbericht. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen.



**Chemisches Z**

1930 Band II.

**Chemie**

**Inhalt.** Wochenbericht. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen.



**CHEMISCHES ZENTRALBLATT**

1969

**Inhalt.** Wochenbericht. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen. — Ueber die Wirkung der Salpetersure auf die Oxidation von Eisen.

**Ueber Beschaffenheit und Zweckmaigkeit von pharmazeutischen Preparaten**  
 Von G. W. GRASSMANN

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Ueber Beschaffenheit und Zweckmaigkeit von pharmazeutischen Preparaten**  
 Von G. W. GRASSMANN

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Ueber Beschaffenheit und Zweckmaigkeit von pharmazeutischen Preparaten**  
 Von G. W. GRASSMANN

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Ueber Beschaffenheit und Zweckmaigkeit von pharmazeutischen Preparaten**  
 Von G. W. GRASSMANN

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

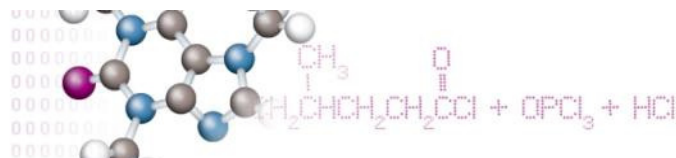
**Ueber Beschaffenheit und Zweckmaigkeit von pharmazeutischen Preparaten**  
 Von G. W. GRASSMANN

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

**Die Aufnahme der Birken Ruthenica und ihre hufige Anwendung als Jodmittel**  
 (vergl. S. 100.)

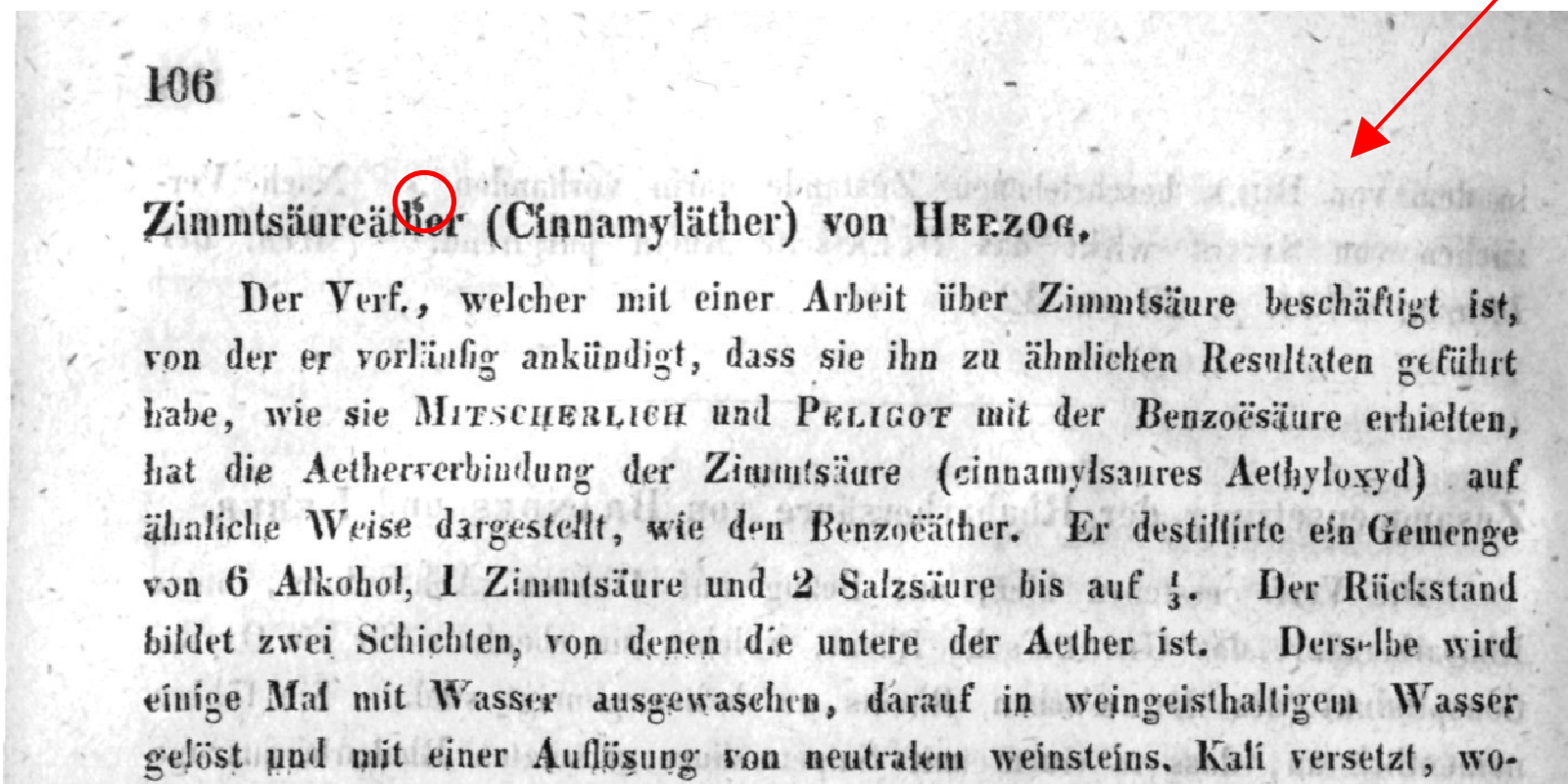


|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | C |
| -9.9889  | -0.1947 | -0.3218 | O |
| -10.2795 | 2.5700  | 0.1393  | C |
| -8.5096  | -1.8624 | -0.6298 | H |
| -7.5171  | -2.3313 | 0.2821  | H |
| -7.0177  | -1.8762 | -1.3207 | H |
| -4.4781  | 1.2434  | 0.1858  | C |



## Challenges OCR (2)

- Bad quality of original source: dirty (blotted, stained) pages  
print from back page





-7.3209 -1.0000 -0.4630 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  
-4.4781 1.2434 0.1858 C



## Challenges OCR (3)

- Tables: extremely small fonts, not recognizable begin / end of columns

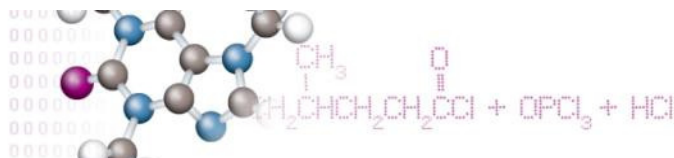
775

| 1 im Jodkalium enthaltenes Jod aufgelöst in folgender Quantität W. | Wirkung durch Stärkemehlösung mit Chlorzusatz                          | Platinechlorid  | Salpeters. Quecksilberoxydul             | Salpeters. Silberoxyd                                    |
|--|--|---|--|--|
| 25000  | dunkelblaue Färbung mit Niederschlag.                                  | nach einigen Min. kirschrothe Färbung, dann schwarzer Niederschlag. | grünlichgelbe Trübung.                   | anfangs schwache Trübung, dann weisslicher Niederschlag. |
| 50000  | blaue Färbung mit Niederschlag.  | nach einigen Min. blassrothe Färbung.                               | minder reichliche grünlichgelbe Trübung. | milchige Trübung.  |
| 100000   | blassblaue Färbung mit Niederschlag.                                   | nach 15 Min. blassrosenrothe Färbung.                               | milchige Trübung.                        | kaum merkliche Trübung.                                  |
| 200000   | sofort keine Färbung; allmählig ganz blassblau gefärbter Niederschlag. | nach 20 Min. ganz blassrosenrothe Färbung.                          | nichts.                                  | nichts.  |





|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |



## Challenges OCR (4)

- Ambiguous old fonts (h=b; c=e; ligations)

**Chlorcalcium und Chlorbley**  
**Zinnchlorürauflösung**

- Spaced text

Proteinstoff des Roggens (H. Ritt-  
hausen) 258.

Protocatechusäure (L. Barth) 360.

Pseudoamylharnstoff (A. Wurtz) 71.

Pseudodiamylenharnstoff (A.  
Wurtz) 72.

Respiration, Untersuchungen über die  
— beim Menschen (M. v. Pettenkofer  
und C. Voit) 289. 161.

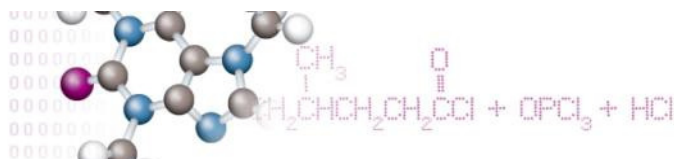
Reten (Berthelot) 818, 825.

Rhodankalium, Verhalten des —  
gegen Quecksilbersalze (J. Philipp) 967.

Specific rules, large German dictionaries and extensive training are applied to correct systematic mistakes of standard OCR process



|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |

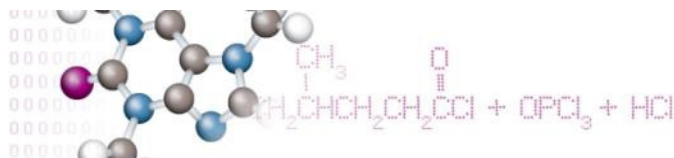


## Main Challenges Annotation

- Obsolete German language
  - Schwefelsaures Natrium, Chlorür, Bromür
- Historical names
  - Pelopeum → Columbium → Niobium
- Different spelling for the same name:
  - Dibrom... ↔ Bibrom...
  - Ätzkali ↔ Aetzkali
  - **Zimmtsäure** ↔ Zimtsäure



|          |         |         |
|----------|---------|---------|
| -7.3209  | -1.0003 | -0.4030 |
| -9.9889  | -0.1947 | -0.3218 |
| -10.2795 | 2.5700  | 0.1393  |
| -8.5096  | -1.8624 | -0.6298 |
| -7.5171  | -2.3313 | 0.2821  |
| -7.0177  | -1.8762 | -1.3207 |
| -4.4781  | 1.2434  | 0.1858  |

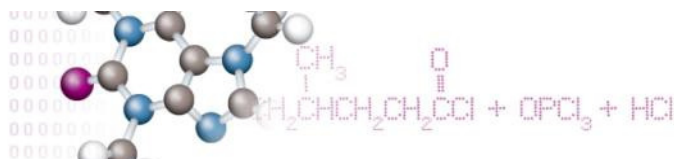


## Solutions in Annotation Process

- Correction of German-specific grammar
- Translation in English of not available chemical names
- Research in old sources:
  - Beilstein
  - Brockhaus Encyclopedia
  - German-English dictionaries of chemistry
  - Meyers Encyclopedia
  - Pierer Encyclopedia
  - References to very old books, journals, articles
    - “Naturwissenschaftliche Exzerpte und Notizen Mitte 1877 bis Anfang 1883”  
by Karl Marx



|          |         |         |
|----------|---------|---------|
| -7.3209  | -1.0000 | -0.4030 |
| -9.9889  | -0.1947 | -0.3218 |
| -10.2795 | 2.5700  | 0.1393  |
| -8.5096  | -1.8624 | -0.6298 |
| -7.5171  | -2.3313 | 0.2821  |
| -7.0177  | -1.8762 | -1.3207 |
| -4.4781  | 1.2434  | 0.1858  |



## Results Annotation Chemisches Zentralblatt

- 120,000 pages covering time period 1830-1907
- 2.4 million chemical names with associated structure
  - 98,000 unique names
  - 47,000 unique structures

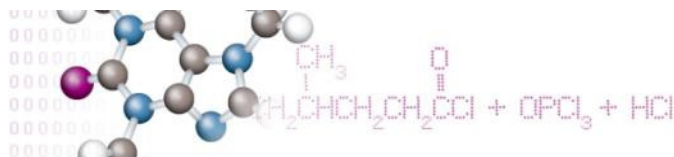
### Quantitative comparison with manually abstracted sample set

- Recall                    51%
- Precision                87%





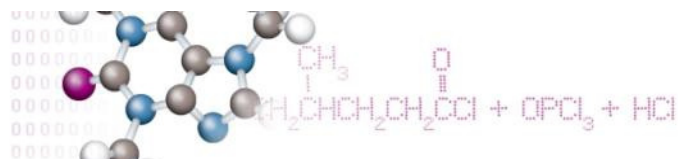
-7.3209 -1.0003 -0.4030 C  
-9.9889 -0.1947 -0.3218 O  
-10.2795 2.5700 0.1393 C  
-8.5096 -1.8624 -0.6298 H  
-7.5171 -2.3313 0.2821 H  
-7.0177 -1.8762 -1.3207 H  
-4.4781 1.2434 0.1858 C



# Federated Search Prototype



-7.3209 -1.0003 -0.4030 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  
-4.4781 1.2434 0.1858 C



# Federated Search Prototype

http://icfedsearch.infochem.de:8080/gmrw30/app/?wicket:interface=:1:17::: - Windows Internet Explorer

http://icfedsearch.infochem.de:8080/gmrw30/app/?wicket:interface=:1:17:::

Favoriten http://icfedsearch.infochem.de:8080/gmrw30/

|      |               |                                     |             |
|------|---------------|-------------------------------------|-------------|
| 1527 | <p>chiral</p> | <a href="#">Lactoflavin</a>         | C17H20N4O6  |
| 1788 | <p>chiral</p> | <a href="#">Riboflavinphosphat</a>  | C17H21N4O9P |
| 1887 | <p>chiral</p> | <a href="#">Flavinmononucleotid</a> | C17H21N4O9P |



|          |         |         |   |
|----------|---------|---------|---|
| -7.2209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |



# Federated Search Prototype

http://icfedsearch.infochem.de/CZB/czb\_1935\_Sachregister\_5020.pdf#search='Lactoflavin' - Windows Internet Explorer

http://icfedsearch.infochem.de/CZB/czb\_1935\_Sachregister\_5020.pdf#search='Lactoflavin'

Verdaulichk. u. Nährwert v. —Heu, —Mehl u. —Blattmehl II 145; Zus. v. —Heu, Frage d. Fütter. v. —Mehl I 2746; Einsäuer. I 325; s. auch *Leguminosen*.

**Luzerne** s. *Wachse*.

**Lycoperdon** s. *Pilze*.

**Lycopin**, Vork. in blaßroten Grapefruits II 2078; Bldg. in Tomaten (Bezieh. zur Temp.) II 2752; Isolier.: aus d. ind. roten Paradiesapfel (*Citrus grandis* Osbeck) I 95; aus Diospyros- u. Arbutusfrüchten II 2686; aus d. Früchten v. *Rosa rugosa* II 707; neuer Farbstoff mit —Spektr. aus *Solanum dulcamara* II 1557; Vork.: in d. Carotinoidfrakt. v. Thiocystisbakterien I 2022; in Butter II 3452; —Geh.: im Menschenfett I 3151; d. menschl. Leber II 2972.

Trenn. v. d. Carotinen durch Adsorpt. I 1066; Absorpt.-Spektr. II 30; (bei gewöhnl. Temp. u. d. Temp. d. fl. Luft) II 2961; Viscositätsmess. I 3530.

**Lycorenin**, über — II 1182.

**Lycorin** (F. 250°), Unters. über — II 1181; Konst. II 3111; Isolier. aus *Eurycles amboinensis*, Salze II 858.

**Lymax**, Verwend. gegen *Hylecoetus dermestoides* u. *Anobium pertinax* II 2309.

**Lymphdrüsen** s. *Drüsen*.

**Lymphhe**, S-Geh. d. Pankreasblutes u. d. Thoracicus— (Beeinfluss. d. S-Haushaltes d. Thoracicus— nach d. Pankreasextirpat.) I 918; (Beeinfluss. d. S-Haushaltes d. Thoracicus— d. Hundes mit unterbundenem Pankreasgang u. d. n. Hundes nach

wachstumswrkg. I 99; (synon. Flavinpräpp.) II 2540; Flavinbilanz im Tierkörper (Flavinausscheid. durch d. Faeces u. durch d. Harn) II 715; (n. Flavinstoffwechsel u. Flavinstoffwechsel bei B<sub>2</sub>-Avitaminose) II 2541; Bldg. v. Flavincarbon-säuren aus Vitamin B<sub>2</sub> im Organism. I 2191; Dehydrier. v. Hexosemonophosphorsäure durch Methylenblau in Ggw. v. Zwischenferment, Flavoprotein u. Co-Ferment I 2033.

**Bestimmung**, Best. I 3425; objektive Meth. zur Best. d. Fluorescenz II 875; s. auch *Enzyme-Atmungsfermente*; *Vitamine-Vitamin B<sub>2</sub>*.

**Einzelne Flavine.**

**Aquoflavin**, Vork. im Harn I 3425; Isolier. aus Harn, Einw. v. Licht, Acetylier. I 722; Frage d. Bldg. aus Vitamin B<sub>2</sub> im Organismus I 2191.

**Aquolumiflavin**, Bldg. aus Aquoflavin, Eigg. I 722.

**Hepaflavin (Hepatoflavin)**, Nichtidentität mit Vitamin B<sub>2</sub> I 3563; Wachstumswrkg. I 99; Wrkg. bei Avitaminose B II 2231.

**Hepatohctoflavin**, potentiometr. Unters. (Vergleich mit Lumilactoflavin) II 58.

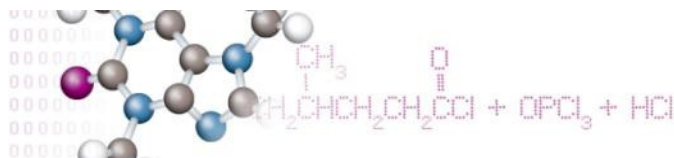
**Lactoflavin** s. *Vitamine-Vitamin B<sub>2</sub>*.

**Leukolactoflavin**, Einw. v. Licht I 722.

**Lumichrom (6,7-Dimethylalloxazin)**, (Zers. 360°), Synth., Eigg., Bldg. aus Lumilactoflavin Methylier. I 723; Reindarst. aus Lactoflavin, Eigg., Lichtabbau I 3551; Bldg.: aus Lactoflavin I 1396; bei d. Bestrahl. v. Lactoflavin (Eigg., Konst.) I 86; (Konst.) I 724; durch Photolyse v.



|          |         |           |
|----------|---------|-----------|
| -7.3209  | -1.0003 | -0.4030 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 |
| -4.4781  | 1.2434  | 0.1858 C  |



## Summary

- Described history, content and importance nowadays of Chemisches Zentralblatt
- Illustrated how the challenges of OCR and annotation process have been solved
- Time period 1830-1907 contains 98,000 unique names and 47,000 unique structures
- Quantitative comparison proves over 50% recall and nearly 90% precision
- Generated structure searchable Chemisches Zentralblatt database is integrated in IC<sub>FEDSEARCH</sub>



|          |         |           |
|----------|---------|-----------|
| -7.2209  | -1.0003 | -0.4030 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 |
| -4.4781  | 1.2434  | 0.1858 C  |

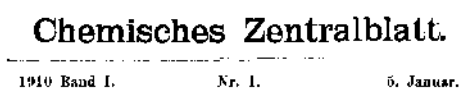
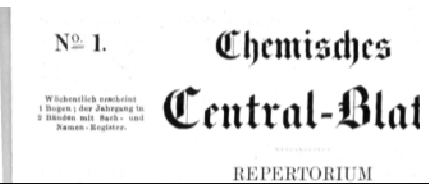


# Outlook



CHEMISCHES  
ZENTRALBLATT  
SINCE 1830

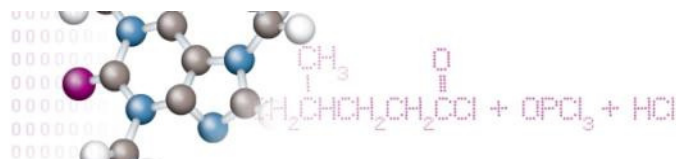
| Chemisches Zentralblatt: | Phase 1, Q2 2009 | Phase 2, Q4 2009 |
|--------------------------|------------------|------------------|
| Pages:                   | 120,000          | 900,000          |
| Time period:             | 1830-1907        | 1830-1969        |
| Unique names:            | 98,000           | Ca. 1 million    |
| Unique structures:       | 47,000           | Ca. 500,000      |
| Recall:                  | 50%              | ?                |







|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | 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 |
| -4.4781  | 1.2434  | 0.1858  | C |



## Acknowledgements

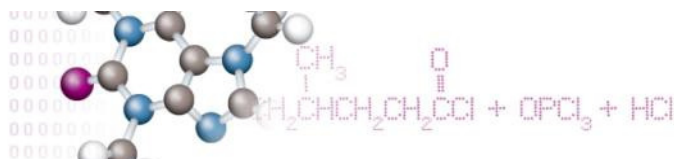
- Prof. Dr. Deplanque, Mr. Heineke and FIZ Chemie Team Berlin
- Ms. Langanke
- Dr. Brändle (ETH Zürich)
- InfoChem Team

# Thank you!

**InfoChem GmbH:** [www.infochem.de](http://www.infochem.de), [www.spresi.com](http://www.spresi.com), [info@infochem.de](mailto:info@infochem.de)



-7.3209 -1.0003 -0.4030 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  
-4.4781 1.2434 0.1858 C



# Chemisches Zentralblatt: Structure

## Abstracts (Referate)

Short reports of national and international progresses in chemistry

→ **Purpose of CZB**

1870-1886  
Wochenberichte

1884  
Figures are comprised in the abstracts

1964  
Numbering of abstracts

## Patents (Patente)

1879  
Introduction (end of abstracts part)

1921  
Patent Numbers country-specific

1938  
Patents from 20 Nations abstracted

1964  
Numbering of patent abstracts

## Indexes (Register)

1830      Autorenregister  
            Sachregister

1863      Systematische Register

1897      Patentregister

1924      Patentrückzitateliste

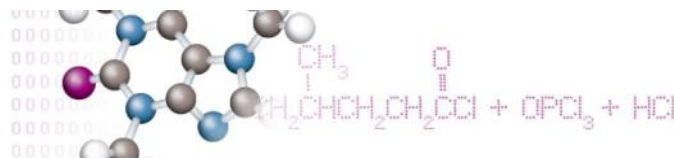
1925      Formelregister

1883      Generalregister

⇒ Aim of all indexes is to facilitate the search of relevant information  
⇒ They are the first attempt of modern search possibilities

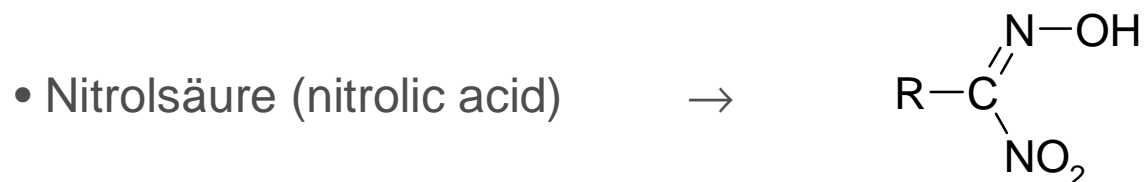


|          |         |         |   |
|----------|---------|---------|---|
| -7.3209  | -1.0003 | -0.4030 | C |
| -9.9889  | -0.1947 | -0.3218 | O |
| -10.2795 | 2.5700  | 0.1393  | C |
| -8.5096  | -1.8624 | -0.6298 | H |
| -7.5171  | -2.3313 | 0.2821  | H |
| -7.0177  | -1.8762 | -1.3207 | H |
| -4.4781  | 1.2434  | 0.1858  | C |



## Challenges Annotation (1)

- Names lack position, valence or stoichiometric information
  - Pimarsäure → is it the R or L form?
  - Platinchlorid → in which oxidation state II, III, IV?
- Chemical names that indicate a chemical class or a mixture



NO solution: information for correct structure is not available in the original source