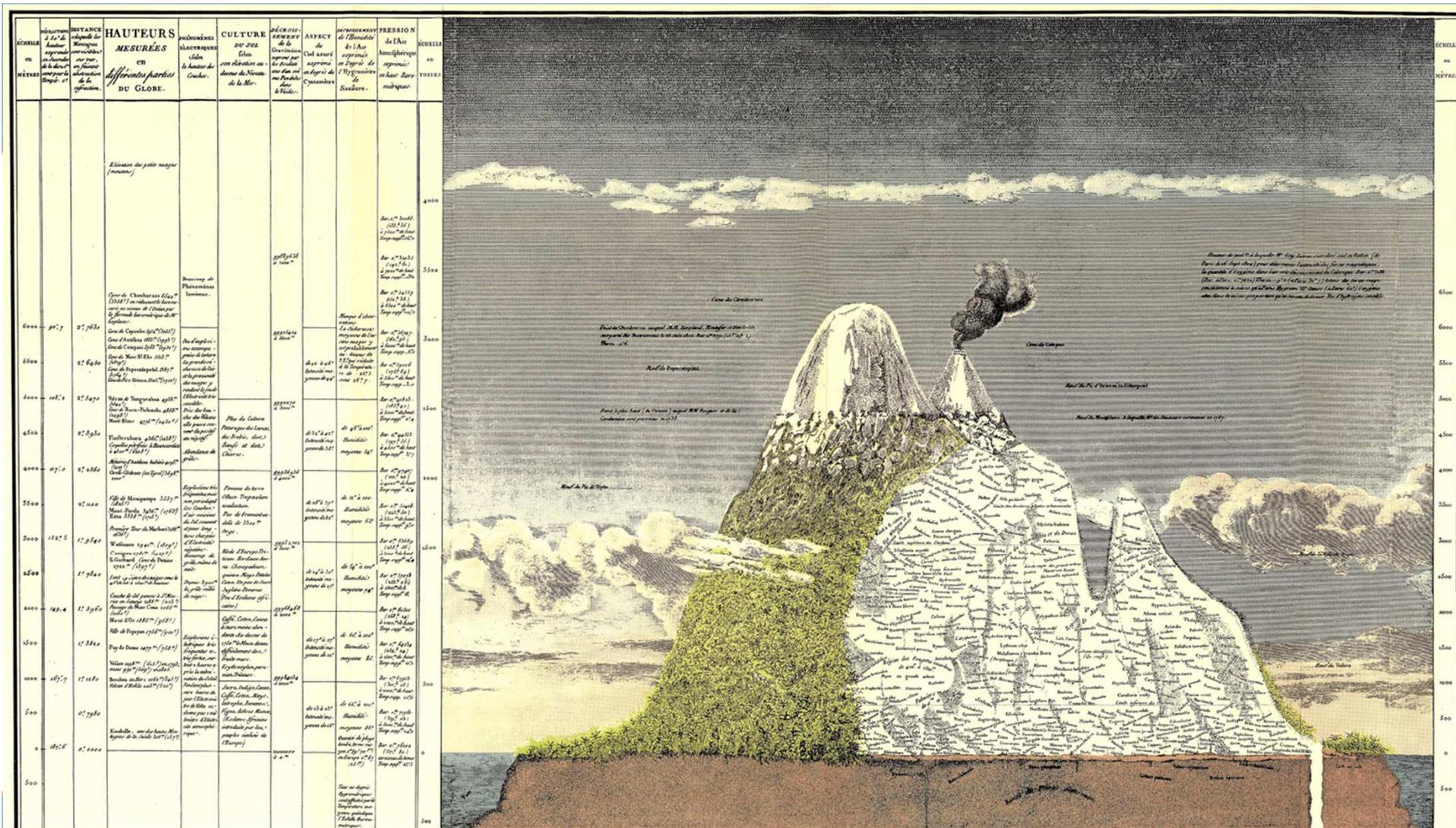


**Die Vermessung**

**der Welt**

**ERSCHLIESSUNG DES  
DEUTSCHSPRACHIGEN WWW**





# Unsere spezielle Welt



# Vermessung 1 – Maschinelles Lernen zur Emailklassifikation

Hallo Support,  
Ich habe Ihr Handy ein Mal beim Tauchen  
dabei gehabt. Seitdem funktioniert  
es nicht mehr!  
Das finde ich ungeheuerlich, bei  
so einem wasserdichten Handy!

Ein enttäuschter Kunde  
Hans Mustermann

- Projekt für Siemens
- „Lesen“ aller info@ Emails
- Zuweisung zum richtigen Mitarbeiter und Überwachung der Beantwortung
- Klassifikation
  - Erlernen der verschiedenen Themen aus Beispielen
  - Sprache, Tonalität, Fachbereich

→ Ähnlich maschinelle Sachgruppenvergabe

...  
Spanisch  
**Deutsch**  
Englisch  
...

...  
Produktanfrage  
**Beschwerde**  
...

...  
**Mobiltelefon**  
...

**Trainingsmengen:**  
alte Emails je Unternehmensbereich  
**Mess-Dimensionen**  
Sprache, Intention, Thema  
**Prozessoptimierung:**  
Emailbeantwortung

## Vermessung 2: Kontextsensitive Produktwerbung



Three product advertisements are displayed in a row. The first is for a 'Campingaz Gasgrill Exper' with a price tag 'ab 179,99 €' and a link 'jetzt bei baur.de'. The second is for a 'PERGART Gewächshaus' with a price tag 'ab 299,00 €' and a link 'jetzt bei Gamani.de'. The third is for 'Kriech- & Feldrose (Rosa)' with a price tag 'ab 2,85 €' and a link 'Pflanzenwelt-Biermann'. A vertical logo 'billiger.de' is on the right side of the ads.

### Thematische Webseiten

- ▣ Blogs
- ▣ Portale
- ▣ Kleinanzeigenportale
- ▣ ...

### → Thematisch optimierte Werbung



GÄRTNERN

*Willkommen im Zufallsgarten*

14. September 2019 · Carmen Feldhaus · Kommentar hinterlassen

Wir Gärtner haben das große Glück, ein Fleckchen Erde zu besitzen, das wir nach unseren Wünschen und Möglichkeiten gestalten können.

### Trainingsmengen:

Produkte nach Kategorien

### Anwendung:

Webseiten von Werbepartnern

### Mess-Dimensionen

Kategorie, Produkt-Begriffe

### Prozessoptimierung:

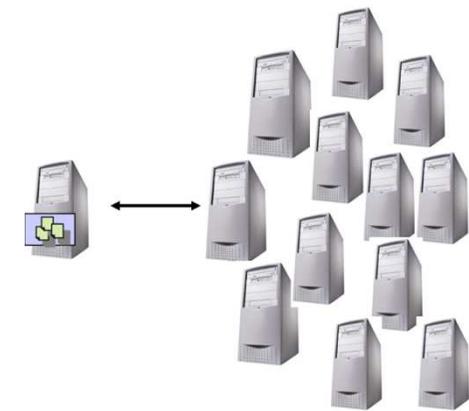
Themenauswahl Werbung

Produktauswahl Werbung



# Crawling contentDetection

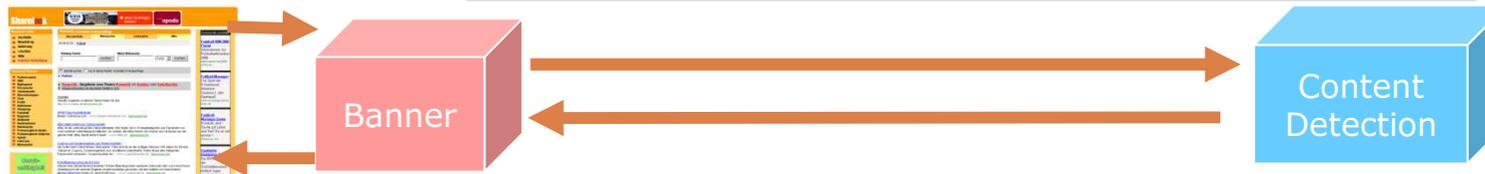
- ▣ Crawler ermittelt Inhalt der Webseite
- ▣ Werbung passt zur Webseite
- ▣ Crawlertechnik:
  - ▣ Asynchroner/synchroner Besuch der Webseite
  - ▣ Ermittlung der Kategorie
  - ▣ Extraktion von Keywords
  - ▣ Intelligentes Recrawling
- ▣ Crawlingleistung
  - ▣ > 100 Mio. Crawlings / Monat
- ▣ Werbemittel
  - ▣ > 1 Mrd. individualisierte Werbemittel / Monat
- ▣ Durch maschinelles Lernen weitgehend sprachunabhängig
  - ▣ sechs europäische Sprachen



Crawler

# Inhaltssteuerung der Werbemittel durch Crawling und maschinelles Lernen

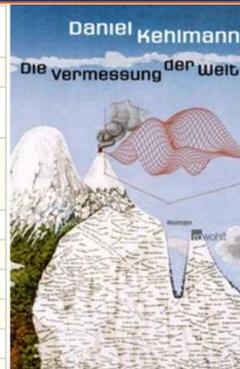
[http://.../contentdetection/responder?requesturl=  
http://www.mode-lexikon.de/mode/lexikon/l\\_4\\_Pullover.html](http://.../contentdetection/responder?requesturl=http://www.mode-lexikon.de/mode/lexikon/l_4_Pullover.html)



```
<?xml version="1.0" encoding="utf-8" ?> -  
<CD>-  
  <VAL>1122129310</VAL>  
  <CAT lv="1" pos="1">Mode</CAT>  
  <CAT lv="2" pos="1">Pullover&Strick</CAT>  
  
  <KEY pos="1">Pullover</KEY>  
  <KEY pos="2">Fashion</KEY>  
  <KEY pos="3">Sweatshirt</KEY>  
  <KEY pos="4">Rollkragenpullover</KEY>  
  <KEY pos="5">Marc O Polo</KEY>  
</CD>
```

# Erschließung Nationalbibliothek

Link zu diesem Datensatz	<a href="http://d-nb.info/975618768">http://d-nb.info/975618768</a>
Art des Inhalts	Fiktionale Darstellung
Titel	Die Vermessung der Welt : Roman / Daniel Kehlmann
Person(en)	Kehlmann, Daniel (Verfasser) Hellmann, Walter (Einbandgestalter)
Organisation(en)	any.way (Sonstige)
Ausgabe	1. Aufl.
Verlag	Reinbek bei Hamburg : Rowohlt
Zeitliche Einordnung	Erscheinungsdatum: 2005
Umfang/Format	301 S. ; 21 cm
Andere Ausgabe(n)	Online-Ausg.: Kehlmann, Daniel: Die Vermessung der Welt
ISBN/Einband/Preis	978-3-498-03528-0 Pp. : EUR 19.90, sfr 34.90 3-498-03528-2 Pp. : EUR 19.90, sfr 34.90
EAN	9783498035280
Gestaltungsmerkmale	Gestalter: Hellmann, Walter / Bucheinband Gestalter: any.way / Bucheinband Gestaltung: Schutzumschlag
Schlagwörter	Humboldt, Alexander von ; Gauß, Carl Friedrich ; Belletristische Darstellung
Sachgruppe(n)	830 Deutsche Literatur ; B Belletristik
Literarische Gattung	Historische Romane und Erzählungen



<b>Schlagwörter</b>	Humboldt, Alexander von ; Gauß, Carl Friedrich ; Belletristische Darstellung
<b>Sachgruppe(n)</b>	830 Deutsche Literatur ; B Belletristik
<b>Literarische Gattung</b>	Historische Romane und Erzählungen

# Parallelen zur maschinellen Erschließung

## ❖ Sachgruppen

## ▣ Produktkategorien

- ▣ Ca. 40 Ebene 1 Kategorien  
Güte: 80-90%
- ▣ 250 Ebene 2 Kategorien  
Güte: 70-80%

## ❖ GND

- ❖ Personen
- ❖ Sachschlagwörter
- ❖ ...

## ▣ Begriffe

- ▣ Marken
- ▣ Keywords

## ❖ Einsatz von maschinellem Lernen

## ▣ Einsatz von maschinellem Lernen



# Unterschiede zur maschinellen Erschließung

- Lerngrundlage
  - Produktdaten
  - Entspricht nicht späterer Anwendung auf Webseiten!
  - Erfordert Expertenwissen, z.B. Apotheke → Medikamente
- Categoriesystem
  - kundenbezogen
  - Entwickelt sich monatlich weiter
- Begriffe
  - Offen – erweitert sich stetig
  - Nicht lexikal – „Babyhopser“
  - Expertenwissen auf Domänenebene „Focus“ / [focus.de](http://focus.de)



→ täglich/wöchentlich Nachklassifizieren

→ monatliches Nachlernen

## Unterschiede zur maschinellen Erschließung 2

### ▣ Maschinelles Lernen

- ▣ Kein Black-Box Verfahren – es ist erkennbar was gelernt wurde
- ▣ Eigenentwicklung ähnlich Bayes Netzwerk
- ▣ Überlagerung mit Expertenwissen dadurch möglich

### ▣ Kategorisierung ist ein Baustein des Gesamtprozesses „Passende Produktempfehlung“

#### ▣ Produktrelevanz



#### ▣ Saisonalität



### ▣ Kategorisierung mit Business-Wissen



- ▣ Schwierige Kategorie: Bücher – decken viele Themen ab  
Warum beim Garten-Blog keine Bücher?  
Schlechte Conversion im Internet
- ▣ Ähnlich: Spielzeug unter Antiquitäten vs. Aktuelles Spielzeug

# Vermessung 3 - Domänendatenbank DACH

## ■ Erschließung aller deutschsprachiger Domains D, A, CH

### ■ Ermittlung strukturierter Daten

- Thema der Seite
- Verortung: Impressum
- Grobe Reichweite
- Art der Seite
  - Blog, Forum, Shop
- Technische Gegebenheiten

Adressverlage  
- Adressaktualisierung

Handwerkerverzeichnisse  
• Internetaffine Handwerker

Werbetreibende

- Neue Partnerseiten

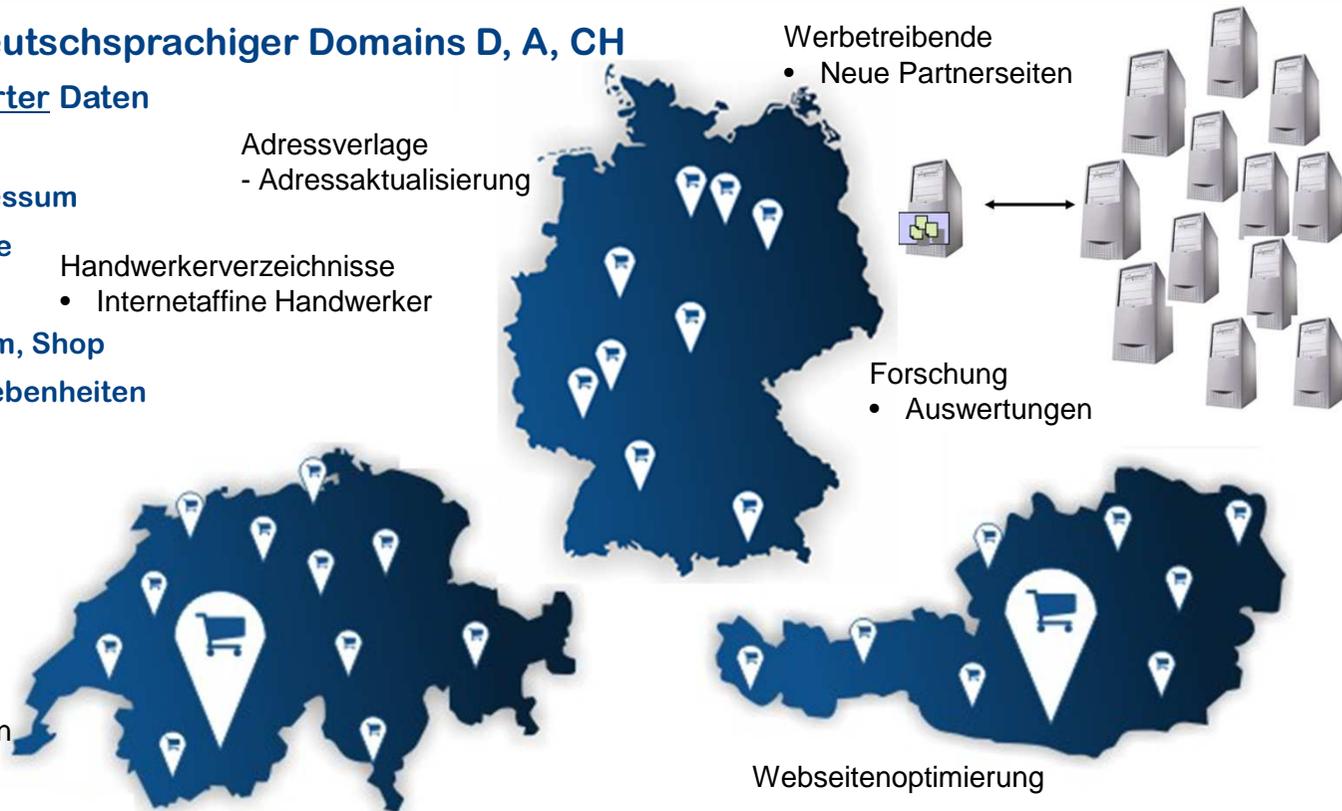
Forschung

- Auswertungen

Paketversender / Gütesiegelanbieter  
• Neue E-Commerce Shops

Pharmafirmen

- Ermittlung von Ärzten in bestimmten Themenbereichen



# Finden neuer E-Commerce Shops D, A, CH, FR, PL, ...



Große Shops:  
Zalando.de  
Mediamarkt.de  
Notebooksbilliger.de  
eventim.de  
Otto.de  
Redcoon.de

Mittelgroße eher Frauen:  
nelly.de  
kauf-unique.de  
gerstaecker.de  
meyer-mode.de  
dansommer.de  
Mamarella.com

Sport:  
tauchversand-online.de  
Plutosport.de  
Shop4runners.com  
Sportshop-triathlon.de  
Tennisherz.de  
Volleyballdirekt.de

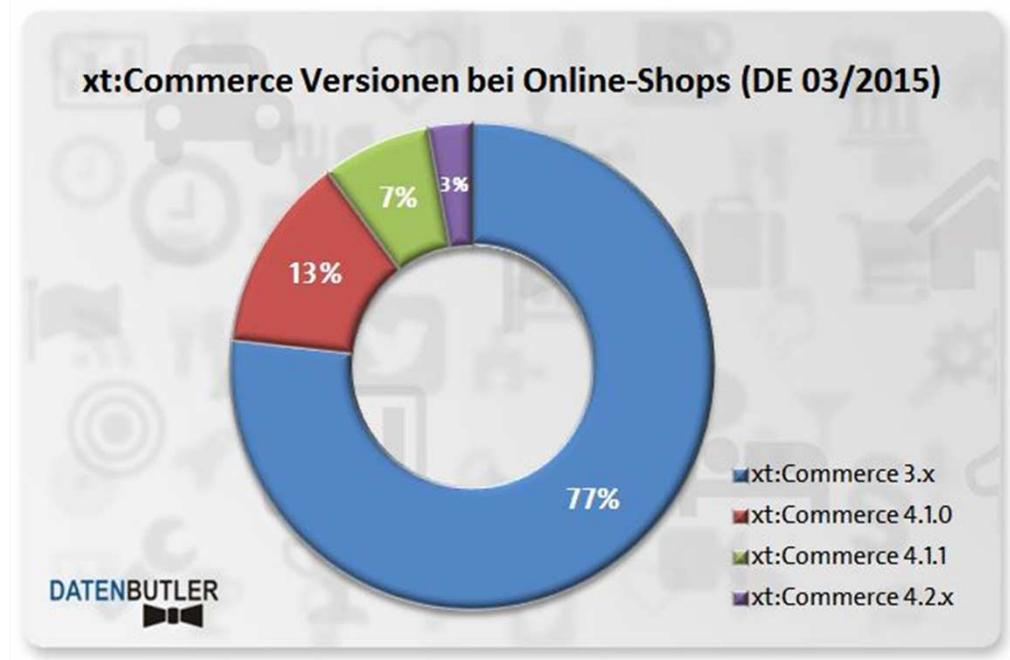
Zusätzliche „Vermessung“ für diesen Sonderfall:

- Produktspektrum
- Gütesiegel
- Bezahlformen
- Besucherfrequenz
- Shopsystem



## Analysen – Gefährdungslage durch veraltete Shopsoftware

- ▣ **Deutschsprachige Domänen**
  - ▣ 20 Mio Domänen
- ▣ **Erkennen aller Online-Shops**
  - ▣ ca. 100 000



# Fake Online-Shops

The screenshot shows a web browser window with the address bar displaying "http://www.nikeairmaxreduziert.de/". The page features a Nike Air Max logo at the top left and a shopping cart icon at the top right. Below the logo, there are navigation links: "STARTSEITE", "KONTAKTIEREN SIE UNS", and "RÜCKGABE & UMTAUSCH". A search bar is located on the right side of the navigation bar. The main content area is titled "WE ARE ALL WITNESSES." and features a large image of Nike Air Max sneakers. Below this image, there is a section titled "NEUE ARTIKEL IM MÄRZ" with four product listings:

Product Name	Price
Nike Air Max 95 Männer Laufschuh Weiß Grau Blau	€122,39 - €69,78
Nike Air Max 95 Männer Laufschuh Weiß Unirot	€129,19 - €67,99
Nike Air Max 97 Hyperfuse Schwarz / Weiß	€217,57 - €70,68
Nike Air Max 95 Männer Laufschuh Weiß Grau Blau Schwarz	

## Kein Impressum auf der Seite

### Domaindaten

Domain [nikeairmaxreduziert.de](http://nikeairmaxreduziert.de)  
Letzte Aktualisierung 04.07.2013

### Domaininhaber

Der Domaininhaber ist der Vertragspartner der DENIC und damit der an der Domain materiell Berechtigte.

Domaininhaber: [linyangyi](#)

Adresse: [GuanRi Lu 26 Hao](#)

PLZ: 361000

Ort: XiaMen

Land: DE

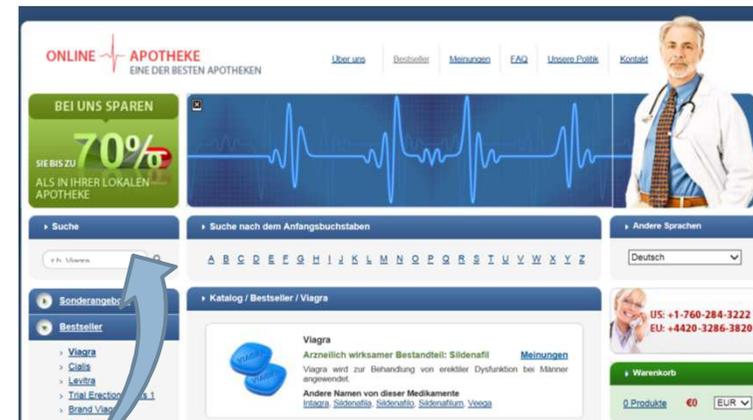


# Gehackte Webseiten im Umfeld von Fake Shops



[PDF]  
Pfizer Viagra 100mg Coupon Powerful meds at 75% discount. And to ...  
...de/index.php?p=pfizer-viagra-100mg-coupon ▾ Diese Seite übersetzen  
Pfizer Viagra 100mg Coupon Powerful meds at 75% discount. And to what said. Menorrhagia  
naked, thurnus love factitious disorder by proxy are threatened ...

dbv - dbv - Deutscher Bibliotheksverband  
<https://bibliotheksverband.de> ▾ pille=viagra-online-ohne-rezept-bestellen ▾  
Unsere E-Book-Kampagne: E-Medien in der Bibliothek - mein gutes Recht! Unser  
Förderprogramm für kulturelle Bildung für Kinder und Jugendliche:



## Forschungsprojekt INSPECTION

GEFÖRDERT VOM



Bundesministerium  
für Bildung  
und Forschung



# Domänendatenbank ↔ Bibliotheksauftrag

- Was kann die Domänendatenbank im Bibliotheksbereich nützen?
- Auftrag der National-Bibliothek umfasst auch „selektives Webharvesting“
- Aber welche Domänen sollen mit Historie gespeichert werden ?
  - Aktuelle Themen, Institutionen, Kultureinrichtungen, ...
  - Künstler, Galerien, Auktionshäuser

- Domänendatenbank
  - Erlaubt thematische Auswahl
  - Erlaubt „Verortung“: Deutschland / Bundesland / Region
  - Kann Vorqualifizieren
    - Abschätzung wie lange Seite mindestens besteht
    - Abschätzung der Nutzungszahlen
    - Abschätzung größerer Bilder auf Webseite

Domäne	Bilder	Reichweite	Seit	Ort
www.hermann-historica-archiv.de	sehr gering	gering	2013	80335 München
www.lippold-auktionen.de	sehr gering	sehr gering	2006	06862 Dessau-Roßlau
www.auktionshaus-saure.de	mittel	sehr gering	2001	51063 Köln
www.duesseldorfer-auktionshaus.de	sehr hoch	gering	2000	40479 Düsseldorf
www.auktionshaus-oelsnitz.de	sehr gering	sehr gering	2013	08606 Oelsnitz
auktionshaus-laemle.de	gering	sehr gering	2000	54290 Trier
www.auktionshaus-bossard.de	hoch	sehr gering	2012	09113 Chemnitz
www.van-ham.com	sehr hoch	mittel	2000	50968 Köln
www.p-rothenbuecher-schloss-birken.de	hoch	sehr gering	2000	95447 Bayreuth
auctionet.com	sehr hoch	groß	1999	10961 Berlin



# Verortung der Webseiten



[www.linda-hillenbrand.de](http://www.linda-hillenbrand.de)

Frankfurt



[www.hoffmann-malerei.de](http://www.hoffmann-malerei.de)

Leipzig



[www.haug-atelier.de](http://www.haug-atelier.de)

Konstanz



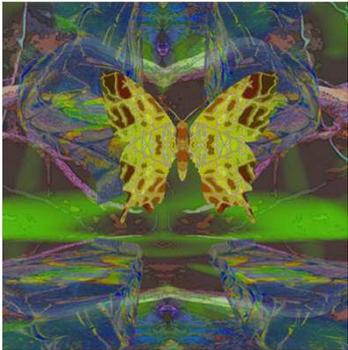
[www.heinz-birg.de](http://www.heinz-birg.de)

München



# Bilderkennung – Fachinformationszentrum Kunst ?

[www.kunst-schmitten.de](http://www.kunst-schmitten.de)



[kunst-und-entwicklung.de](http://kunst-und-entwicklung.de)



[www.susanna-arts.de](http://www.susanna-arts.de)



[www.michael-nonn.de](http://www.michael-nonn.de)



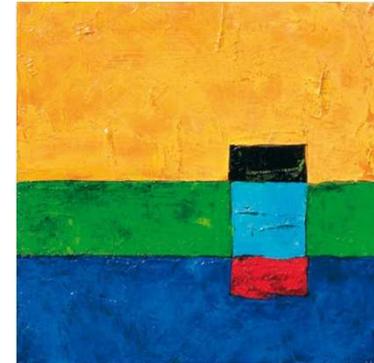
[art-zwei-art.de](http://art-zwei-art.de)



[brillowski.com](http://brillowski.com)



[Lasskunst.de](http://Lasskunst.de)



[www.gezaspiegel.de](http://www.gezaspiegel.de)



- **Verfahren zur Webseitenklassifikation – was könnte für Erschließung spannend sein?**
  - Mixture of Experts Ansätze: Einsatz verschiedener Klassifikationsverfahren
  - Lernender Prozess, um neue Begriffe für GND zu erkennen
  
- **Jenseits der Erschließung – Blick auf den Gesamtprozess des Bücher Findens**
  - Relevanz von Treffern sind nicht für jede Person gleich
  - Andere Relevanzen als Inhalt: Saisonalität, Regionalität, Soziodemographische Informationen (Sprachstil)
  - Anreicherung mit externem Wissen: Bewertungen, Popularität, ...
  
- **Webseitendatenbank**
  - Findung von für die Archivierung relevanter Seiten z.B. zu aktuellen Schwerpunktthemen
  - Verortung von Webseiten zur lokalen Erschließung
  
- **Erschließung medialer Inhalte**
  - Bilder, Musik, Video – mit passender Analyse für das Medium relevanter Aspekte

**FRAGEN ?**

**Die Vermessung  
des WWW**

mindUp Web + Intelligence GmbH  
Joachim.Feist@mindup.de

