

A close-up photograph of a person in a dark suit and tie, pointing their right index finger upwards. The background is blurred. A white rectangular box with a thin border is overlaid on the image, containing the text 'Google für Fortgeschrittene' and 'Search'.

Google für Fortgeschrittene

Search

Was ist Information Retrieval (IR)?

- Ein **Werkzeug** mit dem **Informationen** ausgewählt werden
- **Interaktion** mit Informationen
 - Informationsbedarf des Benutzers muss dem System übermittelt werden
 - Die gefundenen Informationen bzw. eine geeignete Darstellung wird dem Benutzern präsentiert
- Zentrale **Probleme**:
 - Umsetzung des menschlichen Informationsbedarfs in eine für die Maschine verständliche Form
 - Darstellung der maschinengerecht vorliegenden Informationen in eine für Menschen geeignete Form

IR Basiskonzepte

- Wer befasst sich damit?
 - **Informatik:** Technologie, Computerwissenschaften, Datenstrukturen
 - **Linguistik:** Sprache, Konzepte
 - **Informationswissenschaft:** Systeme, Werkzeuge und Strategien des Suchens, Digitale Bibliotheken
- Retrievalsoftware: Volltext-Suchmaschinen, kommerzielle Datenbanksysteme, Suchmaschinen im Web

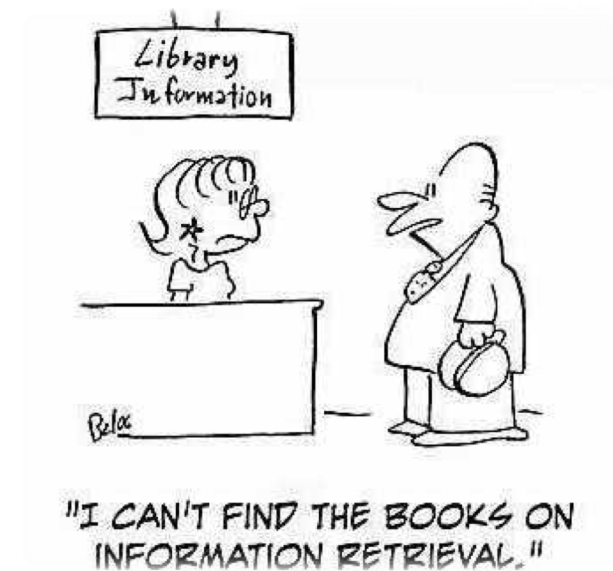
Linguistik und IR

- Sprache
 - Strings
 - Worte
 - Terme
 - Phrasen
- Konzepte
 - Textsammlungen, Kollektionen (Typen, Charakteristiken, statistische Analysen, erkennen von Konkordanzen)
 - Indexierung und Klassifikation
 - Such-
 - Hilfen
 - Prozesse

Geschichte des IR

- Bei der Entwicklung von Informationssystemen spielten folgende Faktoren eine grosse Rolle:
 - Informationsbedarf
 - Technische Machbarkeit
- Computer waren teuer und gross
- Erste Systeme um den Bestand an wissenschaftlicher Literatur zu verwalten
- Bibliothekskontext, Katalogsysteme
 - Klassische Aufgabenstellung
 - Inhaltliche Beschreibungen
- **Information Retrieval**

Websuche ist eine Form von IR



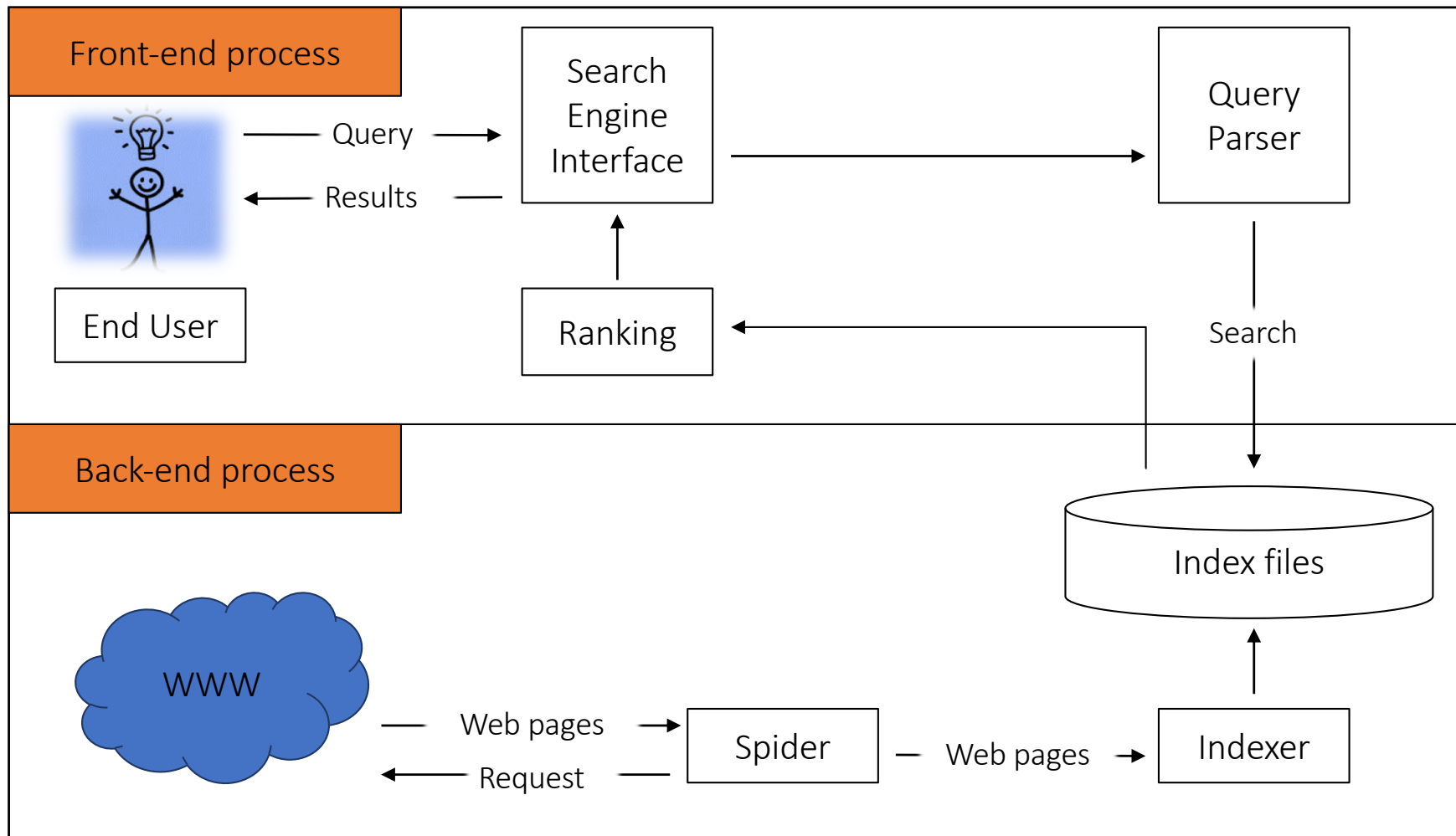
Suche im WWW



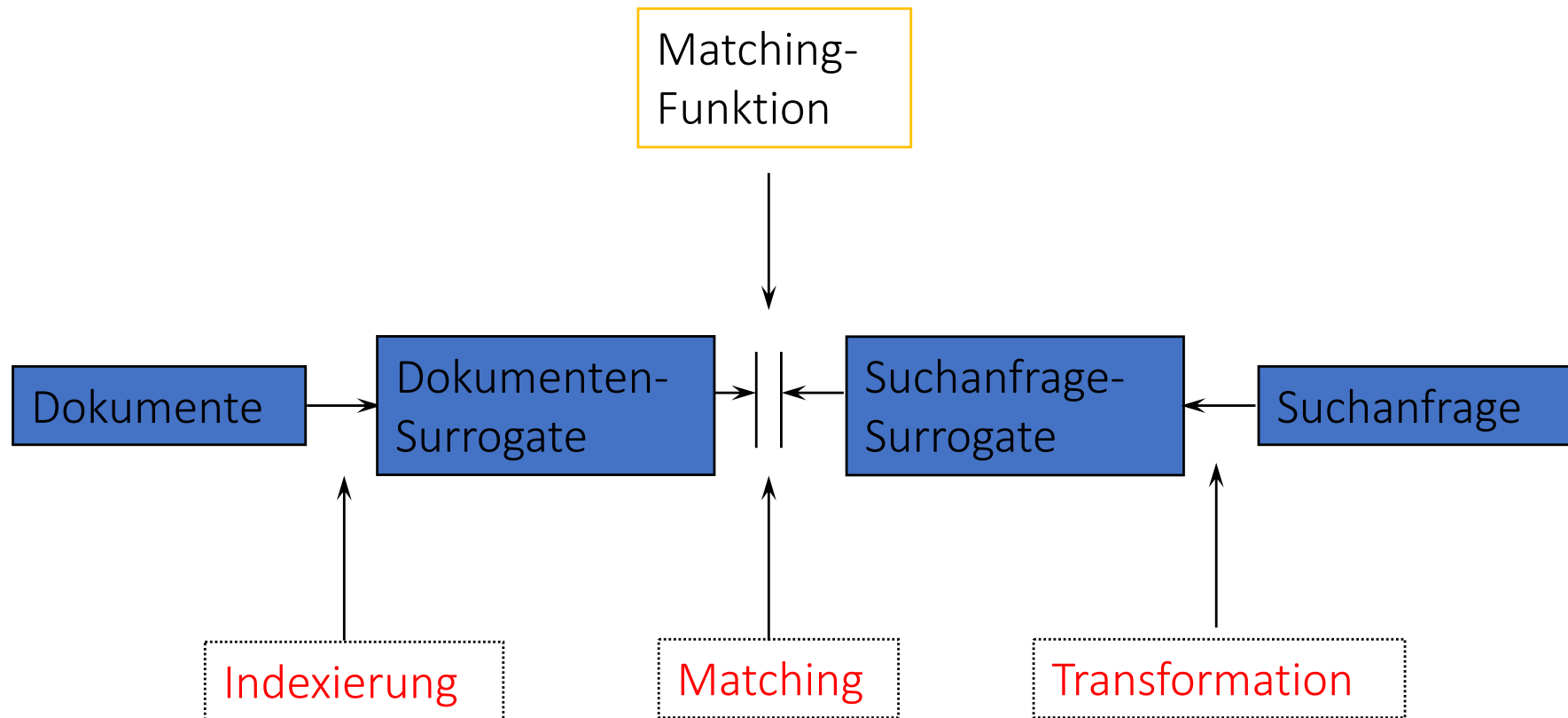
WWW-Suchdienste

- **Suchmaschinen**
 - Zentrale Systemarchitektur
 - Verteilte Systemarchitektur
- **Meta Suchmaschinen**
 - MetaGer (<http://meta.rrzn.uni-hannover.de>)
- **Spezialisierte Suchmaschinen** und Datenbanken
 - Gefilterte Informationen zu speziellen Gebieten
 - Strukturiert durch Portale
 - Informationsflut des Internets wird eingegrenzt
 - Oft breites Zusatzangebot an Materialien und Tools
 - Themenspezifische Suchdienste oder Datenbanken
 - Email-Adressen
 - News

Zentrale Architektur von Suchdiensten



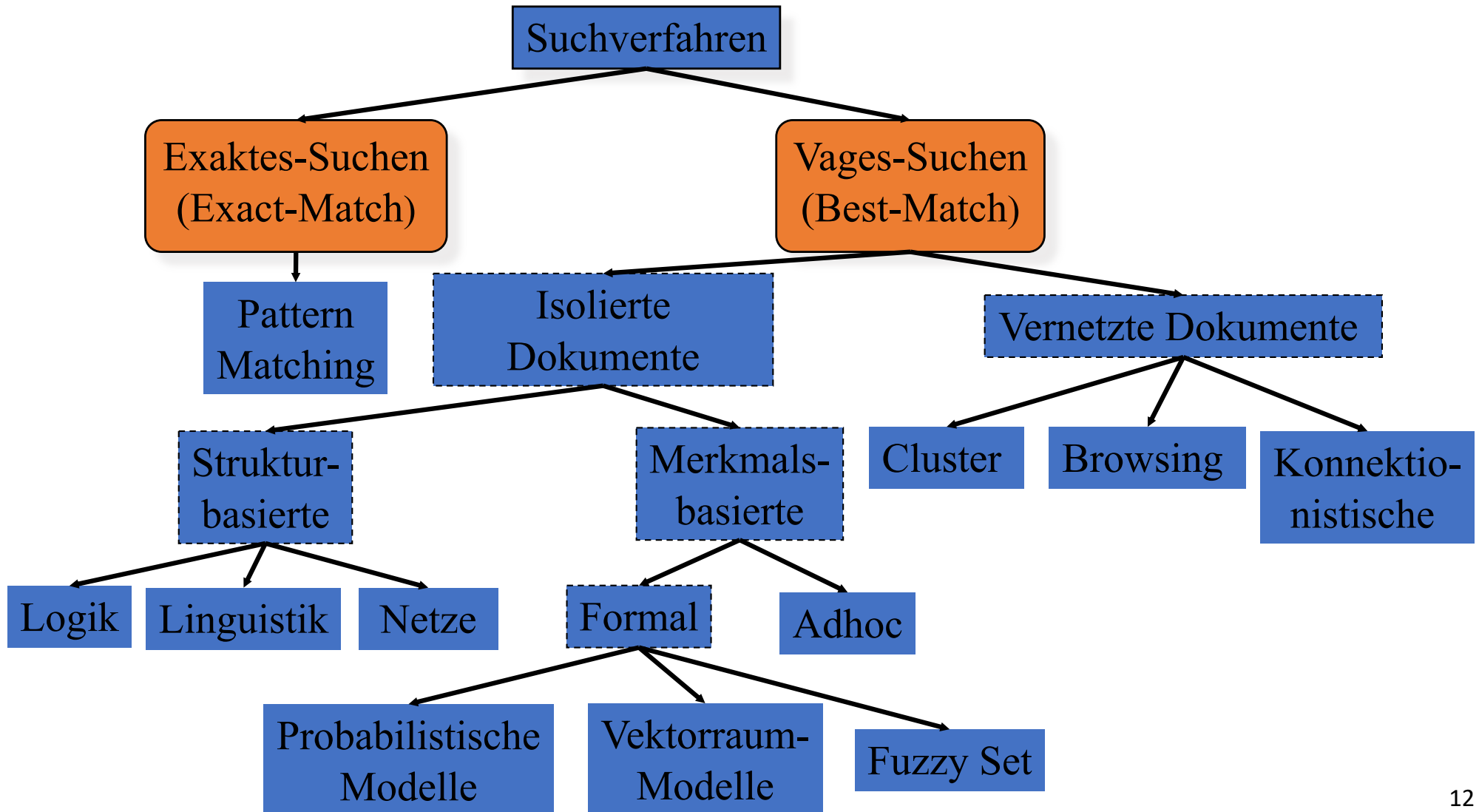
Einfaches Information Retrieval Modell



Arbeitsweise von Suchdiensten

- Strategien von Spidern
 - Parallele Spider
 - Aufteilung der Web nach Domänen, Adressbereichen
 - Startliste von URLs
 - Besuchsintervalle bei der Indexierung
 - Schnelle Indexierung bei ‚bezahlten‘ URL
- Heuristiken
 - Breadth First versus Depth First
 - Volldurchlauf?
- Probleme
 - Last auf Netz und Server
 - Ausschluss der Indexierung durch Robot Exclusion
 - Keine Indexierung dynamisch generierten Seiten, Datenbanken, etc.

Suchverfahren der Maschinen



Hauptaufgaben des Web Information Management

- **Zugang zur Information**
 - Suche (Suchmaschinen, z.B. Google, Altavista)
 - Navigation (Browsers, z.B. IE, Firefox)
 - Filtering (Empfehlungssystem, z.B. Amazon)
- **Organisation der Information**
 - Kategorisierung (Web Directories, z.B. Yahoo!, Open Directory Project, DMOZ)
 - Clustering (Hitliste wird kategorisiert, z.B. Vivisimo)
- **Information Discovery** (Information Mining, z.B. Enterprise Information Portal (EIP) von IBM)

Möglichkeiten und Probleme

- Übersichtsseiten im WWW sind meistens von einzelnen Personen ehrenamtlich zusammengestellt (ähnliche Probleme wie beim persönlichen Nachfragen).
- Suchmaschinen beschränken sich auf die Suche nach einzelnen Wörtern.
- Die meisten Dokumente im WWW sind nicht systematisch strukturiert.
 - Nur HTML-XML-Dokumente, keine animierte Seiten
 - Auf Angebote in Datenbanken und speziellen Informationsseiten können Suchmaschinen nicht zugreifen.
- Im WWW ist es nur bedingt möglich, die Richtigkeit zu überprüfen.

Eigenschaften von Information im WWW

- **“Unendliche” Informationsmenge** (Surface vs. Deep Web)
 - Surface = statische Webseiten (HTML)
 - Deep = dynamisch generierte Webseiten (DB's)
- **Semistrukturierte Dokumente** (kein festes Schema)
 - Strukturierte = HTML tags, hyperlinks, etc.
 - Unstrukturierte = Text
- **Unterschiedliche Formate** (PDF, MS-Word, PS, ...)
- **Multimedia** (Textuelle, Audio, Bilder, ...)
- Grosse Qualitätsunterschiede („Viel Mist“)
- Alles ist publizierbar!
- **Herausforderung** ist: Verteilte Daten, Fluktuation, unterschiedliche Qualität, heterogene Daten

Herausforderungen

- **Verteilte Daten**
 - Auf verschiedenen Plattformen, welche mit unterschiedlichen Bandbreiten vernetzt sind
 - Was ist ein Dokument, was kann indexiert werden?
- **Grosse Fluktuation** der Daten
 - Daten können leicht verändert werden
- **Unterschiedliche Qualität** der Daten
 - Falsche Informationen
 - Veraltet / keine Aktualisierungsautomatik
 - Schlecht geschrieben
 - ‚Unstrukturierte‘ und redundante Daten
- **Heterogene Daten**
 - Format, Multimedia
 - Verschiedene Sprachen, Zeichensätze

How do I do it?

**DENKEN
IST WIE
GOOGELN,
NUR EBEN
VIEL KRASSER.**

Bewusst sein, was effektiv gesucht wird
und in welchem Kontext

Wo und wie könnte die gewünschte
Information aufgeführt sein

Überlegungen vor der Suche können viel
Arbeit ersparen

Stichwort: VORBEREITUNG

Wie Suchen?

- Welche Stichwörter beschreiben das Problem besonders gut?
- Wörter, die **spezifisch für die Fragestellung** sind, aber **so allgemein**, dass sie in jedem **“wichtigen”** Artikel vorkommen.
- Allenfalls erste grob Information über Wikipedia, um anschliessend gezielter zu suchen

- **Beispiel:**
 - Eine Agentur erhält von einem Institut für Rechtsmedizin den Auftrag, eine Informationsbroschüre für Laien zum Thema DNA-Vaterschaftsanalyse zu erstellen. Das Grundlagenmaterial des Kunden für die Arbeit besteht aus einigen wissenschaftlichen Artikeln, den Informationen auf der Website des Instituts und einigen Passagen aus Jahresberichten vergangener Jahre. Aufgabe des zuständigen Redaktors.

Mögliches Vorgehen

1. Verständnis der Thematik fördern – Erst Information

- Beispielsweise über Wikipedia
 - Die verschiedenen Quellenartikel des Wikipedia-Eintrages durchgehen und anschauen
 - Eröffnet neue Ideen zur Weitersuche oder zu effizienten Keywords

Weblinks [Bearbeiten]

 Wikinews: Bundesverfassungsgericht: Heimliche Vaterschaftstests sind unzulässig – Nachricht

- Rechtsmedizin Uni-Mainz [\(Memento vom 18. Mai 2006 im Internet Archive\)](#) Methoden der DNA-Analyse bei Abstammungsbegutachtung und forensischer Spurenkunde (im Internetarchiv)
- Beschlossener Gesetzentwurf des Gesetzes zur Klärung der Vaterschaft unabhängig vom Anfechtungsverfahren [\(PDF; 260 kB\)](#)
- Neue Regeln für Vaterschaftstests Bundestag stimmt Gesetzentwurf zu [\(Memento vom 10. März 2009 im Internet Archive\)](#) von den Seiten des Deutschen Bundestages
- Bundestag verabschiedet Gesetz zur Vaterschaftsfeststellung (22. Februar 2008) [\(Memento vom 4. März 2008 im Internet Archive\)](#)
- Gesetz zur Klärung der Vaterschaft unabhängig vom Anfechtungsverfahren im Bundesgesetzblatt (BGBl. 2008 I S. 441) [\(PDF; 53 kB\)](#)

Einzelnachweise [Bearbeiten]

- ↑ BVerfG, Beschluss vom 18. August 2010 [\(PDF; 34 kB\)](#), Az. 1 BvR 811/09, Volltext.
- ↑ OLG Stuttgart, Beschluss vom 10. August 2009 [\(PDF; 177 kB\)](#), Az. 17 WF 181/09, Volltext; Abs. 5 bis 7.
- ↑ § 17 [\(PDF; 177 kB\)](#) Abs. 3 Nr. 2 Gendiagnostikgesetz.
- ↑ BGH, Urteile vom 12. Januar 2005 [\(PDF; 34 kB\)](#), Az. XII ZR 60/03, Volltext und Az. XII ZR 227/03 [\(PDF; 42 kB\)](#), Volltext.
- ↑ BVerfG, Urteil vom 13. Februar 2007 [\(PDF; 177 kB\)](#), Az. 1 BvR 421/05, Volltext.
- ↑ § 26 [\(PDF; 177 kB\)](#) Abs. 2 Gendiagnostikgesetz.
- ↑ Mitteilung der Gendiagnostik-Kommission (GEKO) vom 10. September 2010 [\(Memento vom 19. September 2010 im Internet Archive\)](#), abgerufen am 10. Mai 2015.
- ↑ Petra Gehring: Heimliche Vaterschaftstests: Blowaffen im Geschlechterkampf [\(PDF; 177 kB\)](#). In: *thema forschung (TU Darmstadt)* 2/2005, S. 31–33.
- ↑ aerztezeitung.de: Männer sind für heimliche Vaterschaftstests [\(PDF; 177 kB\)](#)
- ↑ Dokumentation: Gesetz der Bundesregierung-Entwurf eines Gesetzes zur Klärung der Vaterschaft unabhängig vom Anfechtungsverfahren, in: *Familie, Partnerschaft, Recht (Fachzeitschrift)* 2007, Seite 403
- ↑ Zypries will heimliche Vaterschaftstests verbieten lassen. [\(PDF; 177 kB\)](#) In: *Spiegel Online*. 3. Januar 2005, abgerufen am 21. Januar 2011.
- ↑ Dietmar Hipp: Kuckucksei im Nest. In: *Der Spiegel*. Nr. 7, 2007 [\(online\)](#) [\(PDF; 177 kB\)](#).
- ↑ OLG Stuttgart, Beschluss vom 11. Juli 2008 [\(PDF; 177 kB\)](#), Az. 8 WF 102/08, Volltext.
- ↑ Anlage 2 [\(PDF; 177 kB\)](#) (zu § 10 [\(PDF; 177 kB\)](#) Abs. 1 JVEG).

Mögliches Vorgehen

2. Google Scholar zur Hand nehmen zur Überprüfung von neueren Artikel und der Relevanz der in Wikipedia aufgeführten Artikel
 - Anzahl Zitierungen kann ein Anhaltspunkt sein
 - Suchen sowohl auf Englisch wie auch auf Deutsch durchführen
 - Querlesen der Artikel

The screenshot shows a Google Scholar search interface. The search bar contains the text "vaterschaftstest dna analyse" and a search button. Below the search bar, it indicates "About 216 results (0.04 sec)". The left sidebar contains navigation options: "Articles", "Case law", "My library", "Any time", "Since 2015", "Since 2014", "Since 2011", "Custom range...", "Sort by relevance", "Sort by date", "include patents", "include citations", and "Create alert". The main content area displays search results. The first result is titled "Forensische DNA-Analyse" by RB Dettmeyer, HF Schütz, MA Verhoff, and RB Dettmeyer, published in "Rechtsmedizin" in 2014 by Springer. The abstract discusses the use of multiplex kits for simultaneous analysis of 16 STR systems. The second result is titled "„Wo man mit Blut die Grenze schrieb...“ –Philosophisch-ethische Überlegungen zur Anwendung von DNA-Analysen bei Familienzusammenführungen" by JS Guggenheimer, published in "Migration, Familie und Gesellschaft" in 2014 by Springer. The abstract discusses the ethical implications of DNA testing in family reunification cases. The third result is titled "DNA-auf Spurensuche im Erbgut" by N Podbregar, published in "Genetik" in 2013 by Springer. The abstract discusses the use of DNA testing to identify biological parents. The fourth result is titled "„Richtige“ Kinder: von heimlichen und folgenlosen Vaterschaftstests" by S Schutter, published in 2011 by books.google.com. The abstract discusses the legal and ethical implications of DNA testing in family reunification cases.

Mögliches Vorgehen

3. Mit dem neuen Wissen verfeinerte (Google-) Suchen (nutzen von Operatoren und der Erweiterten Suche) durchführen
 - Bspw. Land einschränken oder ähnliches
4. Herkunft der Infos kritisch hinterfragen und Querchecken
 - Auf was für einer Seite befinde ich mich, habe ich bereits ähnliche Informationen gesammelt , etc.
5. Gesammeltes Sichten und Verarbeiten

Probleme

- Homographen
 - Gleich klingendes Wort
 - Lehre/Leere
- Polyseme
 - Worte mit mehreren Bedeutungen (historisch bedingte Gleichheit für dieselbe Bedeutung)
 - Bank: Sitzgelegenheit oder Geldinstitut
- Flexionsformen
- Derivationsformen
- Komposita
- Vagheit – Unschärfe – Unsicherheit
 - Sowohl in der Anfrage wie auch in der gegebenen Information

Kompositazerlegung/Motivation

- Im Deutschen treten Komposita sehr häufig auf; werden “spontan” gebildet.
- Vollständige Erfassung sämtlicher Komposita quasi nicht möglich.
- **aber:** Rückführung auf Grundformen möglich, endliche Menge von Regeln zur Kompositazerlegung
 - z.B. Einfügung eines -s: Recht + Streit = Recht**s**streit
 - z.B. Einfügung eines -n: Bauer + Kind = Bauern**n**kind
- Dies ermöglicht:
 - Das Auffinden von Dokumenten, die nur einen Teil eines Kompositas enthalten, das in der Anfrage spezifiziert wurde.
 - Das Auffinden von Dokumenten, die Komposita enthalten, in welchen ein Term der Suchanfrage vorkommt.
- Verarbeitungsverfahren analog zur morphologischen Verarbeitung

Grundsätzliche Fragestellungen

- Welche Möglichkeiten zur Spezifikation einer Anfrage gibt es?
- Wie werden diese Anfragen bearbeitet?
- Welche standardisierten "Sprachen" (Protokolle) zur Abfrage von Datenbeständen (Online-Datenbanken, Suchmaschinen) gibt es?

Anfragesprache

- Informationssuche ist die Suche nach Information mittels einer Sprache. Diese Sprache wird als Anfragesprache (Retrievalsprache) bezeichnet.
 - Formale Sprachen
 - Abstrakte Sprache wie z.B. Programmiersprachen
 - Anfragesprachen können wie folgt eingeteilt werden:
 - nach Einsatz und Benutzungsgrad
 - nach Dokumententyp
- Beispiele:
 - Suchmaschinen
 - Anfragen mit Stichworten
 - Faktenretrieval SQL (Structured Query Language)
 - Anfragen auf Tabellen, Datenbanken
 - Anfragen auf Dokumentenstrukturen XML (XQuery)

Anfragen mit Stichworten (1)

- Grundlegende Annahme:
 - Die kleinste Einheit in einem Dokument sind Worte.
- Einfachste Anfrage: ein Wort
 - Je nach System eine Folge von Buchstaben, Zahlen, Sonderzeichen
- Je nach IR-Model muss mindestens eines der Anfrageworte im Dokument enthalten sein.
- Ranking:
 - Gewichtung mittels *term frequency* und *inverse document frequency*.

Anfragen mit Stichworten (2)

- Boolesche Anfragen
 - Aussagenlogische Formeln
 - Operatoren: AND, OR, NOT (oftmals nur AND NOT)
 - Kombination mit regulären Ausdrücken
 - Sing* AND (Beatles OR "Pink Floyd")
- Natürlichsprachliche Anfragen
 - Keine Beschränkung bzgl. der Anfragenformulierung
 - Schwierigkeiten bei der Verarbeitung der natürlichen Sprache

Anfrageverfeinerung durch Operatoren (1/2)

Symbol	Verwendung
+	Suche nach Google+ Seiten oder Blutgruppen Beispiele: <code>+Chrome</code> oder <code>AB+</code>
@	Suche nach sozialen Tags Beispiel: <code>@googler</code>
\$	Suche nach Preisen Beispiel: <code>nikon 400 \$</code>
#	Suche nach beliebten Hashtags zu Trendthemen Beispiel: <code>#throwbackthursday</code>
-	Wenn Sie einen Bindestrich vor einem Wort oder einer Website platzieren, werden Ergebnisse entfernt, die dieses Wort oder diese Website enthalten. Dies kann bei Wörtern mit verschiedenen Bedeutungen nützlich sein, etwa bei "Jaguar", bei dem Ergebnisse zur Automarke und zum Tier ausgegeben werden. Beispiele: <code>jaguar geschwindigkeit -auto</code> oder <code>panda -site:wikipedia.org</code>

- > Mit dem Operator ~ werden ebenfalls nach Synonymen des nachfolgenden Wortes gesucht
- > Mit | können Begriffe separiert werden (entweder das eine oder das andere) Suchmenge wird grösser.

Anfrageverfeinerung durch Operatoren (2/2)

"	<p>Wenn Sie ein Wort oder eine Wortgruppe in Anführungszeichen setzen, enthalten die Ergebnisse nur solche Seiten, auf denen die Wörter in der gleichen Form und in der gleichen Reihenfolge wie innerhalb der Anführungszeichen vorkommen. Verwenden Sie Anführungszeichen nur dann, wenn Sie nach einem ganz speziellen Wort oder einer exakten Wortgruppe suchen. Ansonsten könnten Sie viele hilfreiche Ergebnisse versehentlich ausschließen.</p> <p>Beispiel: <code>"imagine all the people"</code></p>
*	<p>Fügen Sie ein Sternchen als Platzhalter für alle unbekannt Begriffe hinzu.</p> <p>Beispiel: <code>wer den * nicht ehrt, ist des * nicht wert</code></p>
..	<p>Trennen Sie Zahlen durch zwei Punkte ohne Leerzeichen, um Ergebnisse mit einem bestimmten Zahlenbereich zu erhalten.</p> <p>Beispiel: <code>kamera 50..100 \$</code></p>

- > Mit **filetyp:** pdf oder csv usw. kann die Suche auf bestimmte Dokumenttypen eingeschränkt werden.

Weitere hilfreiche Möglichkeiten

Suchoperator	Verwendung
site:	Gibt Ergebnisse für bestimmte Websites oder Domains aus Beispiele: <code>olympiade site:zeit.de</code> und <code>olympiade site:.de</code>
link:	Mit diesem Operator können Sie nach Webseiten suchen, die auf eine bestimmte andere Seite verweisen. Beispiel: <code>link:youtube.com</code>
related:	Mit diesem Operator finden Sie Websites, die einer bestimmten Webadresse ähneln. Beispiel: <code>related:spiegel.de</code>
OR	Mit diesem Operator finden Sie Webseiten, die irgendeinen von mehreren Begriffen enthalten. Beispiel: <code>marathon OR lauf</code>
info:	Dieser Suchoperator liefert Informationen zu einer Webadresse. Dazu zählen unter anderem die im Cache gespeicherte Version der Seite, ähnliche Seiten und Seiten, die auf die angegebene Website verweisen. Beispiel: <code>info:google.de</code>
cache:	Mit diesem Suchoperator können Sie die Seitenversion abrufen, die Google beim letzten Besuch der angegebene Website zwischengespeichert hat. Beispiel: <code>cache:bundesregierung.de</code>

Anfragen auf Dokumentstrukturen

- Integration von inhaltlichen und strukturellen Aspekten in Anfragen
- Auswertung: inhaltlicher Anteil der Anfrage liefert eine Menge von potentiellen Ergebnissen, auf der dann die strukturellen Constraints überprüft werden.
- Unterschiedliche Arten von Text-Struktur erlauben unterschiedliche Anfragen.
 - Festgelegte Struktur
 - Hypertext-Struktur
 - Hierarchische Struktur

Hypertext-Struktur / Beispiele

- "Finde Informationen über Bern, eingeschränkt auf Domäne Schweiz."
 - *SELECT d.url, d.title
FROM Dokument d SUCH THAT d MENTIONS "Bern"
WHERE d.url CONTAINS ".ch";*
- "Finde Dokumente über Informatik, die auf die Universität Bern verweisen."
 - *SELECT d.url, d.title
FROM Dokument d SUCH THAT d MENTIONS "Informatik",
Verweis v SUCH THAT basis = d
WHERE v.anchor CONTAINS "unibe.ch"*
 - **MENTIONS**: Wird durch Suchanfrage an eine globale Suchmaschine (z.B. Google) realisiert; Zwischenergebnis wird weiterverarbeitet
 - **CONTAINS**: Lokale Zeichenkettensuche innerhalb eines Strings

Operationen auf Anfragen

- Initiale Spezifikation des Informationsbedürfnisses oft unzureichend, dementsprechend schlechte Qualität der Ergebnisse nach dem ersten Retrievalschritt
- Reformulierung notwendig, aber durch den Benutzer oft schwierig, da Kenntnisse über Dokumentensammlung kaum vorhanden
- Ziel: automatische Reformulierung
- Basis zur automatischen Reformulierung:
 - Rückmeldungen vom Benutzer über die Relevanz einzelner Dokumente aus dem Ergebnis (Relevance Feedback)
 - Analyse der initial gefundenen Dokumente im Ergebnis (lokal)
 - Analyse des gesamten Datenbestandes (global)

Erweiterung/Modifikation der Anfrage der Termgewichte im VSM

- (Cluster-)Annahmen
 - Gewichtsvektoren der als relevant markierten Dokumente sind einander ähnlich
 - Gewichtsvektoren von irrelevanten Dokumenten sind den obigen Vektoren nicht ähnlich
- Idee: Anfragevektor derart modifizieren, dass er "stärker" in Richtung der relevanten Dokumente zeigt

Stemming und Lemmatisierung

- Motivation:
- Simultane Suche nach allen morphologischen oder orthographischen Varianten.
- Verbesserung des Recalls ohne Verschlechterung der Precision.
- "Einfache Lösung": Benutzer muss durch Verwendung von Trunkationsoperatoren (*) oder Disjunktion über alle Formen selbst dafür Sorge tragen.
- Problem bei Trunkierung: Es werden ungewollte Fortsetzungen erzeugt: auto* findet Auto und Autos, aber auch automatisch, Autor oder Automorphismus usw.

Lemmatisierung

- Reduktion der Wortformen auf ihre Grundform (und weitere Information) durch Nachschlagen in einem elektronischen Wörterbuch.
- Vollformenlexikon: Jede Wortform kann direkt im Lexikon nachgeschlagen werden.
- Grundformenlexikon
 - Wortform wird durch morphologische Regeln auf eine potentielle Grundform reduziert, die dann im Lexikon nachgeschlagen wird.
- Vollformenlexikon
 - Ist aufwendiger hinsichtlich Speicherplatz aber effizienter bei der Verarbeitung

Nachteile von Lemmatisierung

- Erfordert umfangreiches elektronisches Wörterbuch, aufwendig in der Erstellung und Wartung.
- Relativ hohe Anforderungen an Verarbeitungszeit oder Speicherplatz.
- Was passiert mit Wortformen, die nicht im Lexikon gefunden werden?
 - Eigennamen
 - Komposita im Deutschen

Probleme bei der Kompositazerlegung im Deutschen

- Verschiedene korrekte Zerlegungen
 - Wachstube in Wachs + Tube oder wach + Stube
- Simplizia können irrtümlich zerlegt werden
 - Pomade -> po + made
 - Proletarier -> prolet + arier
 - Tangente -> tang + ente
- Grosse Anzahl von Zerlegungsambiguitäten, die nicht korrekt sind:
 - Aluminiumherstellung kann auf 12 versch. Arten zerlegt werden, z.B. alu+mini+umher+stellung
 - Alleinerziehende -> all+ein+erzieh+ende

Suche im WWW



Alternative Suchmaschinen

- **Andere grosse Suchmaschinen** (*Bing/Ask/AOL/Blekkoo*)
- **Duckduckgo** (*!google dem Begriff voransetzen, wenn trotzdem über Google gesucht werden soll, die Suchanfrage bleibt aber anonymisiert*)
- **Wikipedia**
- **Baidu** (*für asiatische Inhalte*)
- **Yandex**
- **Qwant** (*Personensuche in Sozialen Medien*)
- **Swisscows** (*für lokales*)
- **Wolfram Alpha** (*äusserst stark im wissenschaftlichen Bereichen*)
- **CC Search** (*für Inhalte unter der Creative Common Lizenz*)
- **Topsy** (*für das Auffinden von Tweets – durchsucht Twitter*)
- **Metasuchmaschinen benutzen** (*je nach Art, kann thematisch gruppierter gesucht werden*)

Limitationen heutiger Suchmaschinen (1)

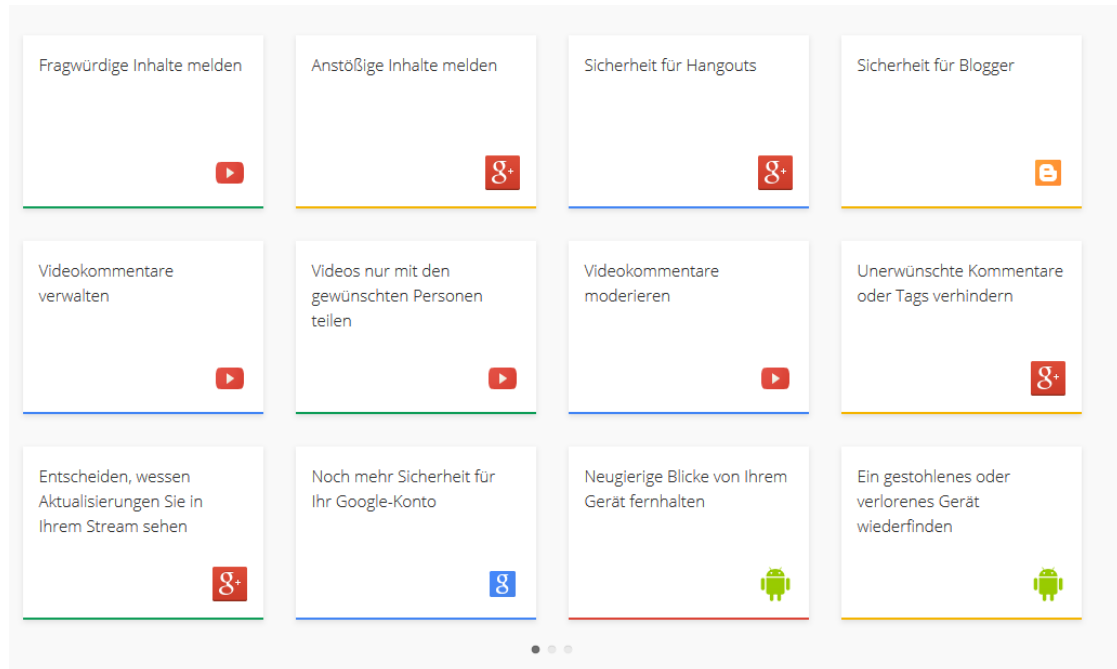
- Suchanfragesprache:
 - Syntax, Lesbarkeit, Verständlichkeit
- Eingeschränktes Verständnis des Dokumenteninhalts (Semantik)
- Ein Sack von Abfrageworte (Zusammenhangslos, Disambiguation) muss verarbeitet werden
- Statistische Gewichtung beim Matching von Suchanfrage zu Dokument
 - Keine Garantie für gute Suchresultate
 - Keine Modellierung des Benutzerkontextes
 - Unterschiedliche Resultate für die gleiche Suchanfrage
 - Ein Benutzer erhält unterschiedliche Resultate zu unterschiedlichen Zeiten und Orten

Limitationen heutiger Suchmaschinen (2)

- Wenig Support für den Suchenden
 - Keine eigentlichen Interaktionen möglich
 - Passive Unterstützung der Suche, keine Empfehlungen
 - Statische Navigation, keine dynamisch generierte Links
 - Schlechte Integration des Retrievals, Empfehlungen und Navigation

Relevanz und Glaubwürdigkeit

Quelle: https://www.google.ch/intl/de_ch/safetycenter/tools/



- Häufig durch User-Community unterstützt, keine Garantie
- Google pflegt einen Relevanz-Algorithmus, der gewährleisten soll das sucherfreundlicher Inhalt höher gewertet wird
- Algorithmus Updates (wie Google Panda, Penguin, Humingbird etc.) bemühen sich darum Qualität des Inhalts zu werten

Sicherheitsangebot von Google

Sicherheitstools von Google

Unsere Produkte bieten Ihnen Tools und Einstellungsmöglichkeiten, mit denen Sie die Internetnutzung Ihrer Kinder steuern können. Die folgenden Abschnitte enthalten Informationen zu SafeSearch, zum sicheren Modus bei YouTube und zum Filtern von Inhalten in Android.

Google SafeSearch

Mit SafeSearch werden Websites gefiltert, die sexuell eindeutige Inhalte enthalten, und aus Ihren Suchergebnissen entfernt. Kein Filter ist absolut sicher, SafeSearch hilft Ihnen jedoch dabei, Inhalte auszublenden, die Sie lieber nicht sehen möchten oder über die Ihre Kinder nicht stolpern sollten.

Standardmäßig ist der moderate SafeSearch-Filter aktiviert, mit dem unangemessene Bilder aus Ihren Suchergebnissen verbannt werden. Wenn Sie Ihre Einstellung ändern und stattdessen die strikte Filterung verwenden, werden sowohl unangemessene Textinhalte als auch Bilder herausgefiltert.

[So ändern Sie die SafeSearch-Einstellungen auf Ihrem Computer](#)

SafeSearch-Sperre

Wenn Sie befürchten, dass andere die strikte SafeSearch-Filterung ohne Ihr Wissen ändern, können Sie die Einstellung über die SafeSearch-Sperre mit einem Passwort schützen. Sobald die Sperre aktiviert ist, ändert sich das Aussehen der Google-Suchergebnisseite, damit der Nutzer weiß, dass SafeSearch aktiviert ist.



Quelle: <http://www.google.ch/intl/de/goodtoknow/familysafety/tools/>

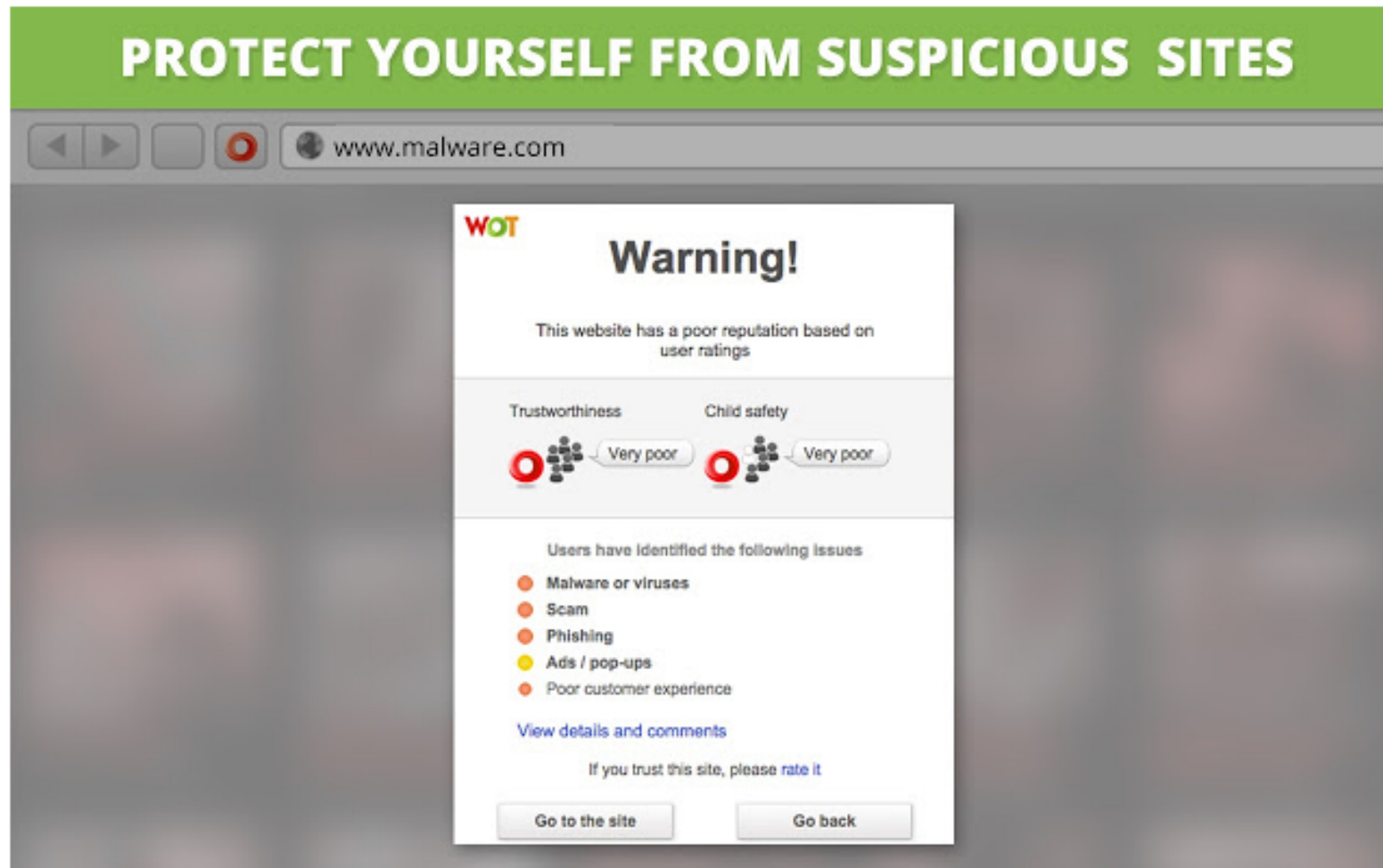
User hat Möglichkeiten sich zu schützen

- Tools für Firefox
 - Https Everywhere
 - WoT (Web of Trust)
 - NoScript
 - BetterPrivacy
 - RequestPolicy



- Ähnliche Tools für andere Browser vorhanden
- Web of Trust auch für Chrome erhältlich

Web of Trust



Information: Ein soziales Phänomen

- Fundamentale **Verschiebung**: progressive Verbreitung von sozialen Aktivitäten,
 - Konversationen, Labeling, Rating und Informationsorganisation in Online Netzwerken
- **Warum** gehen Benutzer Online?
 - Um zu kommunizieren
 - Um informiert zu sein
 - Um sich zu unterhalten
 - Um ein teil eines grossen sozialen Netzes zu sein

Zukunft und Herausforderungen

- Soziale Medien beinhalten **Strukturen** wie Benutzer miteinander interagieren
- Viele Daten warten nur auf ihre Nutzbarmachung
 - “Mining” und “Modellierung”
- Sehr heterogene Daten
- Enormer Umfang von Daten
- Lernen aus Daten und Modelle
 - Lösung von schwierigen Problemen oder für neue Anwendungen
- Websuche ist eine relativ junge wissenschaftliche Disziplin

Trends in „Next Generation Search Engines“

- Support für **Reformulierungen** von Suchanfragen
 - „Query by examples“
 - Automatische Generierung von Suchanfragen (Empfehlungen)
- Bessere **Modellierung des Kontextes**
- Bessere **Personalisierung**
- Bessere Dokumentenanalysen („sentiment analysis“, **Semantik**, etc.)
- **Kontext-sensitives Ranking**

Trends in „Next Generation Search Engines“

- **Spezialisierter / Angepasst an Benutzerbedürfnisse**
 - Spezielle Gruppierungen (Community Suche)
 - *Beispiel:* CiteSeer (CiteSeerx) Scientific Literature Digital Library (Digitale Bibliothek wissenschaftlicher Literatur) ist eine Suchmaschine und Zitationsdatenbank, schwerpunktmässig Informatik
 - Personalisierte Suche
 - Suche innerhalb einer Domäne
- Integration von Suche, **Navigation**, Filtering
- Verbesserung der **Präsentation** der Suchresultate
 - Generierung von Summaries
- **Interaktive** Suche

A photograph of a rolled-up piece of brown paper on a brown background. A white, torn paper strip is placed horizontally across the middle, containing the text "Besten Dank!".

Besten Dank!