

Wie können Websites effektiv gefiltert werden?

Autor(en): **Kester, Harold**

Objektyp: **Article**

Zeitschrift: **Comtec : Informations- und Telekommunikationstechnologie = information and telecommunication technology**

Band (Jahr): **79 (2001)**

Heft 4

PDF erstellt am: **06.08.2024**

Persistenter Link: <https://doi.org/10.5169/seals-876534>

Nutzungsbedingungen

Die ETH-Bibliothek ist Anbieterin der digitalisierten Zeitschriften. Sie besitzt keine Urheberrechte an den Inhalten der Zeitschriften. Die Rechte liegen in der Regel bei den Herausgebern.

Die auf der Plattform e-periodica veröffentlichten Dokumente stehen für nicht-kommerzielle Zwecke in Lehre und Forschung sowie für die private Nutzung frei zur Verfügung. Einzelne Dateien oder Ausdrucke aus diesem Angebot können zusammen mit diesen Nutzungsbedingungen und den korrekten Herkunftsbezeichnungen weitergegeben werden.

Das Veröffentlichen von Bildern in Print- und Online-Publikationen ist nur mit vorheriger Genehmigung der Rechteinhaber erlaubt. Die systematische Speicherung von Teilen des elektronischen Angebots auf anderen Servern bedarf ebenfalls des schriftlichen Einverständnisses der Rechteinhaber.

Haftungsausschluss

Alle Angaben erfolgen ohne Gewähr für Vollständigkeit oder Richtigkeit. Es wird keine Haftung übernommen für Schäden durch die Verwendung von Informationen aus diesem Online-Angebot oder durch das Fehlen von Informationen. Dies gilt auch für Inhalte Dritter, die über dieses Angebot zugänglich sind.

Wie können Websites effektiv gefiltert werden?

Wie gross inzwischen die Bedeutung des Internets für Unternehmen ist, steht angesichts von Vorteilen wie erhöhte Produktivität, verbesserte Kommunikationsmöglichkeiten und Bereitstellung umfangreicher Informationen ausser Frage. Das täuscht jedoch nicht über die Gefahr des Missbrauchs hinweg.

Dank Cyberslacking, dem Zeitvertreib mit Moorhuhnjagd und Shopping während der Arbeitszeit, erleiden nicht nur Unternehmen Verluste in Millionenhöhe, sondern Mitarbeiter können sich auch belästigt

HAROLD KESTER

fühlen, wenn Kollegen pornografische, rassistische oder ähnliche Sites aufrufen. Rechtliche Folgen für das Unternehmen sind nicht ausgeschlossen. Um die Nachteile des Internets zu beheben, müssen individuelle Richtlinien für die Internetnutzung innerhalb von Unternehmen aufgestellt werden. Software-Lösungen zur Umsetzung dieser Richtlinien sind allgemein als Internet-Filter bekannt. Doch Internet Filtering kann seine volle Stärke nur dann ausspielen, wenn die Einteilung der Websites, die im Rahmen der unternehmensweiten Internetregelung gefiltert werden sollen, in bestimmte Kategorien genau vorgenommen wurde. Nur dann lassen sich Auswertungen, wie wer ruft welche Website zu welcher Tageszeit ab oder inwiefern stimmt diese Aktion mit den Richtlinien überein, genau durchführen.

Datenbank-Filtering versus dynamisches Filtering

Für die akurate Bestimmung des Inhalts einer Website gibt es zwei Möglichkeiten. Das dynamische oder Runtime Filtering analysiert den Inhalt einer Website, während sie aufgerufen wird, und legt die Kategorie in Echtzeit fest. Die zweite Methode vergleicht die Adresse der aufgerufenen Website mit einer vorher festgelegten Kontrollliste oder Datenbank. Diese Datenbank- oder Kontrolllisten-Lösung hat eine theoretische Genauigkeit von 100%. Praktisch besteht aber die Gefahr, dass die aufgerufene Website

nicht in der Datenbank verzeichnet ist. Die Datenbank bewertet die Site als negativ und nicht zu filtern, obwohl sie vielleicht nach menschlichem Ermessen geblockt werden sollte. Diese «falschen Negative» sind die Folge des so genannten Underblocking. Die Herausforderung, die sich Unternehmen bei der Entwicklung dieser Lösung stellt, ist das regelmässige und genaue Update der Website-Datenbank.

Dagegen ist der dynamische Filter mit seiner Inhaltsüberprüfung in Echtzeit per definitionem immer aktuell. Ob eine Website seit zehn Monaten oder zehn Minuten existiert, spielt durch die Festlegung der Kategorien genau zum Zeitpunkt der Anfrage keine Rolle. Allerdings sind die dynamischen Filter nicht so akkurat, wie sie sein sollten. Sie sind anfällig für Overblocking und sperren teilweise Sites, die nicht gesperrt werden sollten («falsche Positive»). Kein dynamischer Filter ist heute in der Lage, zwischen einer Site zu unterscheiden, welche beispielsweise die Verwendung von Drogen befürwortet und einer solchen, die soziologische Analysen des Drogenmissbrauchs beinhaltet. Pornografische Sites sind dagegen alles andere als schwer zu entdecken und können daher von dynamischen Filtern im Allgemeinen gut aufgespürt werden. Für andere Kategorien ist der Runtime Filter einfach nicht präzise genug.

Bei der Entwicklung einer Internet-Filtering-Lösung muss der Systementwickler die gegensätzlichen Eigenschaften des dynamischen und Datenbank-Filters abwägen und je nach Priorität optimieren. Denn keiner der beiden Ansätze ist perfekt. Zu strenges und zu lockeres Filtering, Underblocking und Overblocking, bilden eine Metrik aus Recall und Precision. Precision ist der Prozentsatz an gefundenen Datensätzen, die für eine Kategorie relevant sind; Recall ist der Prozentsatz relevanter Dokumente, die

gefunden werden. Die Beziehung von Recall und Precision, den Masseinheiten der Genauigkeit, ist invers: Je grösser der Recall, desto kleiner ist die Precision und umgekehrt. Der Systementwickler muss also für bestimmte Applikationen Prioritäten setzen. Er muss sich entscheiden, ob er mehr Sites finden will, oder ob eine präzisere Kategorisierung für seine Zwecke ausschlaggebend ist, das heisst, dass er zwangsläufig bei einer Metrik Kompromisse eingehen muss.

Anpassung an die unterschiedlichen Märkte

Richtlinien zur Internetnutzung sind für Business- und öffentliche Einrichtungen, Schulen und Haushalte gleichermaßen wichtig. Doch jeder Bereich hat unterschiedliche Anforderungen an das Filtering. Der Unternehmensmarkt schliesst Business, öffentliche Einrichtungen (ausser Schulen) und Non-Profit-Organisationen mit ein. Internetnutzung in Unternehmen besteht aus dem Zugang zu geschäftsrelevanten Informationen, dem Senden und Empfangen von E-Mails und der Koordination von Geschäften im Internet. Unternehmensweite Richtlinien zur Internetnutzung zielen nicht nur auf die Reduzierung von Cyberslacking, sondern auch auf den Schutz von Netzwerk-Ressourcen und die Erschaffung einer kollegialen Arbeitsumgebung. Grundsätzlich erfordert der Unternehmensmarkt mehr Precision als Recall. Der Markt der Schulen und Bibliotheken schliesst öffentliche und private Bildungsinstitute aller Art mit ein, wobei es hierbei um die Informationsbeschaffung zu Lernzwecken und Recherchen geht. Während es hier um den Schutz der Studenten vor anstössigem Material und Junk-Mails geht, liegt der Fokus beim Internet Filtering in privaten Haushalten auf dem Schutz der Kinder. Die Anforderungen an den Internet-Filter sind bei Unternehmen, Schulen und Consumers insofern gleich, da sie alle Genauigkeit des Filterprozesses erwarten und ein Blocking arbeitsrelevanter Sites unerwünscht ist. Overblocking ist in den Bereichen Consumer und Schulen akzeptabel, da der Schutz der Kinder vor an-

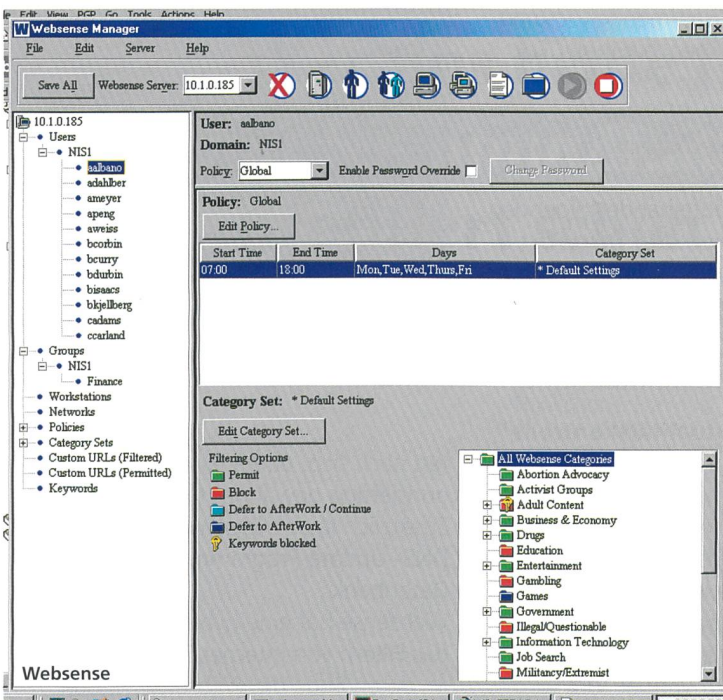
Tools entwickelt, um einzelne virtuell gehostete Sites zu identifizieren und sicherzustellen, dass sie angemessen kategorisiert sind.

Funktionsweise der Master Database

Aufbau und Pflege der Websense Master Database besteht aus mehreren Schritten. An erster Stelle steht das Mining nach den zu sperrenden Sites. Von Websense entwickelte Mining Tools ahmen das Suchverhalten der User nach. Das Ergebnis dieser unterschiedlichen Mining-Prozesse ist eine Liste von Site-Adressen, die dann mit der Hilfe von KILO II (Knowledge Indexing Learning Organization) kategorisiert werden. Zur Kontrolle der von KILO II nicht klassifizierten Sites werden vorsortierte Listen von Web-Analysten erstellt und mit Hilfe des Workbench überprüft.

Hinsichtlich der Mining-Quellen und Tools nutzt Websense unterschiedliche Methoden. Zum einen erwirbt Websense von Unternehmen, die Marktdaten von Endanwendern sammeln, Listen von Domain-Namen der beliebtesten Sites und stützt sich zudem auf Suchmaschinen. Durch die Verwendung auffälliger Ausdrücke, die durch die Analyse von Sites bereits existierender Kategorien gewonnen wurden, fragen Meta-Suchmaschinen gleichzeitig mehrere Suchmaschinen nach einem bestimmten Begriff ab und listen die Ergebnisse nach Relevanz geordnet auf. Bestimmte Websites, so genannte Hubs, sind Directories oder Link-Listen zu anderen Sites. Hubwatcher von Websense ist ein Tool, das diese Adressensammlungen regelmässig durchgeht, um neu gelistete Sites auszumachen und zu überprüfen. Ausserdem können Unternehmen via E-Mail bestimmte Sites für die Aufnahme in die Master Database vorschlagen.

Angefangen bei einer so genannten Seed oder Ausgangs-Webpage, werden mit der Intelligent-Crawling-Technik alle darauf aufgelisteten Links aufgerufen, analysiert und gegebenenfalls als zu sperrende Sites in eine Kategorie aufgenommen. Um dabei effizient zu sein, werden Logik und Algorithmen eingesetzt, damit die richtigen Links einer Prüfung unterzogen werden können. Mit dem Partner Site Collection Program entwickelt



Mit Websense Enterprise 4.2 kann der Zugang zu unerwünschten Websites über 67 Standardkategorien geregelt werden.

stössigem Material Priorität genießt. Dagegen ist Overblocking im Enterprise-Markt inakzeptabel, denn es reduziert die Produktivität und verhindert den Zugang zu legitimen Sites. Auch unterschiedliche Kulturen haben differenzielle Anforderungen an das Filtering. Rassenhass ist zum Beispiel in Deutschland illegal, wodurch eine Maximierung des Recalls bei der Sperrung von Sites in der Kategorie Rassismus erforderlich ist. Manche Kulturen sind restriktiver als andere, was die Kategorie der Erwachsenenunterhaltung betrifft, so werden beispielsweise Subkategorien wie «Badeanzug», «Damenunterwäsche» und «Aufklärungsunterricht» gesperrt. Daher sollte eine Datenbank möglichst umfassend sein. Dem Netzwerkadministrator stehen mit der Lösung von Websense beispielsweise mehr als 65 Kategorien zur Verfügung. Damit lassen sich Sites flexibel freigeben, sperren oder zeitlich eingrenzen. Zusätzlich ermöglicht Websense Enterprise den Administratoren, auf die jeweilige Nutzung abgestimmte Kategorien zu erstellen und sie um handerlesene URLs zu erweitern, damit die Filtering-Ansprüche beispielsweise auch von einzelnen User-Gruppen oder Abteilungen erfüllt werden können.

Websense:

Beispiel einer Filtering-Lösung

Früher wurden Filtering-Datenbanken auf manuelle Art zeitaufwändig aufgebaut. Dagegen wird die Websense

Master Database sehr präzise und hochautomatisch entwickelt, gepflegt und auf dem neusten Stand gehalten. Die Websense Master Database war ursprünglich in 27 Kategorien aufgeteilt. Die Kategorien waren nach ausführlichen Diskussionen mit Unternehmen, Schulen und Anwendern über ihre Filtering-Ansprüche definiert worden. Nach zweijähriger Analyse der Kundenvorschläge wurde die Datenbank weiter verfeinert und Mitte 1999 auf 54 Kategorien erweitert. Seit dieser Zeit wurde sie kontinuierlich erweitert, rekonstruiert und umfasst heute mehr als 65 Kategorien. Mit der Neudefinierung der Kategorien konnten Doppeldeutigkeiten und Überlappungen in Kategorien und Subkategorien minimiert werden. Kontrolle und Flexibilität, sowie der Gebrauch einer klaren, modernen Sprache wurden forciert¹. Je nachdem, inwieweit der Inhalt einer Site der Kategorienbeschreibung entspricht, werden die Sites kategorisiert. Virtuelle Hosts und Webserver, die so konfiguriert sind, um mehrere Websites oder Domains aufzunehmen, haben den Kategorisierungsaufwand verkompliziert. Obwohl die meisten Hosts jeder aufgenommenen Site eine IP-Adresse zuordnen, legen einige virtuelle Hosts für alle Websites und Domains die gleiche IP-Adresse an. So ordnet möglicherweise ein Webserver eine Shopping Site, eine Religions-Site und eine Site mit Erwachsenenunterhaltung einer gleichen IP-Adresse zu. Dagegen hat Websense

¹Die gegenwärtige Kategorienliste der Websens Master Database kann auf der Websens Website eingesehen werden: www.websens.com/products/categories/version4.cfm

Websense derzeit einen Filter, der alle nicht kategorisierten, potenziell pornografischen Sites auflistet. Der Netzwerkadministrator kann die Liste während des Downloads der Master Database zu Websense zurückschicken. Web-Analysten haben dann die Möglichkeit, die Sites zu überprüfen und wenn angebracht, in die Master Database aufzunehmen. Fremdsprachige Sites werden mit denselben Techniken durchsucht, jedoch mit Worten, Hubs und Suchmaschinen, die der jeweiligen Sprache entsprechen.

Klassifizierung der Sites mit KILO II

Nach dem Mining der Websites werden sie mit KILO II klassifiziert. Knowledge Indexing Learning Organization (KILO) ist ein von Websense entwickeltes Software-System zur Klassifizierung von Websites. KILO II ist die zweite Generation algorithmischer Klassifizierer, die allgemein anerkannte Lerntechniken verwendet. Die Software lernt, wie Sites einzuordnen sind, indem sie sich an der bestehenden Websense Database orientiert. Die Effizienz adaptiv trainierter Klassifizierer entspricht der Qualität der Daten ihrer Trainingsgrundlage. Da Websense eine ausführliche, von Menschen entwickelte und überprüfte Datenbank hat, ist die Qualität der Trainingsgrundlage sehr hoch. Entsprechend ist die Genauigkeit einer Klassifizierung mit KILO.

Algorithmische Klassifizierer sind nach Websense nicht ausreichend präzise für dynamisches oder Runtime Filtering. Aber sie können besonders gut für die Entwicklung einer Datenbank für kategorisierte Sites eingesetzt werden. KILO II arbeitet bei der Klassifizierung mit einem Vertrauensfaktor. Liegt dieser über einer bestimmten Schwelle, wird die Site automatisch in die Master Database integriert. Genauso kann die Software festlegen, dass eine Site in keine der Kategorien passt und deshalb nicht in die Master Database integriert wird. Alles, was dazwischen liegt, wird in Ergebnislisten sortiert, die dann an das Workbench der Analysten weitergeleitet und überprüft werden.

Was ist ein Workbench?

Nach dem Mining der Sites und der Bearbeitung mit KILO II beginnt die menschliche Überprüfung. Websense hat ein umfassendes Set an Tools entwickelt, um die Produktivität der Überprüfung zu maximieren. Die folgenden

Tools nennt man zusammengenommen Web Analyst's Workbench. SingleClick ist eine Applikation zur schnellen Bearbeitung der zu prüfenden Website-Listen, welche die Verbreitung der Sites verifiziert, zwischenspeichert und dem Web-Analysten in aufeinanderfolgenden Sequenzen präsentiert. Mit QuickPost können Web-Analysten schneller surfen und meist mit nur einem Mausklick neue Sites in die Kategorien der Master Database aufnehmen.

Qualitäts- und Aktualitätskontrolle

Die Websense Master Database wird regelmässig auf den neusten Stand gebracht und gewährleistet Qualität durch ständige Kontrolle den Analyseprofis und automatische Prozesse. Zwei Tools des Web Analyst Workbench garantieren Qualität und ein gründliches Update der Datenbank. Mit dem Verifier Tool wird die Genauigkeit der Website-Listen des Vortags verifiziert. Das Quality Control Tool sorgt dafür, dass die Sites mit dem ältesten Aufnahmedatum in Listen gesammelt und von den Web-Analysten mit Hilfe des Tool überprüft werden. Die Dimitrius Differential Engine entwickelt eine kontextuelle «Punktzahl» für ein ausgewähltes Set von Sites einer Website. An festgelegten Zeitpunkten entwickelt Dimitrius eine neue Punktzahl für das gleiche Set von Sites und vergleicht die alte mit der neuen Punktzahl. Ist die Differenz grösser als eine zuvor festgelegte Schwelle, wird angenommen, dass sich die Site wesentlich verändert hat. Das hat eine sofortige Markierung der Site zur Folge und zieht eine entsprechende Prüfung durch den Web-Analys-

Info:

Websense
Johanna Severinsson
300 Hillwood Drive
UK-Chertsey, Surrey KT16 0RS
Tel. +44 (0)1932 796-001
Fax +44 (0)1932 796-601
E-Mail: jseverinsson@websense.com

ten nach sich, damit eine korrekte Kategorisierung gewährleistet werden kann. Zusätzlich läuft alle zehn Stunden ein DNS- Update-Programm der gesamten Datenbank.

Websense hat eine hybride Lösung für das Filtern von Internet Sites entwickelt, an deren Ende die Websense Master Database als benutzerfreundliche und präzise Filtering-Datenbank steht. Die hybride Lösung von Websense arbeitet sowohl mit dynamischem Filtering zur Maximierung des Recall wie auch mit einem umfassenden Datenbank-Filtering zur Maximierung der Precision. Mit der bestehenden Database als Training Set wurde die Master Database durch den Einsatz von Metasearch-Technologie und intelligenten Crawling-Techniken optimiert. 11, 4

*Harold Kester, Chief Technical Officer,
Websense, Chertsey*

Summary

Internet filtering: how can Web sites be filtered effectively?

It is quite clear that the Internet offers companies enormous benefits such as increased productivity, improved communications and a means of providing information comprehensively. Nevertheless the danger of misuse is very real. Cyber slacking, spending time playing computer games or shopping during working hours, costs companies large sums of money and calling up sites containing pornographic or racist material at the workplace can be extremely offensive and result in legal action against companies. Companies must establish guidelines if misuse of the Internet is to be avoided. Software solutions for the implementation of guidelines governing Internet use are known generally as Internet filters.

Wer uns heute für **Informatik** und **Kommunikation** kontaktiert, profitiert schon morgen davon.

SOHARD AG – Generalunternehmen für

- Digital Audio Broadcast Solutions
- Globale Informations-Systeme wie Postphone, Bankphone, Fahrgast, Parkplatz
- Flottenmanagement-Systeme für Transportunternehmen, Rettungs- und Pannendienste
- Oracle based Solutions
- Mobile Datenverarbeitung für Aussendienst, Service, Verkauf
- Internet, Intranet, E-Commerce
- Service, Support, Sicherheit



SOHARD AG

Software/Hardware Engineering
Galgenfeldweg 18, CH-3000 Bern 32
Tel. 031 33 99 888, Fax 031 33 99 800
E-Mail: sohard@sohard.ch
Internet: www.sohard.ch



ISO 9001 Reg.-Nr. 10909-02

GIGACOMP AG

The art of solutions



Hochfrequenz und Mikrowellentechnik

Wir beschäftigen uns mit innovativen Produkten für High-Tech Anwendungen von Mobilfunktechnik bis zur Satellitenverbindung. Das Produktangebot umfasst Messtechnik wie Netzwerk- und Spektrum-Analyzer, Power Meter, Signalquellen, Verstärker und Antennen.

Telecom und Datenkommunikation

Wir bieten Produkte und Dienste in den drahtgebundenen sowie drahtlosen Technologien an. Unser Angebot umfasst Messlösungen in der digitalen Informationsübertragung und der Analyse von Signalen in der mobilen Kommunikation.

Optische Messinstrumente

Unser Produktangebot umfasst optische Netzwerk- und Spektrumanalyzer, OTDR's oder Signalquellen, Messzubehör aber auch eine Vielzahl von Komponenten wie Laserdioden oder breitbandige Verstärker.

Dienstleistungen

Vielseitige und flexible Dienstleistungen wie Projektmanagement, Schulung, Hardware-Support und Kalibration ergänzen unser Angebot für innovative Gesamtlösungen.

GIGACOMP AG
Gewerbezone Lätti
CH-3053 Münchenbuchsee
Telefon 0041-(0)31-868 44 55
Telefax 0041-(0)31-868 44 50

info@gigacomp.ch
www.gigacomp.ch

