

Ontologiebasierte Dokumentindizierung und -suche



Oktober 2000

JonnyNewald



Die größte Dokumentensammlung der Welt ist das Internet.

Einführendes Beispiel: Information Retrieval

Stellen Sie sich vor, Sie suchen im **Internet** Informationen über ein bestimmtes Thema.

Suchmaschinen bieten Volltextsuche, Metadaten- und eine hierarchische Kategorisierung.

Dokumente über ein sehr **spezielles Thema**, wie z. B., „Fernsehen“, sind nur sehr schwer exakt auszumachen. Der Begriff „Fernsehen“ hat auch mehrere Bedeutungen und steht in einem gewissen **Kontext**. Meine ich, „Fernsehgeräte“, das „Fernsehprogramm“ oder das Medium „Fernsehen“?

Die Volltextsuche liefert zu viel, die Metadaten- und nicht
verlässliche, die Kategorisierung ist oft zu oberflächlich und sem
antisch nicht
eindeutig.

Nicht nur im Internet stellt sich diese Problematik, auch in Unternehmen.

Ausgangsbasis

Lotus Notes ist

- ✓ ein umfangreiches Dokumentenhaltungs- und Groupwaresystem
- ✓ unternehmensweit einsetzbar
- ✓ in das vorhandene Intranet integrierbar
- ✓ über WEB-Browser benutzbar.

Lotus Notes hält die gesamte Information in **strukturierten Dokumenten**, die sich in **Dokumentendatenbanken** befinden.

Hohes Datenaufkommen bedingt das Problem der Wiederauffindbarkeit.

Knowledgeger bietet weitergehende Hilfsmittel.

Ausgangsbasis (Fortsetzung)

Das Produkt **Knowledgeger** ist

- ✓ eine spezialisierte Lotus Notes - Anwendung der Firma *Knowledge Associates*
- ✓ eine Reihe spezieller Datenbankschablonen, die der Haltung verschiedener Arten von Wissensdokumenten dienen
- ✓ optimiert für die Benutzung über WEB - Browser.

Knowledgeger bietet Unterstützung bei der Bestimmung der Dokument-**Metadaten** und der **Kategorie-Informationen**.

Knowledgeger's Ansätze sind konventionell und einfach.

Knowledgeger ist ein einfaches Wissensmanagement -System.

Wissensmanagement unter Lotus Notes/Knowledgeger

Lotus Notes bietet

- ✓ eine automatische Pfleger der Standard-Metadaten
- ✓ grundsätzliche Möglichkeiten zur Definition spezifischer Ansichten auf Dokumentlisten (Views)
- ✓ eine Suchmöglichkeit über Volltextsuche (auch in File -Attachments)

Knowledgeger bietet

- ✓ verschiedene, thematisch getrennte Datenbanken
- ✓ Anwendergruppenorientierte Navigatoren (getrennt für Administratoren, Manager und normale Mitarbeiter)
- ✓ im Dokument abgelegte, frei editierbare Zusatzfelder zur inhaltlichen Kurzbeschreibung

Ein effiziente inhaltsbezogene Suche bedarf ausgeklügelte Ansätze.

Grundbedürfnisse effizienterer Lösungen

- ✓ Die Beschränkung auf ein **kontrolliertes Vokabular** beim Kategorisieren und Suchen vermeidet Inkonsistenzen und erhöht die Trefferquote.
- ✓ Eine **Projektion des Dokumentinhalts in das Wissensmodell des Unternehmens** läßt sich formalisieren und durch Computerverarbeitung unterstützen.
- ✓ **Ansprechendere Benutzerschnittstellen** fördern die Motivation der Mitarbeiter.

Ein Lösungsansatz ist die Verwendung graphischer Wissensmodelle.

Einfache ontologische Modelle

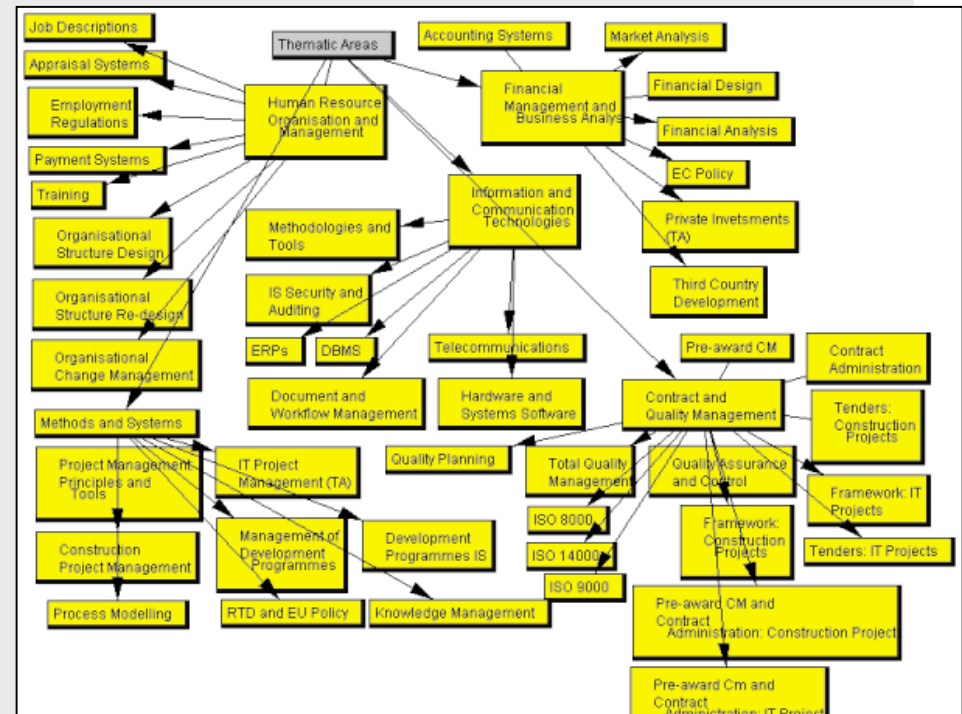
- ✓ Eine abgegrenzte **Wissensdomäne** läßt sich formal und abstrakt durch eine **Ontologie** modellieren.
- ✓ Hierzu überlegt man sich relevante **Konzepte** und **Begriffe**, die eindeutig sind, und die man miteinander in **Beziehung** setzt.
- ✓ Die einfachsten Strukturen sind **Hierarchien**. In vielen Fällen genügt eine solche Struktur.
- ✓ Dokumente werden durch eine bestimmte Auswahl aus der Begriffsmenge indiziert und sind somit in das Wissensmodell des Unternehmens projiziert.

Hierarchische Begriffsmodelle werden jedoch sehr schnell unübersichtlich.

ZukomplexeBegriffshierarchiensollteninTeilmodelleunterteilt werden.

KomplexeBegriffsmodelle

- ✓ KomplexeModellesindunter Umständennotwendig,ihre **Unübersichtlichkeit** kann jedochgebrochenwerden.
- ✓ Abhilfebieten **hypertextmäßig verschachtelteTeilmodelle**.
- ✓ EinsolchesTeilmodellwird dannnichtzusammenmitdem übergeordnetenModell dargestellt.



BeispieleinerkomplexenOntologieohneSchachtelung

DieUmsetzungdesgeschachteltenModellansatzeserfordertneueSoftware.

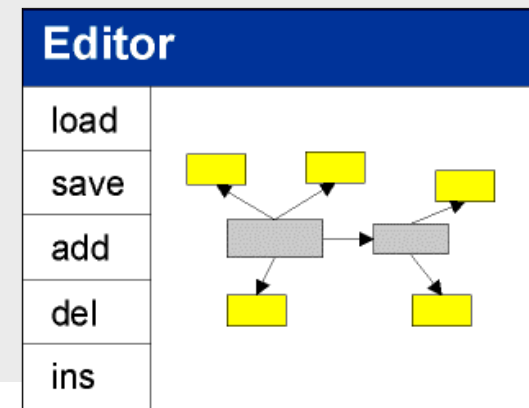
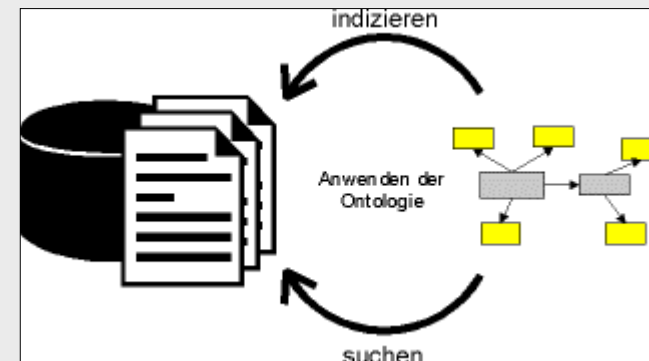
Erforderliche Softwarekomponenten

- ✓ Eine Benutzerschnittstelle zur ontologiebasierten Indizierung und Suche bietet guten Bedienkomfort:

das **Index-Retrieval-Interface (IRI)**

- ✓ Ein graphischer Editor für den komfortablen Entwurf der Ontologien:

der **Ontology Editor**



Beiden Komponenten arbeiten mit derselben zentral gehaltenen Ontologie.

DasIRI -Fenster,einJAVA -Frame,kommuniziertimIndiziermodusmitdemBrowser.

DasIndex -Retrieval-Interface:Indizieren

1. EinLotusNotes - Dokument wirdganznormalimWEB - Browserpräsentiert.
2. NachBetätigungeinesButtons erfolgtdie **Indizierung** über diegraphischeSelektionaus derdurcheinenJAVA -Frame dargestelltenOntologie.
3. DieKategorisierungsinfor - mationenwerdendirektinein speziellesFeldübertragen.

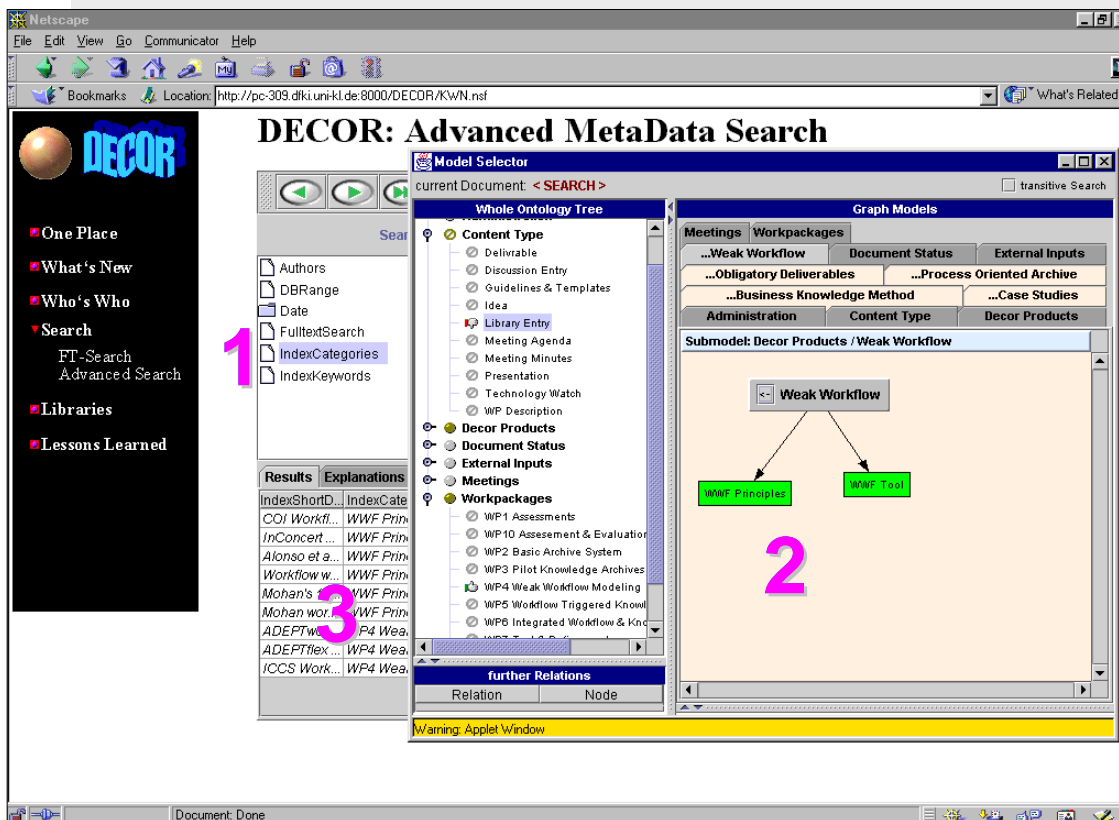
The screenshot displays the ADEPT workflow interface within a Netscape browser window. The interface is divided into several sections:

- Document Metadata:** A form for entering document information. A red arrow labeled '3' points to the 'Short Description' field, which contains the text: 'Please enter a Short Description in t...'. Other fields include 'Categories' (WP4 Weak Workflow Modeling, Workflow, WWF Principles, Technology Watch, WWF Tool), 'Keywords', and 'Status'.
- Model Selector:** A panel on the right showing a 'Whole Ontology Tree' with categories like Administration, Content Type, Decor Products, Document Status, External Inputs, Meetings, and Workpackages. A red arrow labeled '2' points to the 'External Inputs' category.
- Graph Models:** A central area displaying a graph model with 'External Inputs' at the center, connected to various nodes like Workflow, CognitoVision, Agents, UML, Ontologies, SW Methodologies, XML, Java, and SmartFinder.

A green arrow labeled '1' points from the 'Delete Document' button to the document title 'ADEPTworkflow - Advan...'. A red arrow labeled '2' points from the 'External Inputs' node in the graph model to the 'choose' button at the bottom right of the interface.

DasselbeInterfacedientderSucheüberdiegraphischenModelle .

Das Index -Retrieval-Interface: Suchen

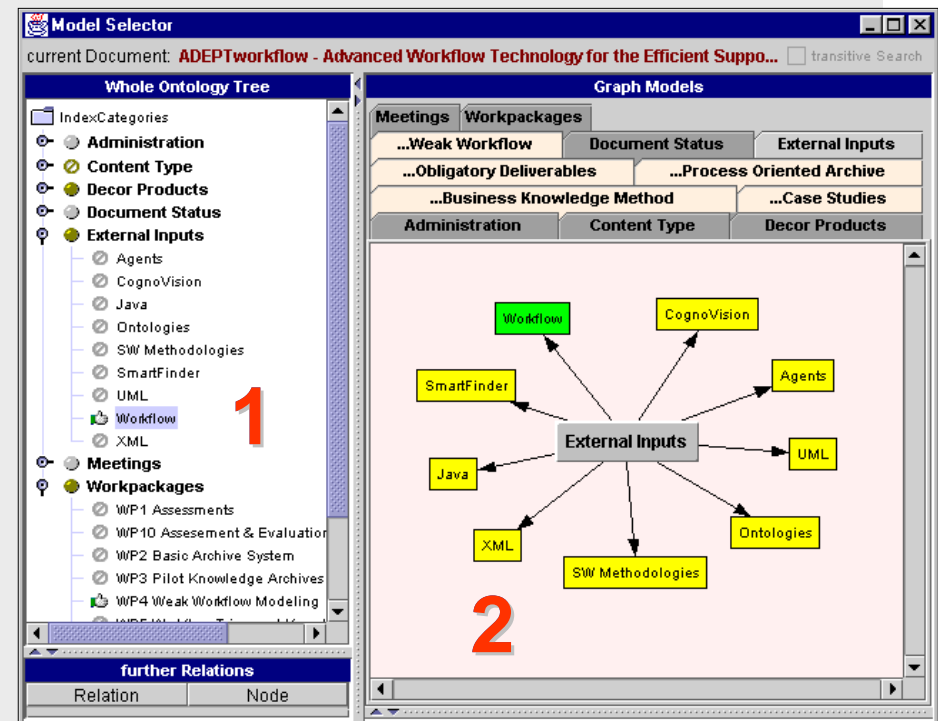


1. Nach Selektion des Vater - konzepts für die graphischen Modelle im Suchapplet befindet sich das IRI im Suchmodus.
2. Jede einzelne Knoten - Selektion führt zu einer neuen Abfrage.
3. Das Abfrageergebnis wird sofort als Liste präsentiert, aus der einzelne Dokumente geöffnet werden können.

Zu jeder Zeit existiert höchstens eine IRI - Instanz bzw. dessen JAVA - Frame.

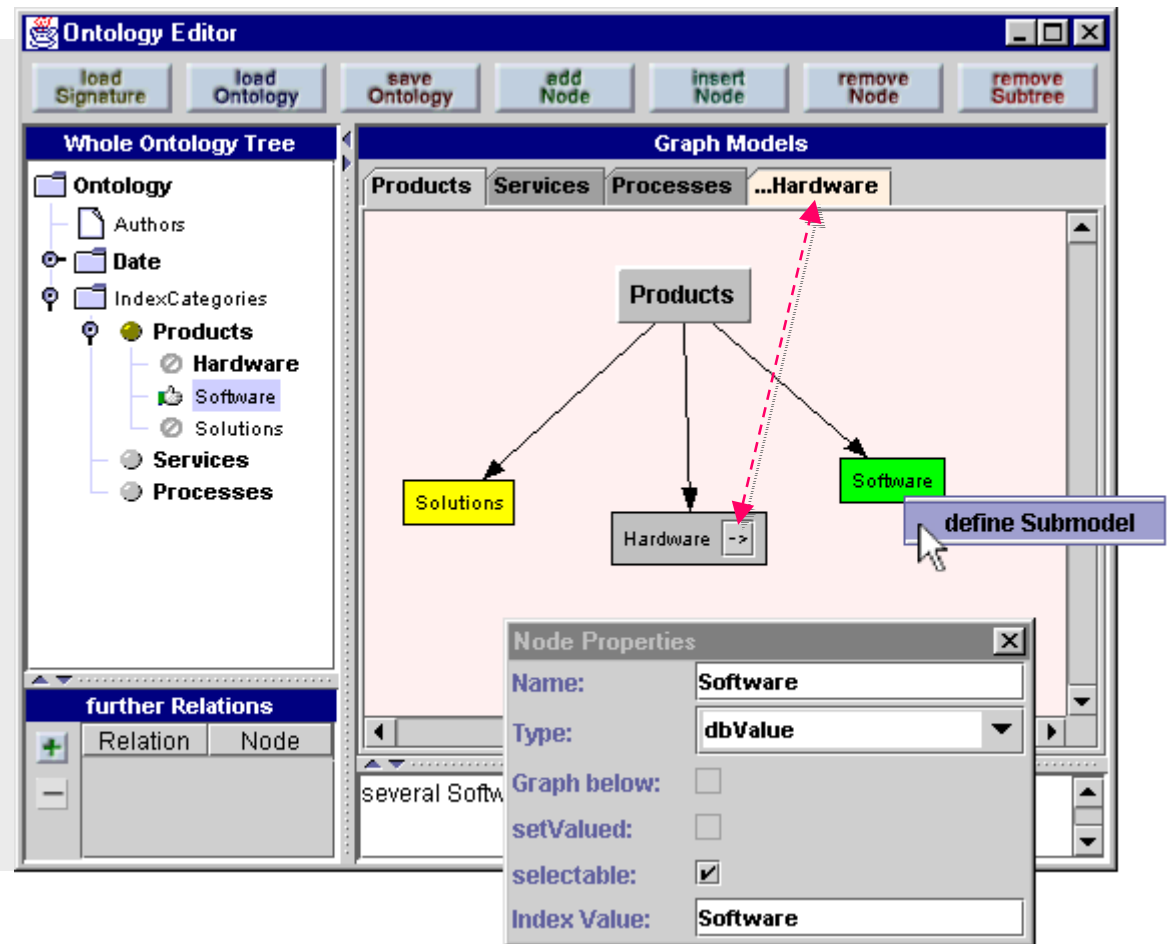
Merkmale des Index -Retrieval-Interface im Überblick

- ✓ reine Java -Komponente (in Knowledgeger eingebettet)
- ✓ komfortable Oberfläche durch Verwendung der JAVA -Klassenbibliothek **SWING**
- ✓ nahtloser Übergang zwischen Such - und Einstell -Anwendung: automatischer **Moduswechsel**
- ✓ einfache Navigation durch komplexe Ontologien mittels einer synchronisierten Baum- und Graph -Darstellung(1, 2)



Der Ontology Editor

- ✓ Der Editor erzeugt graphisch **verschachtelte Teilmodelle**, die über Linkbuttons miteinander verbunden sind.
- ✓ Die Darstellung und Interaktion ist der des **IRI** identisch.
- ✓ Er bietet die Möglichkeit der Definition beliebiger **Querbeziehungen** unter Klärung der **Knotenkommentare**.

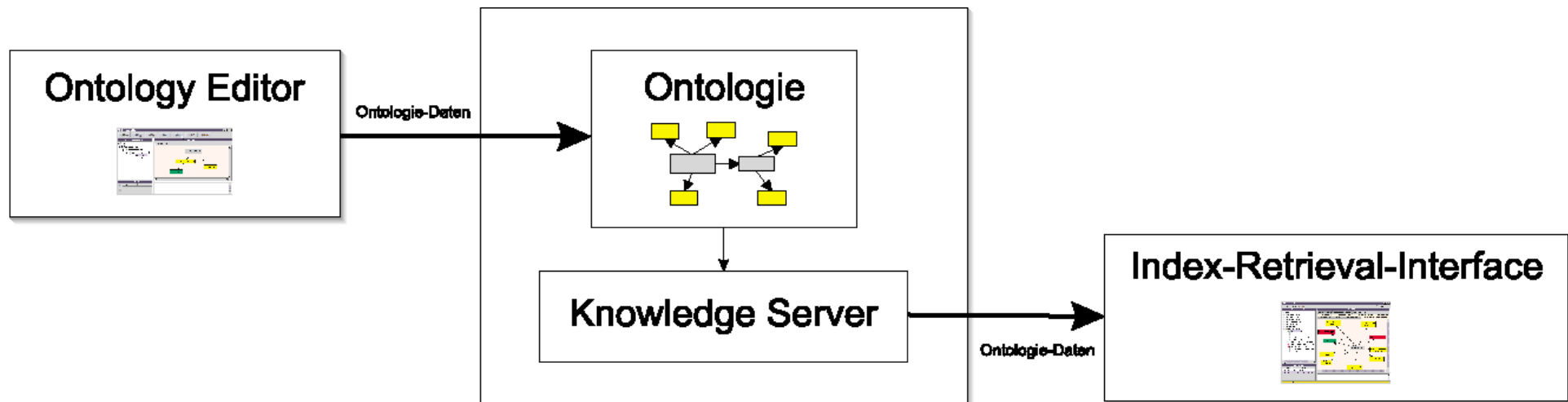


Die Editor -Anwendung läßt sich auf jedem JAVA -fähigen Rechner starten.

Der **KnowledgeServer** ist der zentrale Verwalter der Ontologiedaten.

Der Ontologiestrom

Eine im Editor erstellte Ontologie wird dem **KnowledgeServer** übergeben, sodass das Index-Retrieval-Interface mit den Modell-daten beliefern kann.



Das Konzept der Knowledge Servers stammt aus dem

[KnowNet-Projekt](#).

Ausblick: Ontologiemodifikationen

- ✓ Im Editor umgesetzt ist bereits die **Protokollierung** der getätigten Änderungen.
Nach der Speicherung einer geänderten Ontologie wird eine Änderungsprotokoll-Datei angelegt.
- ✓ Die Datei liefert einen Ansatzpunkt zur weiteren **Verarbeitung**, wie
 - automatische Unterrichtung der Autoren über die Ontologieänderungen
 - „intelligente“ Anpassung der Kategorisierungsinformation der betroffenen Dokumente an die neue Situation

Solche automatischen Reaktionen bedürfen weiterer konzeptioneller Überlegungen.

Ausblick: Beziehungssemantiken

Beispiel:

In einer IT-Beraterfirma existierte ein Teilmodell, das eine Beziehungsstruktur verschiedener Softwareprodukte bestimmt. Man könnte sich eine Beziehungsart „is Incompatible“ vorstellen, die beispielsweise zwischen dem Datenbankprodukt „Informix“ und der Betriebssystemgruppe „MS Windows“ definiert ist.

Solche beliebige Beziehungen können im Editor zwar formal definiert werden, es fehlt jedoch noch eine entsprechende Verarbeitung einer zu definierenden **Semantik** der Beziehungen.

Sie könnte eine Dokument-Abfrage, die aus der Selektion beider Produkte besteht, im Vorfeld abgewehrt werden mit dem Kommentar der Inkompatibilität.

Die Einführung der Beziehungssemantiken ist für eine intelligente Suche unabdingbar.

Die wesentlichen Aspekte dieser Arbeit zusammengefaßt...

Essenz

- ✓ Umsetzung des Prinzips der **hypertextmäßig verschachtelten Teilmodelle**
- ✓ **Vereinheitlichte Oberfläche** in allendrei Anwendungen (Indizieren, Suchen, Editieren)
- ✓ generell vereinfachte Bedienung
- ✓ Einbettung der neuen Komponenten in das vorhandene Programmsystem von **KnowNet**