

Komponentensuche mit Begriffen

Christian Lindig

Institut für Programmiersprachen, Abteilung Softwaretechnologie

Technische Universität Braunschweig, 38106 Braunschweig

lindig@ips.cs.tu-bs.de

Zusammenfassung

Mit Schlüsselwörtern indexierte Komponenten einer Sammlung werden einmalig durch formale Begriffsanalyse zu einem Begriffsverband strukturiert. Dieser erlaubt dann, Komponenten durch die Angabe von Schlüsselwörtern einfach, effizient, inkrementell und mit Unterstützung für den Benutzer zu suchen. Der Suchprozeß garantiert dabei, daß bei einer Suche mindestens eine Komponente gefunden wird. Da der Suchprozeß unabhängig vom Inhalt der verwalteten Komponenten ist, ist er auch für inhomogene und multimediale Sammlungen geeignet. Eine prototypische Implementierung zur Verwaltung einer Unix Online-Dokumentation belegt die Anwendbarkeit des Verfahrens, dessen Komplexität zusätzlich experimentell untersucht wurde. Es zeigt sich, daß das Verfahren eine schnelle und genaue Navigation in Sammlungen mit mehreren tausend Komponenten ermöglicht.

Stichworte: Softwarebibliothek, Wiederverwendung, Inkrementelle Suche, Begriffsverband

1 Einleitung

In der Praxis haben sich Information-Retrieval-Ansätze [11, 7, 9] zum Verwalten von wiederverwendbaren Software-Komponenten als geeignet erwiesen. Im Gegensatz zu formal-logischen Ansätzen [8, 12, 13, 2] sind sie nicht auf spezielle Programmiersprachen festgelegt und ihre Anfragesprache ist einfacher, wenn auch zum Teil mit unklarer Semantik. Die in dieser Arbeit präsentierte Komponentensuche mit Begriffen verbindet eine einfache und klare Anfragesprache mit einer effizienten, inkrementellen Suche, die den Benutzer kontextsensitiv unterstützt. Dazu bestimmt *formale Begriffsanalyse* zu einer indexierten Sammlung einmalig die formalen Begriffe, die dann für spätere Suchen verwendet werden.

Jede Komponente wird zunächst mit einer beliebigen Anzahl von Schlüsselwörtern indexiert. Die berechneten formalen Begriffe der Sammlung bilden einen Verband aus Ober- und Unterbegriffen. Die Suche erfolgt inkrementell durch die Angabe von den Schlüsselwörtern, die die gesuchten Komponenten mindestens aufweisen sollen. Dabei wählt der Benutzer die Schlüsselwörter aus einer präsentierten Menge

aus, wobei der Auswahlprozeß garantiert, daß mindestens eine Komponente selektiert wird – das Ergebnis also nicht leer ist. Mit jedem Schritt verringert sich sowohl die Zahl der selektierten Komponenten, als auch der im nächsten Schritt wählbaren Schlüsselwörter, so daß der Benutzer präzise und schnell zu den gesuchten Komponenten geführt wird.

Der nächste Abschnitt gibt zunächst einen Überblick über das Verfahren in Form eines Beispiels. Die sich daran anschließenden Abschnitte definieren Begriffsanalyse formal und zeigen, wie sie zur Suche eingesetzt wird. Die Vorstellung eines Prototyps, Experimente zur Aufwandsabschätzung und ein Vergleich mit anderen Ansätzen bilden den Schluß.

2 Ein Beispiel: Unix-Dokumentation

Unix-Systeme besitzen traditionell eine Online-Dokumentation ihrer Programme, Systemaufrufe, Bibliotheken und Dateiformate, die zusammen etwa 1800 Dokumente umfaßt. Ein kleiner Ausschnitt aus der Dokumentation der Betriebssystemaufrufe (System Calls) soll demonstrieren, wie Dokumente, stellvertretend für Komponenten bestehend aus Programmcode, einer Signatur und andere für ihre Verwendung relevanten Informationen, indexiert und gesucht werden können. Abbildung 1 zeigt die Namen von 12 Unix System-Calls zusammen mit einer einzeiligen Beschreibung aus ihrer Dokumentation [15] und einer Menge von Schlüsselwörtern als Indexierung. Im Fall der Unix-Dokumentation ist die Indexierung oft naheliegend, da bereits eine Nomenklatur existiert: Dateien sind *files*, Zugriffsrechte *permissions*. Im allgemeinen muß aber zunächst ein Katalog von Schlüsselwörtern und ihrer Semantik erstellt werden.

Nachdem die Komponenten indexiert sind, wird die Sammlung durch eine formale Begriffsanalyse für den späteren Zugriff zu Begriffen strukturiert. Der zugehörige Begriffsverband (siehe nächster Abschnitt) ist lediglich die unterliegende Datenstruktur für die Suche, und ist dem Benutzer nicht sichtbar.

Die Suche nach Komponenten verläuft schrittweise: in jedem Schritt wählt der Benutzer eines der Attribute, die die gesuchten Komponenten *mindestens* aufweisen sollen. Sofort werden ihm alle Komponenten präsentiert, die mindestens die bisher ausgewähl-

Sys.-Call	Kurzbeschreibung	Indexierung
chmod	<i>change mode of file</i>	change mode permission file
chown	<i>change owner and group of a file</i>	change owner group file
stat	<i>get file status</i>	get file status
fork	<i>create a new process</i>	create new process
chdir	<i>change current working directory</i>	change directory
mkdir	<i>make a directory file</i>	create new directory
open	<i>open or create a file for reading or writing</i>	open create file read write
read	<i>read input</i>	read file input
rmdir	<i>remove a directory file</i>	remove directory file
write	<i>write output</i>	write file output
creat	<i>create a new file</i>	create new file
access	<i>determine accessibility of file</i>	check access file

Abbildung 1: Unix-System-Calls mit Kurzbeschreibung und Indexierung

Schritt	Attribut	Komponenten	noch mögliche Attr.
1	–	alle	alle
2	<i>change</i>	chdir chmod chown	<i>directory file group mode owner permission</i>
3	<i>change file</i>	chmod chown	<i>group mode owner permission</i>
4	<i>change file mode</i>	chmod	–
1	–	alle	alle
2	<i>create</i>	creat fork mkdir open	<i>directory file new open process read write</i>
3	<i>create file</i>	creat open	<i>new read open write</i>
4	<i>create file read</i>	open	–

Abbildung 2: Zwei Beispiele für schrittweise Suche

	access	change	check	create	directory	file	get	group	input	mode	new	open	output	owner	permission	process	read	remove	status	write
access	x	x			x															
chdir		x			x										x					
chmod		x							x						x					
chown		x						x						x						
creat				x	x						x									
fork				x							x					x				
fstat						x	x													x
mkdir				x	x						x									
open				x	x							x						x		x
read						x			x								x			
rmdir					x	x												x		
write						x							x							x

Abbildung 3: Kontexttabelle *Unix-System-Calls* für Systemaufrufe des Betriebssystems Unix

ten Attribute aufweisen und gleichzeitig werden ihm die Attribute präsentiert, die er im nächsten Schritt wählen kann. Die zur Auswahl stehenden Attribute sind genau so bestimmt, daß jede Wahl die Menge der selektierten Komponenten verkleinert, sie aber *garan-*

tiert nicht leer ist. Mit anderen Worten: der Benutzer kann keine Attribut-Kombinationen wählen, die keine der Komponenten erfüllt. Bei der Auswahl der Attribute ist die Reihenfolge beliebig: der Benutzer kann sich den gesuchten Komponenten unter verschiedenen Aspekten nähern. Auch bei großen Sammlungen erfolgt die Berechnung des Ergebnisses nach einem Suchschritt praktisch verzögerungsfrei.

Abbildung 2 zeigt zwei Beispiele für schrittweise Suchen; die in Abbildung 1 vorgenommene Indexierung ist in Abbildung 3 nochmals als Tabelle dargestellt. Solange kein Attribut ausgewählt ist, sind alle Komponenten der Sammlung selektiert und alle Attribute stehen noch zur Auswahl. Nachdem der Anwender im ersten Beispiel das Attribut *change* gewählt hat, werden dadurch die Komponenten *chdir*, *chmod* und *chown* selektiert, da sie mit *change* indexiert sind. Für den nächsten Schritt stehen nur noch die Attribute zur Auswahl, die zusammen mit *change* eine nicht leere Selektion ergeben. Der Anwender wählt im Beispiel *file* aus, wodurch sich sowohl die Menge der selektierten Komponenten, als auch der noch möglichen Attribute weiter verkleinert. Es sind jetzt genau die Komponenten selektiert, die mindestens mit *change*

und *file* indexiert sind. Im letzten Schritt wird durch die Auswahl von *mode* noch genau die Komponente **chmod** selektiert. Natürlich kann der Benutzer seine Suche auch schon vorher beenden, und so mehrere selektierte Komponenten erhalten.

Die Zahl der noch möglichen Attribute verringert sich in jedem Schritt im allgemeinen stark, so daß der Benutzer nach wenigen Schritten nur noch aus einer überschaubaren Menge auswählen muß. Könnte er in einem Schritt ein nicht angegebenes Attribut wählen, würde dieses mit den zuvor gewählten Attributen im Widerspruch stehen – keine Komponente weist die entsprechenden Attribute gemeinsam auf.

Die Grundlage für die effiziente Berechnung der selektierten Komponenten und der im nächsten Schritt einer Suche möglichen Attribute ist die vorhergehende Bestimmung der *Begriffe* der Sammlung. Ein Begriff ist ein maximales Paar aus einer Menge von Attributen und der zugehörigen Menge von selektierten Komponenten. So bilden im dritten Schritt des zweiten Beispiels *create*, *file* und **creat**, **open** einen Begriff – allerdings bildet nicht jedes Paar aus Anfrage und selektierten Komponenten einen Begriff, da die zu einem Ergebnis gehörende Anfrage nicht alle die den Komponenten des Ergebnisses gemeinsamen Attribute enthalten muß.

3 Formale Begriffsanalyse

Formale Begriffsanalyse wurde von R. Wille begründet und untersucht die Beziehung zwischen einer Menge von Objekten und diesen Objekten zugeordneten Attributen [16, 1].

Definition 1 *Ein formaler Kontext ist ein Tripel $(\mathcal{O}, \mathcal{A}, \mathcal{R})$ aus einer endlichen Objektmenge \mathcal{O} , einer endlichen Attributmengemenge \mathcal{A} und einer Relation $\mathcal{R} \subseteq \mathcal{O} \times \mathcal{A}$ zwischen ihnen. $(o, a) \in \mathcal{R}$ wird gelesen als: Objekt o besitzt das Attribut a .*

Die Relation zwischen Objekten und Attributen läßt sich anschaulich in einer *Kontexttabelle* darstellen. Dies ist genau die schon in Abbildung 3 gezeigte Tabelle mit System-Calls als Objekten und Schlüsselwörtern als Attributen.

Definition 2 *Für Objekt-Mengen $O \subseteq \mathcal{O}$ und Attributmengen $A \subseteq \mathcal{A}$ eines Kontextes $(\mathcal{O}, \mathcal{A}, \mathcal{R})$ werden die Mengen der ihnen gemeinsamen Attribute $\omega(O)$ bzw. Objekte $\alpha(A)$ definiert:*

$$\begin{aligned}\alpha(A) &:= \{o \in \mathcal{O} \mid \forall a \in A : (o, a) \in \mathcal{R}\} \\ \omega(O) &:= \{a \in \mathcal{A} \mid \forall o \in O : (o, a) \in \mathcal{R}\}\end{aligned}$$

Zum Beispiel besitzen die Objekte **chdir** und **chmod** nur *change* als gemeinsames Attribut, also gilt $\omega(\{\text{chdir}, \text{chmod}\}) = \{\text{change}\}$.

Definition 3 *Ein formaler Begriff eines Kontextes $(\mathcal{O}, \mathcal{A}, \mathcal{R})$ ist ein Paar (O, A) aus Objekt- und Attributmengen mit $O \subseteq \mathcal{O}$, $A \subseteq \mathcal{A}$, wobei gleichzeitig gilt:*

$$\alpha(A) = O \quad \text{und} \quad \omega(O) = A$$

Für einen Begriff $b = (O, A)$ bezeichnen $\pi_o(b)$ die Objekte (Umfang) und $\pi_a(b)$ die Attribute (Inhalt) des Begriffs: $\pi_o(b) := O$ und $\pi_a(b) := A$. Die Menge aller Begriffe eines Kontextes wird mit $B(\mathcal{O}, \mathcal{A}, \mathcal{R})$ bezeichnet.

Die Objekte eines Begriffs (in einem gegebenen Kontext) sind synonym für eine Menge von Attributen und umgekehrt – die Objektmenge eines Begriffs bestimmt die Attributmengemenge eindeutig. Ein Beispiel für einen Begriff des Kontexts Unix-System-Calls ist das Paar $(\{\text{chmod}, \text{chown}\}, \{\text{file}, \text{change}\})$. In diesem Kontext sind die Attribute *file* und *change* ein Synonym für die Objekte **chmod** und **chown** – und umgekehrt.

Zur Berechnung aller Begriffe $B(\mathcal{O}, \mathcal{A}, \mathcal{R})$ eines Kontextes existieren mehrere Algorithmen [3], deren Komplexität im ungünstigsten Fall exponentiell ist, da ein Kontext maximal 2^n mit $n = \min(|\mathcal{O}|, |\mathcal{A}|)$ Begriffe enthalten kann. Der typische Aufwand ist aber wesentlich günstiger und wird unten noch genauer untersucht.

Definition 4 *Zwei Begriffe $(O_1, A_1), (O_2, A_2) \in B(\mathcal{O}, \mathcal{A}, \mathcal{R})$ eines Kontextes werden durch die Unterbegriffs-Relation \leq geordnet:*

$$(O_1, A_1) \leq (O_2, A_2) \stackrel{\text{def}}{\iff} O_1 \subseteq O_2 \quad (\stackrel{\text{def}}{\iff} A_1 \supseteq A_2)$$

Die Relation wird gelesen als: (O_1, A_1) ist Unterbegriff von (O_2, A_2) . Die Struktur aus Kontext und Ordnungsrelation wird abgekürzt als $\mathcal{B}(\mathcal{O}, \mathcal{A}, \mathcal{R}) = (B(\mathcal{O}, \mathcal{A}, \mathcal{R}), \leq)$.

In dem gegebenen Beispiel ist der Begriff $(\{\text{mkdir}, \text{creat}, \text{fork}\}, \{\text{create}, \text{new}\})$ ein Oberbegriff des Begriffes $(\{\text{fork}\}, \{\text{create}, \text{new}, \text{process}\})$. Ein Oberbegriff hat größeren Umfang und kleineren Inhalt als seine Unterbegriffe; speziell enthält ein Unterbegriff mindestens die Attribute aller seiner Oberbegriffe und ein Oberbegriff mindestens die Objekte aller seiner Unterbegriffe.

Satz 1 (Hauptsatz der Begriffsanalyse [16])

Sei $K = (\mathcal{O}, \mathcal{A}, \mathcal{R})$ ein Kontext, dann ist $\mathcal{B}(\mathcal{O}, \mathcal{A}, \mathcal{R})$ ein vollständiger Verband, der Begriffsverband von K . Infimum und Supremum sind gegeben durch:

$$\begin{aligned}\bigwedge_{i \in I} (O_i, A_i) &= \left(\bigcap_{i \in I} O_i, \omega\left(\bigcup_{i \in I} A_i\right) \right) \\ \bigvee_{i \in I} (O_i, A_i) &= \left(\alpha\left(\omega\left(\bigcup_{i \in I} O_i\right)\right), \bigcap_{i \in I} A_i \right)\end{aligned}$$

Dieser Satz besagt, daß mehrere Begriffe einen eindeutig bestimmten kleinsten Oberbegriff und größten

schritt davon ausgegangen werden, daß die Anfrage diese Attribute bereits enthält.

Abbildung 5 zeigt graphisch, wie das Ergebnis einer Anfrage in einem Begriffsverband bestimmt wird. Für die Attribute a_1 und a_2 einer Anfrage werden die zugehörigen größten Begriffe bestimmt – sie sind mit a_1 und a_2 beschriftet. Von diesen Begriffen wird das Infimum gebildet; in der Grafik geschieht dies durch Verfolgen der stärker gezeichneten Linien zum ersten gemeinsamen Begriff (Pfeil). Dieser Begriff enthält alle die Anfrage erfüllende Objekte o_1, o_2 und o_3 . Dies bedeutet auch, daß diese Objekte durch Begriffe des induzierten Unterverbandes (grau) eingeführt werden.

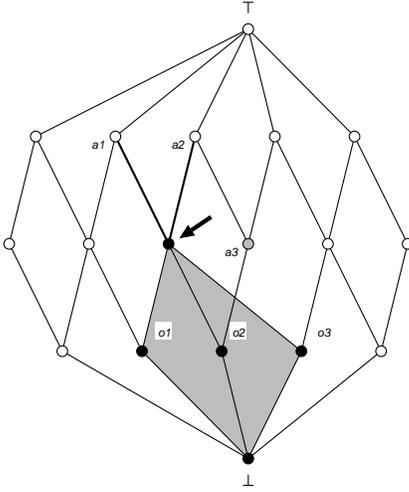


Abbildung 5: Suche im Begriffsverband

Aus welchen *sinnvollen* Attributen kann ein Benutzer auswählen, um seine Anfrage zu spezialisieren, so daß das Ergebnis der neuen Anfrage nicht leer ist? Die Objekte des Begriffes zu einer Anfrage enthalten Attribute, die nicht zu diesem Begriff gehören. In Abbildung 5 ist a_3 Attribut des Objektes o_2 , aber a_3 ist nicht Attribut des bisherigen Infimums (Pfeil), da $\mu(a_3)$ kein Oberbegriff des Infimums ist. Alle Attribute, die nicht Attribut des Infimums sind, aber Attribut eines Objektes des Infimums, können eine Anfrage weiter spezialisieren. Das Attribut a_3 kann die Anfrage aus a_1 und a_2 also spezialisieren: das neue Ergebnis besteht nur aus dem Objekt o_2 .

Definition 6 (Sinnvolle Attribute) Sei $A \subseteq \mathcal{A}$ eine Anfrage in $\mathcal{B}(\mathcal{O}, \mathcal{A}, \mathcal{R})$. Dann bezeichnet $\langle\langle A \rangle\rangle$ die Menge der sinnvollen Attribute für die Spezialisierung von A :

$$\langle\langle A \rangle\rangle := \{a \in \mathcal{A} \mid \emptyset \subset [A \cup \{a\}] \subset [A]\}$$

Die Definition von $\langle\langle A \rangle\rangle$ zu einer Anfrage A ist also so gewählt, daß jedes Attribut aus $\langle\langle A \rangle\rangle$ A derart erweitert, daß das Ergebnis der neuen Anfrage eine echte Spezialisierung der alten, aber nicht leer ist. Die

Attribute der Objekte von $[A]$, die noch nicht durch die Anfrage A selbst impliziert werden, sind die Kandidaten für die Spezialisierung der Anfrage: sinnvolle Attribute werden durch Begriffe eingeführt, die mit dem Infimum zur Anfrage unvergleichlich sind oder ein echter Unterbegriff des Infimums sind und deren Objekte zum Ergebnis der Anfrage beitragen.

Satz 4 Sei $A \subseteq \mathcal{A}$ eine Anfrage in $\mathcal{B}(\mathcal{O}, \mathcal{A}, \mathcal{R})$. Dann gilt für die Attribute $\langle\langle A \rangle\rangle$ zur Spezialisierung von A :

$$\langle\langle A \rangle\rangle = \left(\bigcup_{o \in [A]} \omega(o) \right) \setminus \pi_a \left(\bigwedge_{a \in A} \mu(a) \right)$$

Der Benutzer kann nach jeder Anfrage A ein Attribut $a_1 \in \langle\langle A \rangle\rangle$ auswählen, um seine Anfrage weiter zu spezialisieren. Unabhängig von dem gewählten Attribut ist garantiert, daß das neue Ergebnis nicht leer ist. Zur algorithmischen Bestimmung von $\langle\langle A \rangle\rangle$ werden die Attribute der Objekte aus $[A]$ vereinigt, wobei die Attribute ausgeschlossen werden, die bereits in dem durch die Anfrage A induzierten Begriff enthalten sind. Da $\omega(o) = \pi_a(\sigma(o))$ für $o \in \mathcal{O}$ gilt, kann $\omega(o)$ ebenfalls leicht aus dem Begriffsverband gewonnen werden, so daß eine Implementierung der Suche nicht zusätzlich den Kontext benötigt.

Das Ergebnis $[A_1]$ der neuen Anfrage $A_1 = A \cup \{a_1\}$ wird schließlich inkrementell durch Berechnung eines neuen Infimums $[A_1] = \pi_o(b_{[A]} \wedge \mu(a_1))$ bestimmt. Dabei ist $b_{[A]} = \bigwedge_{a \in A} \mu(a)$ bereits durch die vorhergehende Berechnung von $[A]$ bekannt. Anschließend kann der nächste Zyklus zur Spezialisierung der Anfrage durchlaufen werden.

Die Berechnung des Ergebnisses $[A]$ einer Anfrage und die exakt verbleibenden Möglichkeiten $\langle\langle A \rangle\rangle$ zur Spezialisierung der Anfrage lassen sich also elegant auf die Begriffe des Begriffsverbandes zurückführen. Der Begriffsverband kann als eine übersetzte Form aller Wahlmöglichkeiten des Benutzers aufgefaßt werden. Der Verband enthält gegenüber dem zugrundeliegenden Kontext keine wesentlich neuen Informationen, erlaubt aber eine effizientere Bearbeitung von Anfragen. Wie bei der Übersetzung von Programmen ist dies von Vorteil, wenn die Ausgangsdaten (Quelltext/Kontext) selten geändert, aber oft abgearbeitet (ausgeführt, durchsucht) werden.

5 Prototyp

Das beschriebene Verfahren wurde als Prototyp zur Verwaltung einer Unix Online-Dokumentation implementiert. Insgesamt 1658 Dokumente wurden mit mehreren Attributen aus einer Menge von 92 Attributen versehen – der zugehörige Begriffsverband enthält 714 Begriffe. Abbildung 6 zeigt die Schnittstelle des Prototyps zum Benutzer, nachdem die Attribute *change* und *file* ausgewählt wurden.

In der mittleren Liste werden die ausgewählten Attribute A angezeigt und rechts die noch wählbaren At-

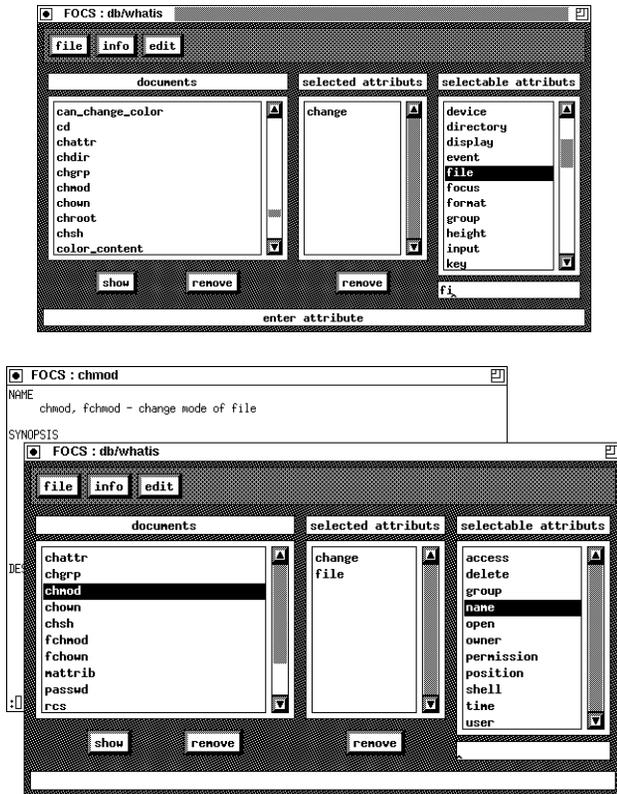


Abbildung 6: Zwei Schritte einer Suche mit dem Prototyp

tribute, also die Menge $\langle\langle A \rangle\rangle$. Die linke Liste enthält die durch die Attribute A ausgewählten Dokumente $\llbracket A \rrbracket$; durch Anklicken kann das entsprechende Dokument angesehen werden, wie es im Hintergrund für das Kommando *chmod(1)* zu sehen ist. Ein rechts ausgewähltes Attribut wird in die mittlere Liste übernommen und führt zu einer Verringerung der ausgewählten Dokumente und der noch wählbaren Attribute. Die Auswahl eines Attributs geschieht entweder direkt mit der Maus oder durch schriftliche Eingabe, wobei die Angabe eines Präfixes schon genügt, um den Listen-Cursor an das entsprechende Attribut springen zu lassen. Ein bereits ausgewähltes Attribut kann durch nochmaliges Auswählen mit der Maus zurückgenommen werden, woraufhin die Darstellung entsprechend aktualisiert wird. Dabei kann eine völlig neue Situation entstehen, wenn die Menge der dann noch ausgewählten Schlüsselwörter neu ist. Die Reaktion des Prototyps auf die Aktionen des Benutzers ist verzögerungsfrei – der Benutzer kann sehr schnell durch den Datenbestand navigieren. Durch die Präsentation seiner verbleibenden Wahlmöglichkeiten wird er dabei stark kontextsensitiv unterstützt.

Das Berechnen des Begriffsverbandes aus einer Indexierung von Objekten mit Attributen benötigt 150 Sekunden auf einer SPARCstation ELC. Diese Berechnung ist allerdings nur nötig, wenn sich die Ausgangsdaten geändert haben, ansonsten wird der vor-

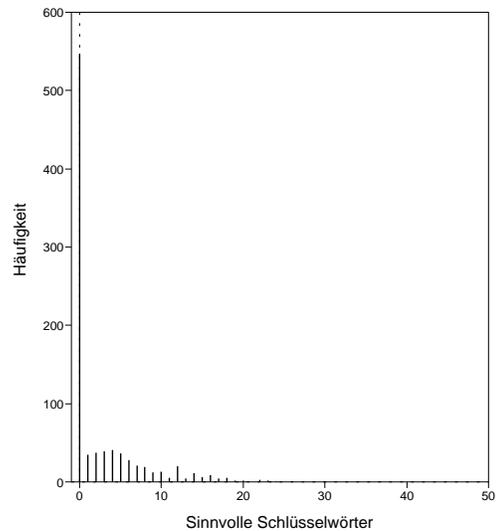
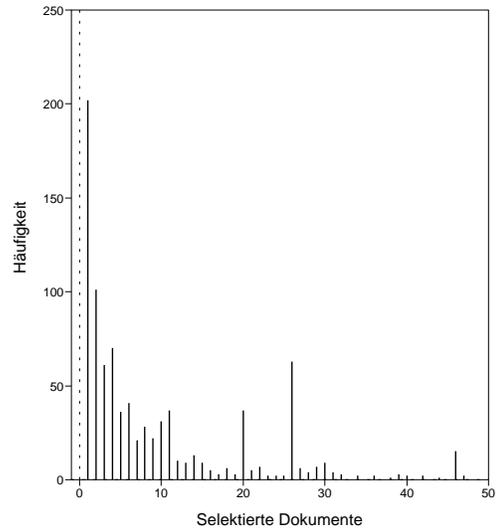


Abbildung 7: Verteilung von Dokumenten und Schlüsselwörtern nach 2 Suchschritten

berechnete Begriffsverband benutzt.

Für die mit dem Prototyp verwalteten Dokumente wurde untersucht, wie sich die Zahl der selektierten Dokumente $\llbracket A \rrbracket$ und der dann noch wählbaren Schlüsselwörter $\langle\langle A \rangle\rangle$ nach zwei Suchschritten ($|A| = 2$) verteilen. Von den $\binom{92}{2} = 4186$ theoretisch möglichen Kombinationen von Schlüsselwörtern sind 903 tatsächlich erreichbar. Für diese wurden jeweils die Zahl der selektierten Dokumente und im dritten Schritt wählbaren Schlüsselwörter bestimmt. Abbildung 7 zeigt die beiden Verteilungen: oben für die Dokumente, unten für die Schlüsselwörter.

Bereits nach zwei Suchschritten ist die Zahl der selektierten Dokumente in 52% aller Fälle von anfänglich 1658 auf weniger als 6 gefallen. Gleichzeitig ist die Zahl der noch zusätzlich wählbaren Schlüsselwörter in 68% aller Fälle kleiner als 2. Dies zeigt, wie das Verfahren den Suchraum eingrenzt und damit die Navi-

gation zielgerichtet und schnell wird.

Eine empirische Untersuchung von *Precision* und *Recall* beim Einsatz der begriffsbasierten Suche zur Verwaltung von Komponenten oder ihrer Dokumentation steht noch aus. Es ist allerdings zu bemerken, daß das Ergebnis $[A]$ einer Anfrage A exakt und vollständig ist; *Precision* und *Recall* hängen also nur von der Qualität der Indexierung ab.

6 Berechnung von Begriffsverbänden

Unglücklicherweise kann die Zahl der Begriffe in einem Begriffsverband exponentiell mit der Zahl der Objekte und Attribute einer Sammlung wachsen. Für die praktische Anwendbarkeit der Suche muß deshalb untersucht werden, ob dieses Verhalten typisch ist, oder ob es vielmehr eine pathologische Ausnahme ist. Dazu wurden formale Kontexte zufällig generiert und die zugehörigen Begriffsverbände berechnet. Die erzeugten Kontexte wiesen eine variable Anzahl von Objekten auf, die mit je fünf Attributen versehen waren, um eine facettrierte Klassifikation mit 5 Facetten nachzubilden [11]. Die fünf Attribute wurden dabei zufällig aus einer variabel großen Menge von Attributen ausgewählt. Abbildung 8 zeigt die Größe der Verbände und den Aufwand zu ihrer Berechnung durch den Algorithmus von Ganter [3] auf einer SPARCstation ELC. Die Zeitkomplexität des Algorithmus ist $O(n \cdot |A| \cdot |\mathcal{O}|^2)$, wobei n die Zahl der Begriffe von $(\mathcal{O}, \mathcal{A}, \mathcal{R})$ ist.

Der Zeitaufwand bei den untersuchten Kontexten ist im Wesentlichen durch die Zahl der Objekte, also die Größe einer Sammlung, bestimmt. Die Größe des Wortschatzes zur Indexierung hat dagegen geringeren Einfluß. Die gemessenen Rechenzeiten liegen zwischen einigen Sekunden für mehrere hundert Objekte und steigen bis zu mehreren Minuten für Sammlungen von fast 2000 Objekten an. Die Größe der entstandenen Begriffsverbände wird ebenfalls durch die Zahl der Objekte dominiert, allerdings entstehen auffallend große Verbände bei vielen Objekten mit vergleichsweise wenigen Attributen zur Auswahl. Die Zahl der Begriffe liegt typischerweise bei einigen tausend – das präsentierte Verfahren ist damit für kleine und mittlere Sammlungen mit mehreren hundert bis einigen tausend Objekten ideal. Damit ergibt sich im typischen Fall eine polynomiale Zeit-Komplexität für die Bestimmung aller Begriffe eines Kontextes.

7 Andere Ansätze

Die begriffsbasierte Suche ist ein Information-Retrieval-Ansatz, im Gegensatz zu den formallogischen Verfahren [8, 12, 13, 2]. Information-Retrieval-Ansätze werden von Maarek et al. [7] und Wood et al. [17] wegen ihrer oft leichteren Bedienung und Durchschaubarkeit bevorzugt. Die facettrierte Klassifikation [11] ist ein bekannter Information-Retrieval-Ansatz, dessen Indexierung aus einem Vek-

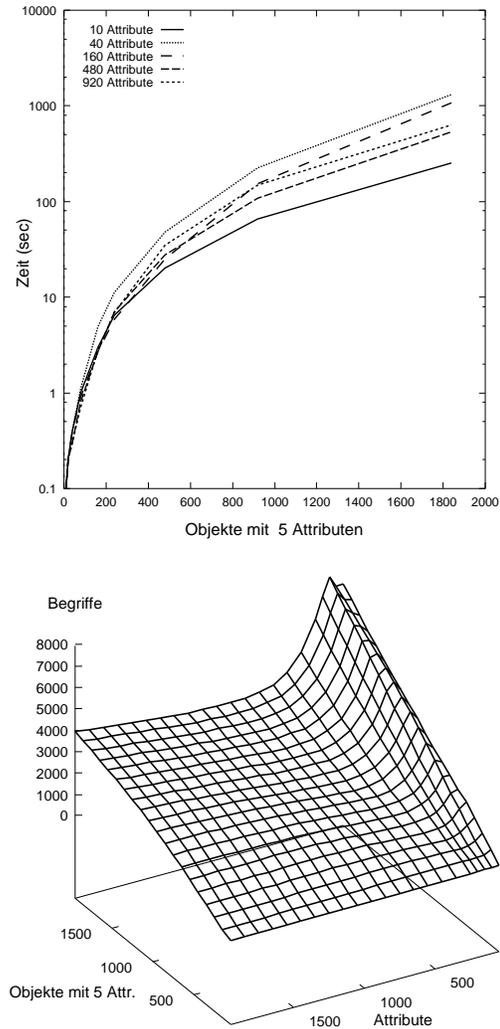


Abbildung 8: Zeitaufwand zur Berechnung von Begriffsverbänden und ihre Größe

tor von Schlüsselwörtern für jede Komponente besteht. Damit ist bei der Suche die Reihenfolge der Schlüsselwörter signifikant und sogar fest vorgegeben; zusätzlich sind Platzhalter nur sehr eingeschränkt zugelassen. Der Benutzer kann sich einer gesuchten Komponente deshalb nur über *einen* Pfad nähern, wobei er nicht kontextsensitiv unterstützt wird. AIRS [9] ist eine Verallgemeinerung dieses Ansatzes; durch manuell spezifizierte Beziehungen zwischen Komponenten wird die Navigation erleichtert, allerdings auf Kosten der Wartbarkeit der Sammlung.

GURU [7] verwaltet ebenfalls die Dokumentation eines Unix-Systems. Es indexiert den Inhalt von Dokumenten und Anfragen statistisch und erlaubt so natürlichsprachliche Anfragen. Zusätzlich stehen Dokumente in einer hierarchischen Beziehung, durch die sie ebenfalls inspiziert werden können. Die natürlichsprachliche Anfrage ist zunächst für einen Benutzer angenehm, allerdings bleibt das Ergebnis undurchschaubar, da das System nicht den Sinn einer Anfrage

verstehen, sondern statistisch arbeiten. Im Gegensatz dazu ist die nicht-operationale Semantik einer Anfrage beim begriffsbasierten Ansatz dem Benutzer einsehbar.

Godin et al. [5] schlagen ebenfalls Begriffsverbände als Navigationsgrundlage vor. Bei ihnen erhält ein Benutzer aber nur eine Teilmenge aller sinnvollen Attribute einer Anfrage zur Unterstützung – nur die, die durch Begriffe unmittelbar kleiner als der Ergebnisbegriff eingeführt werden. Sie schlagen eine Benutzeroberfläche vor, die dem Benutzer den Begriffsverband darstellt. Der vorgestellte Prototyp benutzt dagegen den Verband nur als für den Benutzer unsichtbare unterliegende Datenstruktur.

8 Ausblick

Die begriffsbasierte Suche erscheint nicht nur für die Verwaltung von Softwarekomponenten attraktiv. Damit nicht für jede weitere Anwendung die Suche vollständig neu implementiert werden muß, entsteht zur Zeit ein anwendungsunabhängiger Kern, der die zentralen Algorithmen der begriffsbasierten Suche implementiert. Dazu wird die Sprache TCL/TK [10] um begriffsbasierte Kommandos erweitert, die dann zusammen mit dem graphischen Toolkit Tk einen Ausgangspunkt für weitere begriffsbasierte Anwendungen bietet.

Zur Zeit muß bei einer Veränderung der Sammlung anschließend der gesamte Begriffsverband neu berechnet werden. Durch den Einsatz inkrementeller Algorithmen [4] zur Berechnung von Begriffsverbänden kann der Aufwand nach einer Änderung deutlich gesenkt werden, womit das Verfahren weiter an Attraktivität gewinnt.

9 Ergebnis

Mit Schlüsselwörtern indexierte Komponenten bilden einen formalen Kontext, aus dem durch formale Begriffsanalyse ein Begriffsverband berechnet werden kann. Dieser erlaubt dann eine effiziente, inkrementelle Suche von Komponenten zu implementieren, die den Benutzer durch die vollständige Präsentation seiner Wahlmöglichkeiten unterstützt. In Verbindung mit der Effizienz des Verfahrens ergibt sich eine schnelle und sichere Navigation durch eine Sammlung von Komponenten. Da zwischen Komponenten keine Konsistenz erfordernden expliziten Verbindungen spezifiziert werden, ist die Wartung einer Sammlung einfach. Die nach einer Veränderung der Sammlung nötige Neuberechnung des Begriffsverbandes hat typischerweise polynomialen Aufwand, wie Experimente gezeigt haben. Eine prototypische Implementierung als Teil von NORA [14, 6] demonstriert die praktische Anwendung und die Eignung des Verfahrens zur Verwaltung einer Unix Online-Dokumentation. Da das Suchverfahren unabhängig vom Inhalt der verwalteten Komponenten ist, können auch inhomogene und mul-

timediale Sammlungen begriffsbasiert verwaltet werden.

Danksagung: Peter Dettmer implementierte den Prototyp in C++. Gregor Snelting, Franz-Josef Grosch und Bernd Fischer trugen durch Diskussionen und Anmerkungen viel zum Verständnis des Verfahrens bei.

Literatur

- [1] Davey, B. A. & Priestley, H. A. (1990), *Introduction to Lattices and Order*, 2nd ed., Cambridge University Press, Cambridge, GB, Formal Concept Analysis, 221–236.
- [2] Fischer, B., Kievernagel, M. & Struckmann, W. (1995) VCR: A VDM-based software component retrieval tool. In Proc. ICSE-17 Workshop *Formal Methods Appl. in SW Eng. Practice*, Seattle, 30–38.
- [3] Ganter, B. (1986), Algorithmen zur Formalen Begriffsanalyse, in B. Ganter, R. Wille & K. E. Wolff, Hrsg., 'Beiträge zur Begriffsanalyse', BI Wissenschaftsverlag, Mannheim, pp. 241–254.
- [4] Godin R. & Missaoui R. (1994), An incremental concept formation approach for learning from databases. *Theoretical Computer Science*, 133 (1994), 387–419.
- [5] Godin, R., Gecsei, J. & Pichet, C. (1989), Design of a Browsing Interface for Information Retrieval, in N. J. Belkin & C. J. van Rijsbergen, Eds, '12th Int. SIGIR Conference', Cambridge, Massachusetts, 32–39.
- [6] Krone, M. & Snelting, G. (1994). On the inference of configuration structures from source code. In *Proc. 16th International Conference on Software Engineering*, Sorrento, Italy, 49–57.
- [7] Maarek, Y. S., Berry, D. M. & Kaiser, G. E. (1991), 'An Information Retrieval Approach For Automatically Constructing Software Libraries', *IEEE Transactions on Software Engineering* **SE-17**(8), 800–813.
- [8] Moorman Zaremski, A. & Wing, J. M. (1993), Signature Matching: A Key to Reuse, in D. Notkin, Ed., 'Proc. of the 1st ACM SIGSOFT Symposium on the Foundation of Software Engineering', Los Angeles, 182–190.
- [9] Ostertag, E., Hendler, J., Prieto-Díaz, R. & Braun, C. (1992), 'Computing Similarity in a Reuse Library System: An AI-Based Approach', *ACM Transactions on Programming Languages and Systems* **1**(3), 205–228.

- [10] Ousterhout, J. K., (1991), An X11 Toolkit Based on the TCL Language, In *Proc. 1991 Winter USENIX Conf.*, 105–115.
- [11] Prieto-Díaz, R. (1991), ‘Implementing Faceted Classification for Software Reuse’, *Journal of the ACM* **34**(5), 89–97.
- [12] Rittri, M. (1991), ‘Using types as search keys in function libraries’, *Journal of Functional Programming* **1**(1), 71–89.
- [13] Rollins, E. J. & Wing, J. M. (1991), Specifications as Search Keys for Software Libraries, in K. Furukawa, ed., ‘Logic Programming: Proc. of the 8th Int. Conference’, Paris, 173–187.
- [14] Snelting G., Fischer, B., Grosch, F.-J., Kievernagel M. & Zeller, A. (1994). Die inferenzbasierte Softwareentwicklungsumgebung NORA. *Informatik—Forschung und Entwicklung*, 9(3):116–131
- [15] Sun Microsystems. *SunOS Reference Manual*. Sun Microsystems, Inc., SunOS 4.1.1, 1990.
- [16] Wille, R. (1982), Restructuring lattice theory: An approach based on hierarchies of concepts. in I. Rival, ed., *Ordered Sets*, Reidel, 445–470.
- [17] Wood, M. & Sommerville, I. (1988), ‘An Information Retrieval System for Software Components’, *SIGIR Forum* **22**(3), 11–25.

A Beweise

A.1 Satz 2

Beweis: Davey, B. A. & Priestley, H. A. (1990), *Introduction to Lattices and Order*, 2nd ed., Cambridge University Press, Cambridge, 221–236.

A.2 Satz 3

$$\begin{aligned}
\pi_o\left(\bigwedge_{a \in A} \mu(a)\right) &= \pi_o\left(\bigwedge_{a \in A} (\alpha(a), \omega(\alpha(a)))\right) \\
&= \bigcap_{a \in A} \alpha(a) = \bigcap_{a \in A} \{o \in \mathcal{O} \mid (o, a) \in \mathcal{R}\} \\
&= \{o \in \mathcal{O} \mid \forall a \in A : (o, a) \in \mathcal{R}\} \\
&= \{o \in \mathcal{O} \mid \omega(o) \supseteq A\} = \llbracket A \rrbracket
\end{aligned}$$

A.3 Satz 4

1. Sei $a_1 \in \llbracket A \rrbracket$. Zu zeigen ist:

- $a_1 \in \bigcup_{o \in \llbracket A \rrbracket} \omega(o)$
- $a_1 \notin \pi_a(\bigwedge_{a \in A} \mu(a))$

(a) Sei $a_1 \in \llbracket A \rrbracket$. Also existiert $o_1 \in \llbracket A \cup \{a_1\} \rrbracket \subset \llbracket A \rrbracket$ mit $\omega(o_1) \supseteq A \cup \{a_1\} \supseteq \{a_1\}$. Dies ist äquivalent zu $a_1 \in \omega(o_1)$ und daraus folgt $a_1 \in \bigcup_{o \in \llbracket A \rrbracket} \omega(o)$.

(b) Annahme: $a_1 \in \pi_a(\bigwedge_{a \in A} \mu(a))$. Daraus folgt $\mu(a_1) \geq \bigwedge_{a \in A} \mu(a)$ und damit $\bigwedge_{a \in A \cup \{a_1\}} \mu(a) = \bigwedge_{a \in A} \mu(a)$. Also gilt $\llbracket A \rrbracket = \llbracket A \cup \{a_1\} \rrbracket$, was im Widerspruch zur Voraussetzung $\llbracket A \cup \{a_1\} \rrbracket \subset \llbracket A \rrbracket$ steht.

2. Sei $a_1 \in \bigcup_{o \in \llbracket A \rrbracket} \omega(o)$ und $a_1 \notin \pi_a(\bigwedge_{a \in A} \mu(a))$. Zu zeigen ist:

- $\llbracket A \cup \{a_1\} \rrbracket \neq \emptyset$
- $\llbracket A \cup \{a_1\} \rrbracket \subseteq \llbracket A \rrbracket$
- $\llbracket A \rrbracket \neq \llbracket A \cup \{a_1\} \rrbracket$

(a) Aus $a_1 \in \bigcup_{o \in \llbracket A \rrbracket} \omega(o)$ folgt, daß $o_1 \in \llbracket A \rrbracket$ existiert, mit $a_1 \in \omega(o_1)$ und damit $o_1 \in \alpha(a_1)$. Betrachte nun:

$$\begin{aligned}
\llbracket A \cup \{a_1\} \rrbracket &= \bigcap_{a \in A} \alpha(a) \cap \alpha(a_1) \\
&\supseteq \{o_1\} \neq \emptyset
\end{aligned}$$

(b) Betrachte: $\llbracket A \cup \{a_1\} \rrbracket = \bigcap_{a \in A} \alpha(a) \cap \alpha(a_1) \subseteq \bigcap_{a \in A} \alpha(a) = \llbracket A \rrbracket$.

(c) Annahme: $\forall o \in \llbracket A \rrbracket : \omega(o) \cap a_1 = \emptyset$. Daraus folgt $a_1 \in \omega(\llbracket A \rrbracket) = \pi_a(\bigwedge_{a \in A} \mu(a))$. Dies ist ein Widerspruch; also gilt: $\exists o_1 \in \llbracket A \rrbracket : \omega(o_1) \cap a_1 = \emptyset$ und damit $o_1 \notin \llbracket A \cup \{a_1\} \rrbracket$. Also $\llbracket A \rrbracket \neq \llbracket A \cup \{a_1\} \rrbracket$.