

ICOT Technical Memorandum: TM-1041

TM-1041

知識ベース管理システムKappaにおける
検索の実現方式と
集合を使ったレコード読み

五十嵐 努、杉崎 元(三菱)

May, 1991

© 1991, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03)3456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

知識ベース管理システム Kappa における 検索の実現方式と集合を使ったレコード読み

基本ソフトウェア開発第一部 · 五十嵐 努 · 杉崎 元

まえがき

知識ベース管理システム Kappa (Knowledge APPLICATION oriented Advanced database management system) は逐次推論マシン P S I (Personal Sequential Inference machine) 上に構築されており、高度な知識処理のための環境を提供することを目的としている。そのためにはこれまでのデータベース管理システムでは困難であった複雑なデータ表現を可能にし、実世界の事象を無理なく表わせる必要がある。そこで Kappa ではデータモデルに非正規形モデル (nested relational model) を採用している。

非正規形モデルは属性 (項目) の構造化や繰り返しを許すが、Kappa ではさらに項 (term) の格納も許すこととこれまでにないデータ表現が可能になっている。

現在では遺伝子配列や、辞書など、従来のデータベースで取り扱うのは困難であった情報をも格納できる知識ベース管理システムとして国内、及び海外の大学、研究機関などで使用されている。

本論文では第 I 部で非正規形モデル採用に伴い、高度な機能を持つことになった検索機能の実現方式について述べ、第 II 部で検索結果である集合を用いたレコードの読みについて、その処理方式を述べる。そして最後に検索機能と集合の今後の課題について考察する。

第 I 部 検索の実現方式

1. 機能の概要

本章では Kappa で提供している検索機能の概要を示す。

1.1 検索の種類

Kappa が提供する検索機能は、検索対象に応じて以下の 3 種類がある。

- インデックス検索
索引である B+tree を使用して検索する。
- データ検索
B+tree を使用せず検索対象すべてのレコードを読み検索する。

● 項検索

インデックス／データ検索を組み合わせて単一化による検索を行う。

これらの 3 種類の検索はいずれも検索結果として、集合を作成する。ユーザはこの集合を使用して他の集合との演算や、レコード読みを行う（詳しくは第 II 部を参照）。

以下では、この中のインデックス検索と、データ検索について述べる。

1.2 検索条件

インデックス検索とデータ検索で用いる検索条件は共通で、検索条件内には種々のオプションを指定でき、多様な検索が 1 回のコマンド実行で行える。

検索条件の基本となるのは検索したい属性と検索値を指定した单一検索条件である。また検索条件内に論理演算子を用い、複数の单一検索条件を指定することもできる（複合検索条件）。

Kappa の検索条件文法の主な項目を Backus Normal Form (BNF) 記法で表記したものを図 I-1.1 に、使用例を図 I-1.2 に示す。

```

Condition ::= " " <検索条件>
<検索条件> ::= <条件>
| [<検索条件>] <論理演算子> <検索条件>
| "(" <検索条件> ")"
% <検索条件>と略できるのは<論理演算子>が"not"のときのみ
<条件> ::= <属性名><比較演算子><値 1>
| <属性名><等号><検索値>
| <属性名><等号>"["<検索値>","<検索値>]"
<論理演算子> ::= "or" | "and" | "dif" | "not"
<比較演算子> ::= "<" | ">"
| "<=" | ">=" | "="
<等号> ::= "="
<テーブル名> ::= アトム
<属性名> ::= アトム
<検索値> ::= <部分値>
| "[" <部分値>(" , " <部分値>)" ")
<部分値> ::= <値>
| <曖昧検索値>
<値> ::= 整数 | アトム | 文字列
<曖昧検索値> ::= <値 1>
| <ワイルドカード>
| <曖昧検索値><ワイルドカード>
| <曖昧検索値><ワイルドカード><曖昧検索値>
| <ワイルドカード><曖昧検索値>
| <ワイルドカード><曖昧検索値><ワイルドカード>
<値 1> ::= 文字列
<ワイルドカード> ::= "*" | "?" |
% ワイルドカード"*"を連続して指定することはできない

```

図 I - 1. 1 検索条件

○ 検索条件解析フェーズ



○ 検索実行フェーズ

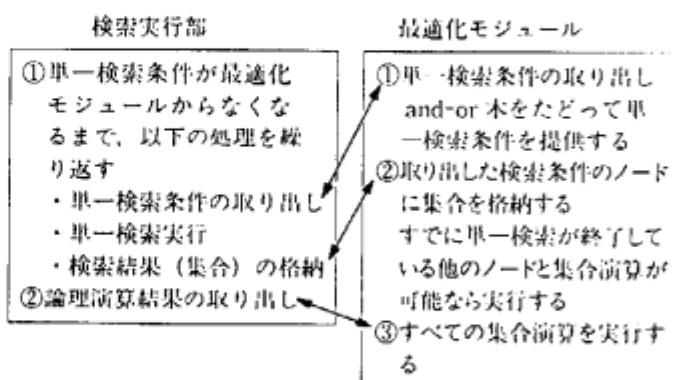
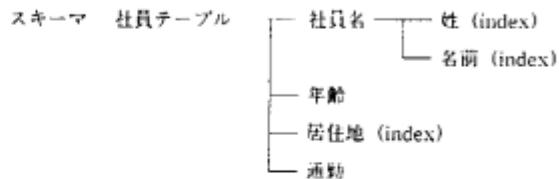


図 I - 2.1 処理フェーズと各モジュールの処理

社員テーブルに対する検索条件の例



- (1) 名前（インデックス属性、属性名：名前）の中に“美”がある人か、“太郎”という名前の人の探し。


```
Condition = "(名前 = (" * 美 * ", " 太郎"))"
: search_index (Kappa, 社員テーブル, Condition, Set)
```
- (2) 居住地（属性名：居住地）が鎌倉市で、年齢（属性名：年齢）が25歳以下の人か、通勤手段（属性名：通勤）が自動車の人を探す。


```
Condition = "((居住地 = " 鎌倉市 ") and
(年齢 <= 25)) or (通勤 = " 自動車")"
: search_data (Kappa, 社員テーブル, Condition, Set)
```

図 I - 1.2 検索条件の使用例

2. 実現方式

本章では検索処理の実現方式について述べる。

2.1 実現方式の概要

検索コマンドは、解析部によるコマンド指定の構文解析、検索最適化モジュールによる検索条件の格納、実行部による実際の検索処理という手順で実行される。

各モジュールの関係と検索処理の流れを図 I - 2.1 にまとめる。

2.2 検索条件の解析

解析部は検索条件の解析と最適化モジュールへの格納を行う。その実現方式を以下に示す。

(1) 検索条件の解析

ユーザ指定の検索条件は項として指定される。項は PSI 上ではスタックベクタ (PSI のデータ型の一つ。以下 SV とする) で表現される。入力された検索条件は SIMPOS (Sequential Inference Machine Programming & Operating System = PSI の OS) によって構文解析され、検索条件内の演算子などが誤って指定されていた場合はここでエラーになる。

構文解析された検索条件は図 I-2.2 に示した通り、それが「単一検索条件であっても、複合検索条件であっても、2または3要素の SV になっている。さらに、その SV の第2、3要素も2または3要素の SV に…」と再帰的な構造をとる。

- (1) Condition = " (attr1 = "ABC") (2) Condition = " (attr2 <= 10)
 → { -, attr1, "ABC" } → { = <, attr2, 10 }

(3) Condition = " (attr3 = [abc . def . ghi])
 → { =, attr3, [abc . def . ghi] }

(4) Condition = " (attr1 = jkl) and (attr2 = mno))
 → { and, { =, attr1, jkl }, { =, attr2, mno } }

(5) Condition = " (not (attr3 = okm))
 → { not, { =, attr3, okm } }

(6) Condition = " ((attr1 = qaz) or ((attr2 = wsx) diff (attr3 = edc))

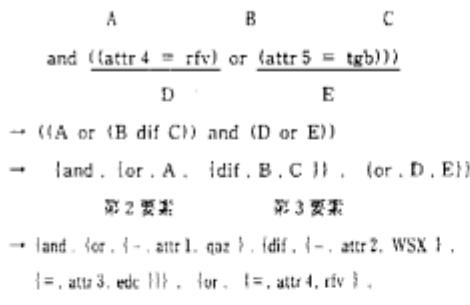


図1-2-2 検索条件の構文解析

一見すると常に要素数の同じ SV を処理するので、その SV が複合検索条件を示すのか、単一検索条件を示すのか判断が難しそうであるが、検索条件内の演算子を判定の基準とすることでその問題を解決した。

構文解析を行うと、S V の第1要素は必ず演算子になる。この演算子を基準に以下の第2、3要素の内容を判定するようにしたことにより、処理手順が明確にできた。その様子を図 1-2,3 に示す。

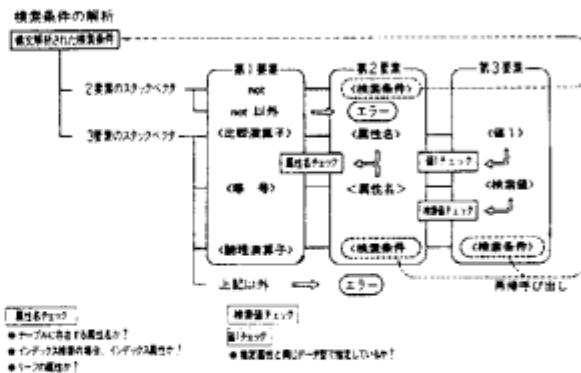


図 1-2.3 検索条件の解析

図 I-2.3中の網掛け部分は SVの要素を示す。図中の<検索条件>は図 I-2.2に示した通り複合検索条件か、単一検索条件かの判定がつかないので、再帰呼び出しで処理を行う。

解析部の行う条件文法解析処理は、項構造を解析するループが基本になっている。また〈比較演算子〉、〈等号〉を検出すると、それは單一検索条件であるとして、SV第2、3要素の内容が確定する。矩形枠の処理はその内容の検索条件文法との適合状態、スキーマとの整合性を解析する処理となる。

(2) 検索条件の格納

検索条件は図1-2.4に示す通り解析部が解析処理で検索条件を分解しながら最適化モジュールに通知し、and-or木に格納していく。and-or木の生成については次節で説明する。

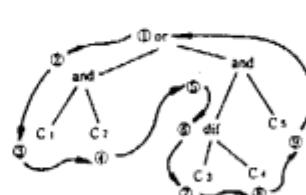
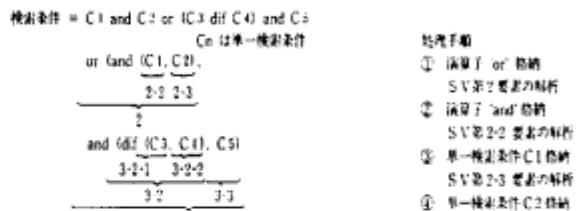


図 I-2-4 and or 杰への格納手順

2.3 and-or木

and-or木についてその実現方式を述べる。

2.3.1 and-or木を採用した背景

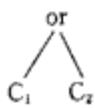
and-or木は検索条件（項構造）をヒープベクタ（P S Iのデータ型）化しKappa内部に保持できるようにしたものである。通常、このような場合SIMPOSの機能を使用し木構造を生成するが、Kappaでは以下の要求を満たすため独自に木構造を生成した。

- メモリ性能の向上
- 速度性能の向上
- 検索処理の最適化のための機能

論理的にはand-or木は検索条件内の論理演算子をノードに、単一検索条件をリーフに持つ木構造で、物理的にはヒープベクタである。and-or木は検索最適化モジュールに生成／管理され、その内部に保持されている。図I-2.5に論理的な構造図を示す。

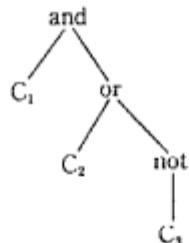
例1

条件 = C₁or C₂
★ and-or木



例2

条件 C₁and (C₂or (not C₃))
★ and-or木



図I-2.5 and-or木の論理構造

2.3.2 and-or木の生成

2.2節で述べた通り、and-or木の生成は検索条件の解析と並行して行っている。その理由は解析した検索条件をand-or木生成のために再び展開するのでは効率が悪いためである。このとき解析部から最適化モジュールには論理演算子や単一検索条件が1つ1つ通知されてくるので規則性が無いように見えるが、通知される情報には以下の規則があり、それを元に木を構築していく。

- (1) notを除く論理演算子(リンク)が通知された場合その下に2つのノードが存在する。
notの場合、その下のノードは1つである。
- (2) 単一検索条件が通知された場合は、そのノードから下に木の展開はない。
- (3) 最初に通知された情報が単一検索条件なら、その検索条件は単一検索条件である。

- (4) 論理演算子が通知されたら、次に通知されるのはそのノードの左邊で、左邊が決定した後に、右邊の通知がある。

すなわち解析部は、最適化モジュールに論理演算子と単一検索条件を解析した順番に通知する。and or木はバージングされた検索条件に木の制御情報（位置情報、最適化情報等）を付加してヒープベクタに格納したものであるから検索条件を解析した順番に木を構成していくけばよいことになる。

3. 今後の課題

Kappaの検索機能の今後の課題について考察する。

(1) 最適化の実装

現在Kappaでは検索処理の最適化をほとんど実装していない。機能検討、設計では最適化の内容まで詰めていたが、開発時間の関係で最適化処理は実装できなかつた。

and-or木など最適化機能を実装するための環境は整っているので、今後機会があれば実装していきたい。

特に知識処理システムを構築するとなれば、検索条件も複雑化し最適化が必要になってくることが予想される。

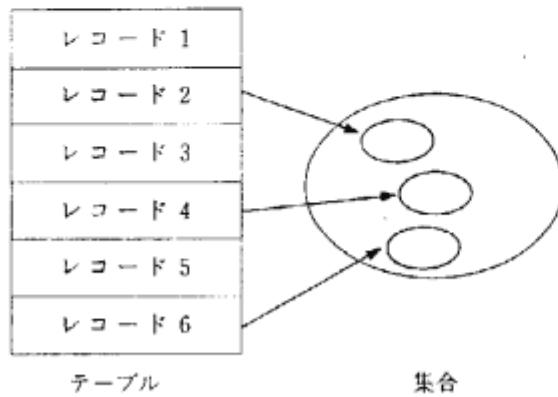
(2) 検索コマンドの統一

1章で説明した通り、Kappaの検索コマンドは現在3種類あるが、これを統一しユーザは検索対象について何も意識せずコマンドを発行できるようにしたい。Kappa内部で検索対象に応じて処理を切り分けるようにすることにより実現可能である。

第II部 集合を用いたレコード読み

1. 集合の概要

Kappa では、データ操作の効率化をおこなうためテーブルの内部表現として集合の概念を採用している。集合とは検索条件を満たすレコードの集まりであり一時テーブルとして使うことができる。



図II-1.1 テーブルと集合の関係

また集合演算（積、和、差、否定）を行うことによって更にレコードの絞り込みを行うことができる。集合では絞り込みを行ったレコードの情報として以下の情報を集合要素としている。

- (1) レコード位置情報 (R I D : レコード番号)
- (2) 属性情報
 - 属性名に対応した番号
 - オカレンス情報 (繰り返しデータの時に選択された繰り返し番号)

集合を用いたレコード読みでは(1)~(3)の情報に従ってレコードを読むことにより第I部で述べた検索条件や集合演算の結果にあったレコードの内容を読むことができる。

第II部ではこれらの情報をどのような構造で実現しているかについて述べた後、その情報を用いてレコードを読む手順について述べる。

2. 実現方式

2.1 集合の内部構造

Kappa では検索や集合演算の結果、レコード位置情報 (+属性情報) を集合要素として効率良くストリング (PSI のデータ型の 1つ) に格納している。通常は 16bit ストリングに格納しているが、処理の高速化を図るためにレ

コードの位置情報のみの場合は 32bit ストリングに格納している。

このように集合要素をストリングに保持することによって実際にレコードを読むことなく索引検索や集合演算を行うことができる。

集合の内部構造をBNF 記法で表現し、それぞれの意味を以下に述べる。

(1) 集合の定義

```
<集合> ::= "(" [<R I D要素> [, <R I D要素>] ] ")"
<R I D要素> ::= "(" "R I D" "[" , "]" ["<subRID>
    (" , "<subRID> ) " ] ")""
<R I D> ::= レコード番号
<subRID> ::= ("<positive>[;"<negative> ] ")"
<positive> ::= "[" "<path> { , , <path> } " ] "
<negative> ::= "{" "<path> { , , <path> } " } "
<path> ::= <attr_id>";"<occ>";"<occ>
<attr_id> ::= 属性情報
<occ> ::= オカレンス情報
```

(注)・各 R I D要素は R I D の昇順に並んでいる。

・各 path 情報は属性情報の昇順に並んでいる。

(2) R I D要素の意味

```
<R I D要素> ::= "(" <R I D>, "[" , "]" [<subRID>
    (" , "<subRID> ) " ] ")""
各 subRID 間は or 関係を表している。
```

<R I D要素> = (R I D, subRID₁, subRID₂, ..., subRID_n)
<=> subRID₁ or subRID₂ or ... or subRID_n

(3) subRID の意味

```
<subRID> ::= ("<positive>[;"<negative> ] ")"
<positive>, <negative> 間は dif 演算を意味している。
- subRID = (<positive>;<negative>)
    <-> positive dif negative
```

(4) positive の意味

```
<positive> ::= "{" "<path> { , , <path> } " }
<path> ::= <attr_id>";"<occ>";"<occ>
<attr_id> ::= 属性情報
各 path 間は and 条件を意味している。
<positive> = {path1, path2, ..., pathn}
<=> path1 and path2 and ... and pathn
```

(5) negative の意味

```
<negative> ::= "(" "<path> { , , <path> } " )"
<path> ::= <attr_id>";"<occ>";"<occ>
<attr_id> ::= 属性情報
各 path 間は or 条件を意味している。
<negative> = {path1, path2, ..., pathn}
```

$\leq \Rightarrow p_1 \text{ or } p_2 \text{ or } \dots \text{ or } p_m$

(6) positive と negative の関係

positive dif negative

$\leq \Rightarrow (p_1 \text{ and } p_2 \text{ and } \dots \text{ and } p_n)$

 dif $(q_1 \text{ or } q_2 \text{ or } \dots \text{ or } q_m)$

$\leq \Rightarrow (p_1 \text{ and } p_2 \text{ and } \dots \text{ and } p_n)$

 and not $(q_1 \text{ or } q_2 \text{ or } \dots \text{ or } q_m)$

$\leq \Rightarrow (p_1 \text{ and } p_2 \text{ and } \dots \text{ and } p_n)$

 and not $(q_1 \text{ and } \dots \text{ and not } (q_m))$

(7) その他

subRID = ($\langle \text{positive} \rangle ; \langle \text{negative} \rangle$) とする。

$\langle \text{positive} \rangle ; \langle \text{negative} \rangle \leq \Rightarrow \text{all dif negative}$

=not(negative) を意味する。

$\langle \text{positive} \rangle ; \langle \rangle \leq \Rightarrow \text{positive dif } \phi$

= positive を意味する。

%all := レコード全体 (繰り返しを含む)

2.2 レコード読み

次に集合を用いてレコードを読む手順を示す。

- (1) 集合から集合要素を取り出す。
- (2) レコード位置情報 (R I D) に従いレコードを読む。
- (3) レコード内容を使い属性情報とオカレンス情報を従って正規形に分解する。
- (4) 選択情報を満たすタブルを選択する。
- (5) 選択されたタブルを繰り返し構造 (非正規形) にする順番を決める。
- (6) 非正規形にまとめる。
- (7) 非正規形にまとめた情報に従い内容を返す。

このような処理が必要な理由は、非正規関係の意味を正規関係と同じとしているためである。

以上の処理手順を以下に示すテーブルを例に考え、説明する。

スキーマ情報		レコード内容	
テーブル	属性 1'	属性 1	属性 2
	a	10	
		20	
' は繰り返しをあらわす	b	30	

検索条件 = (属性 1 = a or 属性 2 = 30)

とすると検索結果の集合には以下の情報が格納される。

- レコードの位置情報 → レコード 1
- 属性情報 → 属性 1, 属性 2
- オカレンス情報 → 1, 3

以上の情報を集合の内部構造に合わせて表記すると以下の様になる。

集合の内部構造 := { (1, ({属性 1; 1}, {属性 2; 3})) }

- (1) 集合から集合要素(レコード位置情報)を取り出す。

- (2) レコード位置情報に従い内部的にレコードを読む。

レコード内容

属性 1	属性 2
a	10
	20
b	30

属性 1	属性 2
1	1
	2
2	3



- (3) レコード内容を元に

オカレンス情報を繰り返しのない正規形に分解する。

* は条件を満たすものを示している。

属性 1	属性 2
1 *	1
1 *	2
1 *	3 *
2	1
2	2
2	3 *

- (4) 選択情報を満たすタブルのみ取り出す。

正規形になっているオカレンス情報の中から選択条件を満たすものだけを取り出す。

属性 1	属性 2
1 *	1
1 *	2
1 *	3 *
2	3 *

- (5) 非正規形にするため、繰り返し構造にまとめる順番を決める。

現在 Kappa では繰り返し構造にまとめる順番をスキーマ上の属性順にしている。

この例の場合は属性 1, 属性 2 の順で繰り返し構造にする。

- (6) (5)で決定した順番に従い非正規形にまとめる。

正規形になっているオカレンス情報を属性 1, 属性 2 の順で繰り返し構造にまとめていく。

属性 1	属性 2
1 *	1
1 *	2
1 *	3 *
2	3 *

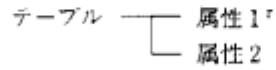
テーブルの状態



属性1	属性2
1 *	1
1 *	2
1 *	3 *
2	3 *

↓

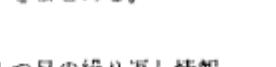
テーブルの状態



属性1	属性2
1 *	1
1 *	2
1 *	3 *
2	

↓

テーブルの状態



属性1	属性2
1 *	1
	2
1 *	3 *
2	

1つ目の繰り返し情報
2つ目の繰り返し情報

これ以上1つにまとめることは不可能なので2つのオカレンス情報を分けています。こうすることで正規形の意味を損なうことなく効率良く情報を返す事が出来ます。

(7) 处理の結果テーブルの状態は



となり、属性1、属性2共に繰り返し構造の属性となる。
また、ユーザに返すレコードは下図の繰り返し情報に従い、
下図の通り2つのレコードに分ける。

属性1	属性2
a *	10
	20
a *	30 *
b	

1つ目のレコード
2つ目のレコード

3. 今後の課題

検索によって集合をつくり、それを使ってレコードを読むという処理は、膨大な情報をもつたテーブルを利用する場合の基本的な手法になる。また、条件によって選択した集合要素を有効に利用し、さらに要素の絞り込みを行うため各種の集合演算が行われることが考えられる。以上のことから次に示す2つの課題をあげることができる。

(1) 検索の高速化

現在の検索ではレコードの選択情報となる集合要素を1要素ずつ集合に格納している。そのため集合要素を格納するため必要となる前後の処理に要する時間が要素数分必要となるため、処理効率は良いとはいえない。そこで集合要素をバッファリングし、一括して格納することにより集合要素を格納するための前後処理の発生回数を減らし、検索全体の処理時間を短くすることが可能である。

(2) 集合要素のチェック機能

集合を用いて処理をしている際に障害が発生した場合、集合に格納されている要素が正当なものかチェックする機能がないため障害の解析に必要以上の時間が取られる。そこで、集合要素のチェック機能を実装し、障害発生時にチェック機能を使い集合要素の確認をおこなうことにより集合要素の正当性を効率良くチェックし、障害解析に要する時間を短くすることが出来る。

むすび

現在 Kappa は I C O T (Institute for new generation Computer Technology)にツールとして登録されており、非正規形モデルのもつ自由なデータ表現という長所を生かし I C O T 内部の研究室だけでなく I C O T 参加各社の研究部門や国外の研究施設でも使用されている。

しかし I , II 部の今後の課題であげたように Kappa にはまだ実装できていない機能や高速化の余地がある。今後はこれらの機能を充実させ、性能向上をおこなうと共にユーザから寄せられる要望を積極的に取り入れ、より良い知識ベース管理システムとなるよう努めていきたい。

(参考文献)

- [1] 根本, 小澤, 丹羽
「知識ベース管理システム Kappa 試作システム
のデータモデル」
昭和 63 年 3 月 M T C 技術論文集
- [2] I C O T
「Kappa 説明書 1.1 版」
- [3] I C O T
「Kappa 機能設計書」
- [4] M T C
「原始コマンドマネージャ詳細設計書」
- [5] 横野
技ノート N-90-102
「Kappa での集合演算アルゴリズム」
平成 2 年 4 月

執筆者紹介



五十嵐 劳

1986 年 4 月入社。
関東学院大学工学部卒
(社) 情報処理学会会員
日英機械翻訳システム MEL-TRANJ/E (操作系) 開発
を担当。 Kappa-II, Tool Kappa の操作系コマンドマ
ンピュレータ, 原始コマンドマネージャの開発を担当。 現在
拡張関係代数の開発を担当している。



杉崎 元

1986 年 4 月入社。
湖北短期大学電子工学科卒
(社) 情報処理学会会員
Kappa-I 原始コマンド・マ
ネージャ開発を担当。
Kappa-II, Tool Kappa の
パッファ・マネージャ, 原始コ
マンド・マネージャの開発を担当。
現在拡張関係代数の開発を担
当している。