

データベース処理のデータフロー言語

北海道大学・工学部

田中謙

1. 序論

現在著者等が開発中のデータストリーム処理方式データベースマシンのためのデータフロー言語について試案を述べる。このマシンはサーチエンジン(SEE: Search Engine)とソートエンジン(SOE: Sort Engine)と名付けた2種の機能モジュールを用いて、データベース処理のボトルネックとなるデータの転送と、サーチ、ソートのデータベースにおける基本演算の処理とを重畠させ、データベース処理の高速化を図ろうとするものである。⁽¹⁾⁽²⁾⁽³⁾ 本論文では、このマシンのアーキテクチャ、SEEとSOEのアーキテクチャに関する説明は省略する。

我々は、このマシンのシステム言語を3つのレベルに分けている。最上位レベルは、更新演算を含むように拡張した関係代数言語で、第2レベルは、これより直接記述のレベルが低く、よく似ている原始関係代数言語、そして最下位レベルが、データベースマシンの各構成モジュールの機能とよく対応のとれてるデータストリーム言語である。

CODASYL DML インターフェースや QBE インターフェース等は、システム言語の最上位にある関係代数言語との論理的なマッピングの問題なので、ここでは触れないでおく。

関係代数言語とデータストリーム言語の間に原始関係代数言語を置いたのは、オブジェクトコードの最適化や、スケジューリングの問題をやさしくする限り、関係代数言語と原始関係代数言語との間のマッピングの問題として論理的に扱えたからである。本論文では、最適化およびスケジューリングの問題を避け、原始関係代数言語とデータストリーム言語の間の変換について述べる。さらに、更新演算についても本論文では触れないでおく。

2. 原始関係代数言語

原始関係代数は関係代数と同じ演算を持つが、属性を属性名で参照することができる、つまりに指標を用いる。

$R(A, B, C)$, $S(D, E, F)$, $T(G, H, I)$ に対して、関係代数式の例：これに対する原始関係代数式をいくつか列挙する。

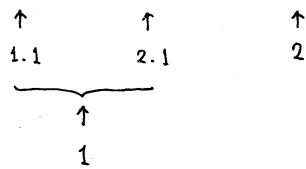
$$R[A, B] \quad [1, 2] R$$

$$R[B = F] S \quad [2=3] R S$$

$$R[C = 'v'] \quad \{3 = 'v'\} R$$

$R[A = B]$ $\{1=2\} R$ $(R[B=F]S)[C,D]$ $[3.1, 1.2][2=3]RS$

注、 $R \in 1, S \in 2$ で表わし,
 $A, B, C \in 1.1, 2.1, 3.1$ で
 表わす。

 $((R[B=F]S)[E=I]T)[C,D,H] [3.1.1, 1.2.1, 2.2][2.2=3]$ $[2=3]RS T$ 

指標 i_0, i_1, \dots, i_n に対して、 i_0 を \downarrow , i_1, \dots, i_n をタブルと呼ぶ。指標の \downarrow , \uparrow とタブルを $\downarrow(\sigma)$, $\uparrow(\sigma)$ で表わす。

通常の結合演算 $R[X=Y]S$ の他に、半結合演算

$$R[X=Y]S \triangleq \{ r \mid r \in R \text{ s.t. } \exists s \in S \quad r.X = s.Y \}$$

を考慮する。

原始関係代数言語では、プロジェクションをとらなければならず、関係中のタブルの各値を見るにはやきなとする。

3. データストリーム言語

データタイプ α のデータストリームとは、データタイプ α に属するデータの有限系列のことである。タイプ α のデータ

a_0, a_1, \dots, a_n がこの順に並んだデータストリームを、
 $\langle a_0, a_1, \dots, a_n \rangle$ で表わす。データタイプ α のデータストリームの集合を α^* で表わす。
 リストを以下のようく定義する。

- (1) a_i をタイプ α_i のデータとするとき, (a_0, a_1, \dots, a_n) はタイプ $(\alpha_0, \alpha_1, \dots, \alpha_n)$ のリストである。
- (2) b_i をタイプ β_i のデータとするリストとするとき,
 (b_0, b_1, \dots, b_n) はタイプ $(\beta_0, \beta_1, \dots, \beta_n)$ のリストである。
- (3) (1)(2) を有限回適用して得られるすべてのリストである。

タイプ α のリストのストリームの集合 α^* で表わす。データタイプ α の要素のみを用いて得られるリストのストリームのすべてを $\llbracket \alpha \rrbracket$ で表わすことにする。 i を正整数とするとき,

$$\cdot i(a_1, a_2, \dots, a_n) = \begin{cases} a_i & i \leq n \\ \text{undefined} & \text{otherwise.} \end{cases}$$

と定義する。たとえば、例えは

$$\cdot 2.3 (a, (b, c), (e, (f, g), h)) = (f, g)$$

である。

データストリーム言語で扱うデータタイプには、関係名のタプル \mathbb{R} 、データベースに現われる値のタイプ \mathbb{D} 、非負整数のタイプ \mathbb{N} と、ストリーム型 ω^* 、 \mathbb{N}^* 、 \mathbb{L}^* がある。

以下にデータストリーム言語の関数演算子を列挙する。

$! : \mathbb{R} \rightarrow \mathbb{N}^*$

関係 \mathbb{R} のタブル数を n とするとき、

$! \mathbb{R} = \langle 0, 1, 2, \dots, (n-1) \rangle$

である。

$? : \mathbb{R} \rightarrow (\mathbb{N} \rightarrow (\mathbb{N}^* \rightarrow \mathbb{D}^*))$

$? \mathbb{R} i \langle i_1, i_2, \dots, i_j \rangle$

i が \mathbb{R} の属性数より大きいとき空ストリーム、
そうでなければ \mathbb{R} の i 番目のタブルのオイ属性
の値を a_{ij} とするとき

$\langle a_{i1}, a_{i2}, \dots, a_{ij} \rangle$

とする。

$\uparrow : \mathbb{D}^* \rightarrow \mathbb{N}^* \quad (\downarrow : \mathbb{D}^* \rightarrow \mathbb{N}^*)$

$\uparrow \langle a_0, a_1, \dots, a_n \rangle \quad (\downarrow \langle a_0, a_1, \dots, a_n \rangle)$

a_0, a_1, \dots, a_n を昇順(降順)にソートした
ものを $a_{i1}, a_{i2}, \dots, a_{in+1}$ とするとき、

$\langle i_1, i_2, \dots, i_{n+1} \rangle$

を値とする。

$S_\theta : \mathbb{D} \rightarrow (\mathbb{D}^* \rightarrow \mathbb{N}^*)$

$S_\theta v \langle a_0, a_1, \dots, a_n \rangle$

$a_i \in b_i$ となる i を順に $i_0 < i_1 < \dots < i_k$ とするとき

\star

$\langle i_0, i_1, \dots, i_k \rangle$

を値とする。

$R_\theta : \mathbb{D}^* \rightarrow (\mathbb{D}^* \rightarrow \mathbb{N}^*)$

$R_\theta \langle a_0, a_1, \dots, a_n \rangle \langle b_0, b_1, \dots, b_m \rangle$

$a_i \in b_i$ となる $0 \leq i \leq \min(n, m)$ かつ $i_0 < i_1 < \dots < i_k$ となるとき

$\langle i_0, i_1, \dots, i_k \rangle$

を値とする。

$X_\theta^L : \mathbb{D}^* \rightarrow (\mathbb{D}^* \rightarrow (\mathbb{N}, \mathbb{N})^*)$

$X_\theta^R : \mathbb{D}^* \rightarrow (\mathbb{D}^* \rightarrow (\mathbb{N}, \mathbb{N})^*)$

$X_\theta^L \langle a_0, a_1, \dots, a_n \rangle \langle b_0, b_1, \dots, b_m \rangle$

$(X_\theta^R \langle a_0, a_1, \dots, a_n \rangle \langle b_0, b_1, \dots, b_m \rangle)$

$i = 0, 1, 2, \dots, m$ の順に, $a_{j_i} \in b_i$ となる最小(最大)

の j_i を求め (i, j_i) を出力する。 $\langle a_0, a_1, \dots, a_n \rangle$

は $a_0 \leq a_1 \leq \dots \leq a_n$ の場合に限る。

$A : \mathbb{L}_D \rightarrow (\mathbb{N}^* \rightarrow \mathbb{L}_D)$

$A < a_0, a_1, \dots, a_n > < i_0, i_1, \dots, i_m >$

$j = 0, 1, 2, \dots, m$ の順に, $i_j \leq n$ であれば a_{i_j} を
出力する。

$\epsilon : \mathbb{N}^* \rightarrow (\mathbb{N}^* \rightarrow (\mathbb{N}^* \rightarrow (\mathbb{N}, \mathbb{N})^*))$

被演算項である 3 → のタイプ \mathbb{N} のデータストリーム
の各々のオフセット要素を i_1, i_2, i_3 とする。これらオフセット
要素が $i_2 \leq i_3$ のとき, これらに対して, サイズト
リーム

$< (i_1, i_2), (i_1, i_2+1), \dots, (i_1, i_3) >$

を出力する。

$1 : \mathbb{L}_D \rightarrow \mathbb{L}_D$

$s \in \mathbb{L}_D$ は $\{1\}$, $1s = s$ である。

$\cdot : \mathbb{N} \rightarrow (\mathbb{L}_D \rightarrow \mathbb{L}_D)$

$i \geq 1$ の時, i は前記定義した ϵ の等しい。

$P : \mathbb{L}_D \rightarrow (\mathbb{L}_D \rightarrow \mathbb{L}_D)$

$P < a_0, a_1, \dots, a_n > < b_0, b_1, \dots, b_m >$

$= < (a_0, b_0), (a_1, b_1), \dots, (a_{\min(n,m)}, b_{\min(n,m)}) >$

$a_i, b_i \in K, c \notin$,

$$K = \lambda xyz. \quad Pxyz \quad , \quad C = \lambda xyz. \quad xz yz$$

と定義する。これらを用いた式がデータストリーム言語のプログラムである。

ここに述べたデータストリーム言語の各演算子は、著者等が開発中のデータベースマシンの各モジュールによって直接実行される。例えば↑, ↓はSOEで処理され, ?, x^L , x^R , AはSEEで処理される。⁽³⁾これらの事柄の詳細は省略する。

4. 原始関係代数言語とデータストリーム言語の変換

原始関係代数言語のある式において、指標の指す属性を持つ関係を $r(\sigma)$ で表わすこととする。原始関係代数言語の各演算子は、データストリーム言語の演算子を用いて以下のように展開できる。x, yはタクト \sqsubseteq_N の変数とする。

関係 R ! R

射影 $[\sigma_1, \sigma_2, \dots, \sigma_n]$ $\lambda x. (\ ? r(\sigma_1) \sigma_1 x, ? r(\sigma_2) \sigma_2 x, \dots, ? r(\sigma_n) \sigma_n x)$

選択 $\{\sigma \circ v\}$ $\lambda x. Ax S_v v ? r(\sigma) \sigma x$

制約 $\{\sigma \circ \sigma'\}$ $\lambda x. Ax R_\sigma C (? r(\sigma) \sigma) (? r(\sigma') \sigma') x$

半結合 $\langle \sigma \circ \sigma' \rangle$ $\lambda xy. Ax . 1 x^L A ? r(\sigma') \sigma' y \uparrow ? r(\sigma') \sigma' y ? r(\sigma) \sigma x$

結合 $[\sigma \circ \sigma']$ $\lambda xy. K (Ax . 1) (Ay . 2) C (\uparrow ? r(\sigma) \sigma x . 1) (. 1)$

$$X_0 \stackrel{?}{=} r(\sigma) \circ x \uparrow ? r(\sigma) \circ x \quad ? r(\sigma') \circ' y$$

$$X_0 \stackrel{\Delta}{=} \lambda x y. \in (.1 X_0^L x y) (.2 X_0^L x y) (.2 X_0^R x y)$$

上述の関係を用いることにより、原始関係代数言語の任意の式をデータストリーム言語に変換することが可能となる。このときのオブジェクトプログラムは、中間結果を関係として求めるこではなく最終結果を与える。

上の定義において、例えば「 $\alpha\beta\alpha'$ 」のオブジェクトコードで、母式の x と y を入れ換えたものが「 $\alpha\beta\alpha'$ 」のオブジェクトコードと考えることができます。しかし一般に両者は処理に要する時間が異なる。ここでオブジェクトコードの最適化、スケジューリングの問題が生じる。しかし、今の例で後者のオブジェクトコードを必要とするのであれば「 $\alpha\beta\alpha'$ 」 xy とする所を「 $\alpha\beta\alpha'$ 」 yx とすればよい。このことは他の演算子についても同様である。最適化とスケジューリングの問題は、データストリーム言語のレベルより一段上に引き上げられ、原始関係代数言語のレベルで議論することができる。したがって上述のように、各演算子のオブジェクトコードを固定してしまってよい。

5. 結論

本論文では、現在著者等が開発中のデータストリーム処理方式データベースマシンのデータフロー言語として、原始関係代数言語と、マシンの構成要素の機能に対応した演算子からなるデータストリーム言語の2つを示し、前者の演算子の後者の演算子を用いた一意な展開形を与えた。これにより、検索処理に限っては、関係代数とデータストリーム言語の美しい対応関係がとれた。

今後の課題としては、更新処理を含めた拡張がある。さらには、データストリーム言語の体系、そのものの単純化が課題の一つとしてある。このためには、組み合せ論理で用いられるエンゼネーターを用いた体系化等が考えられる。

参考文献

- (1) Y. Tanaka, T. Nozaka, and A. Masuyama, "Pipeline Searching and Sorting Modules as Components of a Data Flow Database Computer," Proc. IFIP Congress 80, 1980 (Tokyo), pp 427-432.
- (2) 田中 譲, "データストリーム処理方式のデータベースコンピュータ," 情報処理研究, 記号処理 12-14, 1980, pp 97-103.
- (3) 田中 譲 他, "データストリーム処理方式のデータベースコンピュータ," 情報処理全国大会, 1980, pp 493-494.