

連想記憶 (Hashing) とその応用

東大 理学部 後藤 英一
(情報科学科)

1. 任意多倍長整数用システム

今日の大多数の計算機システムはハードウェア、ソフトウェアのいずれの面からも、整数特に任意多倍長整数の取扱いは極めて不便にできている。FORTRAN を始めとする数値計算指向の言語で“整数型”変数、定数と称するものはいずれも固定された有限な範囲 (-2^m から $2^m - 1$ の場合が多く m は 16, 32, 48, 64 等計算機によって異なる) の整数しか直接には取扱えない。“整数”とは本来任意多倍長整数を意味するはずであるが、FORTRAN で、任意多倍長整数を取扱おうとすると、文法糖 (Syntactic Sugar, 分り易く、使い易いプログラム言語表記法。4 則演算は $K=I+J$, $K=I * J$ など書くのはその一例) の甘味のないサブルーチン Call の形に書く外ない。任意多倍長浮動小数点演算に関しては Wyatt たちの文法糖 FORTRAN Precompiler [4] などがあるが、“整数”につ

いてはまだそのようなものは作られた話は聞いていない。

任意の倍長整数が自由に使えるシステムとしては、数式処理システム [1, 2] がある。これは、 $Z := X + Y$ とした時 X と Y が数式、例えば有理係数の変数有理式を表わす場合にも、それに対する演算を正しく行うものである。そのごく特別な場合として“整数”を正しく取扱えなければ全く用をたさない。特に Hearn の REDUCE 2 システム [2] は計算機間の納設が容易であり、筆者たちのグループで日本の大学の代表的大型計算機への納設をほぼ完了しているので、“整数”と母関数などの数式処理を必要とされる場合にはその試用をおすすめしたい。また筆者に FORTRAN と ALGOL と LISP のいずれの言語の混用も許し、しかも“整数”は常に任意の倍長という言語処理系(ソフト) FLATS を提案した。[7] このような任意の倍長整数処理系をもつソフトウェアシステムは、“整数”演算には確かに便利ではあるが、整数の長さをソフトで実行時に検査するので、計算速度は遅い。これを速くするには、どうしてもデータ型(整数の長さ)を実行時の検査をハードウェア化する必要があり、その一つの実現法を [8] に示した。またこのようなハードウェアを持つ FLATS 計算機的设计、製作にも着手した。

以上を要約すれば、“整数”を合理的かつ効率よく処理

することは、今日の計算機のハードとソフトでは不可能というのが筆者の結論である。

2. 連想記憶 (Hashing) とその応用

アドレス (番地) を指定してそこにあるデータを読み書きする普通の記憶装置に対し、データの内容、これを鍵と呼ぶ、に対応する記憶を検索する記憶方式は連想記憶

(Associative Memory 又は Content Addressed Memory) と呼ばれている。人工知能とデータ・ベースの諸問題は連想記憶の構成法の問題に帰着する部分が多く、連想記憶には数多くの未解決な困難が残されている。しかし1個の鍵に対して高々1個の記憶内容しか対応しないことが保証されている場合 (0個の場合には、空が対応するとみれば高々1個は常に1個とも言える) には、これを高速度で検索するハード・ウェア (Time Complexity の意味での時間は鍵の個数 n に無関係で $O(1)$) 並びに Hashing と呼ばれるソフト・ウェア技法 (時間は統計的平均が $O(1)$) が知られている。Hashing のソフト・ウェアとハード・ウェア技法に関しては解説 [9] などを参照いただくとして、ここではその応用について簡単にふれてみたい。

Hashing は最初コンパイラやアセンブラなどの名前表の

検索の高速化の為に考案された。この場合の鍵は名前を表わす文字列である。筆者らはこれを拡張して、順序付 n 組、順序付有向木、(有限) 集合などのデータ構造の内部(計算機の記憶装置上) 表現の唯一化に Hashing を応用することと、この種の唯一表現法が、数式処理に必要ないくつかの基本アルゴリズムの高速化に有効であることを示した。[5, 6]

一例を示すと加算の可換性により数式の外部(紙上) 表現には $X + Y + Z = Y + X + Z = Z + Y + X \dots$ と多数の同値な表現がある。この為に従来の数式処理システムではごく簡単な2個の式の同値性の検査に99%の時間を費やしている。例えば、 n 個の変数の和である2式の同値性検査時間は、変数の単純比較法で $O(n^2)$; Sorting を利用すると、

Sorting に $O(n \log_2 n)$, 比較に $O(n)$ がかかる。ところが、 $X + Y + Z$ を集合 $\{X, Y, Z\}$ で表現すれば、 $\{X, Y, Z\} = \{Y, X, Z\} \dots$ 等、同一集合に対しては唯一の内部表現しか作られないので ([5, 6] 参照, 唯一表現の作成時間は n 元集合に対し $O(n)$) 2個の式の同値性の検査は2個のポイントの比較のみ, 時間 $O(1)$ でできる。99倍長型数の場合, その内部表現に n 進法

$$I = K_0 + K_1 r + K_2 r^2 + \dots + K_{n-1} r^{n-1}$$
 係数の順序付 n 組

$(K_0, K_1, \dots, K_{n-1})$ の唯一表現を使えば, 2整数の同値性

の検査 "if $I = J$ then …" は常に $O(1)$ ができる。
 (これは普通にやると $O(n)$ がかかる。なお、唯一表現の作成
 時間は $O(n)$)。

唯一内部表現をもつデータ構造はプログラム言語における変数名の拡張とみなして、これに値を付値することができる。(所要時間 $O(1)$)。これを利用して、筆者らは、HLISP システムに連想計算機能を組込んだ。(FLATS にも組込む予定) これは連想的に計算せよと指定した関数に対しては、一度計算した関数系は Hash 表に記録し、その検索を自動的に行うものである。(表が一杯になると忘れることもあるが、この時にだけ、再計算を行う。) この計算法を使うと、漸化式などに基づく recursive なアルゴリズムがプログラムに何ら手を入れることなく自動的に高速化できる。本文付録に再録した文献 [6] には、この計算法による $n!$ 、Fibonacci 数と二項係数の計算の高速化の例を示したが、ここに一例を付加しておく。

加法整数論的関数である制限分割数 $g_{n,m}$ (整数 $n \geq 1$ を $m \geq 1$ を越えない整数の和として表わす方法の数) は次の漸化式から計算できる。

$$g_{n,1} = g_{1,m} = 1$$

$$m \geq n > 1 \text{ ならば } g_{n,m} = 1 + g_{n,n-1}$$

$1 < m < n$ ならば $g_{n,m} = g_{n-m,m} + g_{n,m-1}$

この漸化式はそのままの形で recursive なプログラムになる。〔3〕 しかし、時間は Ramanujan-Hardy の漸近式により $O(e^{\pi\sqrt{2n/3}})$ となる。連想計算法を使うと初回は $O(n)$ と $O(n^2)$ の間、次回以降は $O(1)$ となる。組合せ論的数は漸化式による以外には計算法が知られていない場合が多く、また一連の数を連用する場合も多いので、連想計算法が有効であろう。興味のある方は、東大大型計算機センターの HLISP でこの方法の試用をお勧めしたい。

参 考 文 献

- [1] The Mathlab Group: "MACSYMA Manual",
MIT, Cambridge Mass, 1974, 1976.
- [2] A.C. Hearn: "REDUCE-2 User's Manual",
University of Utah, Salt Lake City, Utah, 1973.
- [3] D.W. Barron: "Recursive Techniques in
Programming" Macdonald, London and Elsevier,
New York (1968)
邦訳書もある。
- [4] W.T. Wyatt, P.W. Lozier and P.J. Orser:
"A Portable Extended Precision Arithmetic
Package and Library with Fortran Precompiler",
ACM. Trans. Mathematical Software vol 2
(1976) pp 209 - 231
- [5] M. Sassa and E. Goto: "A Hashing Method
for Fast Set Operations", Inf. Proc. Letters
5 (1976) 31 - 34
- [6] E. Goto and Y. Kanada: "Hashing Lemmas on
Time Complexities with Applications to Formula
Manipulation", ACM-SYMSAC 76,

Yorktown Heights N. Y (Aug. 1976)

本文付録に再録。^{*}

- (7) 後藤英一: "言語 FLATS と SP 試論",
数理科学 昭和51年 8月-9月号.
- (8) 後藤英一, 井田哲雄: "データ構造に関する一考察,
"大きな数"の高効率処理法", 情報処理学会
プログラミング・シンポジウム 昭和52年1月
- (9) 後藤英一, 井田哲雄: "ハッシングプロセッサ(解説)",
情報処理学会誌 18巻 395-401 (1977)
- (10) E. Goto, T. Ida and T. Gunji: "Parallel Hashing
Algorithms", Inf. Proc. Letters. vol 6 (1977) 8-13
- (11) T. Ida and E. Goto: "Performance of a Parallel
Hashing with Key Deletion", To appear in Proc.
IFIP Congress 77, Toronto Aug. 1977.
- (12) 後藤英一, 佐々木健昭: "計算機による数式処理の現状
(解説)" 情報処理学会誌に掲載予定.

* (編集者注) 既発表の論文の再録はなるべく遠慮する申し合せであるが、これは国際学会の予稿という非公式出版物に載ったものであり、入手の困難性をも考えて再録した。縮小しすぎて読み難い点をお許し願いたい。

Hashing Lemmas on Time Complexities with Applications to Formula Manipulation

EIICHI GOTO * ** AND YASUMASA KANADA *

* DEPARTMENT OF INFORMATION SCIENCE, UNIVERSITY OF TOKYO, TOKYO, 113, JAPAN

** INSTITUTE FOR PHYSICAL AND CHEMICAL RESEARCH, WAKOSHI, SAITAMA, 351, JAPAN

I. Introduction and Summary

Johnson[1] and Horowitz[2] applied sorting to improve time complexity of multiplication of univariate polynomials. Their results may be regarded as applications of the following LEMMA:

Sorting LEMMA. The time complexity of sorting of N items is $O(N \log_2 N)$ and that of binary search of sorted N items is $O(\log_2 N)$.

In this paper, time complexities of operation on "sets" and "ordered n-tuples" based on a hashing table search technique are presented as "Hashing LEMMAS" and are applied to formula manipulation. Unique normal forms for multivariate symbolic formulas resulting in $O(1)$ time complexity for identity checks are presented. The logarithmic factor $\log_2 N$, characteristic to sorting algorithms, is shown to all disappear from time complexities of polynomial manipulations. Actual implementation of the hashing technique is outlined and actual timing data are presented in the appendix.

II. Hashing LEMMAS on Sets and n-Tuples.

(2.0) Denotations and Conventions:

In case x represents a set or an n-tuple, $|x|$ means the number of elements.

Sets are denoted by underscored capital letter(s). Specially,

\underline{INT} is the set of (all) integers;
 $\underline{INT}_0 = \underline{INT} - \{0\}$, i.e., integers except 0;
 \underline{INT}_+ is the set of positive integers.

A BNF metaobject is denoted by embracketing a set in the underscoring notation between "<" and ">", with optional commentary un-underscored letters. This convention enables us to use both BNF and set notations. E.g., $\underline{BIT} = \{0,1\}$ and $\langle \underline{\text{Binary digit}} \rangle ::= 0|1$ are equivalent definitions, where "." means the end of a BNF definition.

In order to present algorithms precisely and concisely, Lisp with three additional data types $\langle \text{ordered n-TUPLE} \rangle$, $\langle \text{SET} \rangle$ and $\langle \text{ASSociator} \rangle$ are used in this paper. $\langle \text{INTEger} \rangle$, $\langle \text{SYMBol} \rangle$, i.e., nonnumeric atoms and $\langle \text{CONS} \rangle$, i.e., data created by Lisp functions "cons" or "list" are the three data types of ordinary Lisps. (Floating point numbers and arrays are omitted because of irrelevance to this paper.) Since the time complexity of high precision arithmetic is not the theme of this paper, the time complexities of arithmetic operations on $\langle \text{INT} \rangle$'s are assumed to be $O(1)$ for the sake of simplicity.

$\langle \text{IDentifiables} \rangle$ are defined as:

$\underline{ID} = \underline{INT} \cup \underline{SYM} \cup \underline{TUP} \cup \underline{SET} \cup \underline{ASS}$; $\langle \text{CONS} \rangle \notin \underline{ID}$.

While $\langle \text{ASS} \rangle$'s are denoted as $\langle \text{ASS} \rangle ::= (. \langle \text{ID} \rangle)$, $\langle \text{TUP} \rangle$'s and $\langle \text{SET} \rangle$'s are denoted in accordance with ordinary mathematical notations:

$\langle \text{TUP} \rangle ::= (\langle \text{ID} \rangle, \dots)$; $\langle \text{SET} \rangle ::= \{\langle \text{ID} \rangle, \dots\}$, where "..." means nonzero repetition of the same metaobject. Specially the 0-tuple $()$ and the null set $\{\}$ are regarded equivalent to NIL , i.e., $() \equiv \{\} \equiv \text{NIL}$.

$\langle \text{CONS} \rangle$ is printed as $\text{cons}[A;()] = (\text{A}\ \text{A})$ with extra blanks (A's) at both ends to discriminate them from a $\langle \text{TUP} \rangle$ printed as (A) .

(2.1) A function "tcons" appends an $\langle \text{ID} \rangle$ to a $\langle \text{TUP} \rangle$, e.g.,

$\text{tcons}[A;()] = (A)$, $\text{tcons}[(A,B);(C)] = ((A,B),C)$.

Lisp functions "car", "cdr", "cadr" etc. work on $\langle \text{TUP} \rangle$'s as on $\langle \text{Lisp LIST} \rangle$'s, e.g.,

$\text{car}[(A,B)] = A$, $\text{cdr}[(A,B)] = (B)$, $\text{cadr}[(A,B)] = B$.

$\langle \text{TUP} \rangle$'s are uniquely represented in the machine by making use of hashing for speed:

LEMMA 1. The time complexities of functions "tcons", "car" and "cdr" on $\langle \text{TUP} \rangle$ are all $O(1)$.

(2.2) A function "setup" transforms a $\langle \text{TUP} \rangle$ into a $\langle \text{SET} \rangle$ with the corresponding elements; "tupset" does the converse, e.g.,

$\text{setup}[(A,B)] = \{A,B\}$ or $\{B,A\}$;
 $\text{setup}[(A,B,B)] = \{A,B\}$ or $\{B,A\}$;
 $\text{tupset}[\{A,B\}] = (A,B)$ or (B,A) .

Specially for $t \in \underline{TUP}$, $\text{tupset}[t] = t$ (a coercion rule). Although the ordering of elements of a $\langle \text{SET} \rangle$ is irrelevant to its identity, the ordering of the elements of the $\langle \text{TUP} \rangle$ used first to define a $\langle \text{SET} \rangle$ establishes a "canonical order" among the elements of the $\langle \text{SET} \rangle$. Whenever the canonical order is needed, it can be retrieved by performing $\text{tupset}[\langle \text{SET} \rangle]$. $\langle \text{SET} \rangle$'s are represented uniquely in the machine by making use of hashing for speed:

LEMMA 2. For $t \in \underline{TUP}$, $s \in (\underline{SET} \cup \underline{TUP})$, the time complexities of $\text{setup}[t]$ and $\text{tupset}[s]$ are $O(|t|)$ and $O(1)$, respectively.

(2.3) For $x \in \underline{ID}$ the function "ass" yields an $\langle \text{ASS} \rangle$: $\text{ass}[x] = (.x^*)$. (* means actual datum represented by the variable). Conversely, for $a = \text{ass}[x] \in \underline{ASS}$ the function "key" gives the $\langle \text{ID} \rangle$, x : $\text{key}[a] = x$ and the pseudo-function $\text{assign}[a;v]$ assigns a value v , of any type, to $\langle \text{ASS} \rangle$, a . The value is $\text{assign}[a;v] = v$ and the assigned value can be retrieved as the value of the function $\text{value}[a] = v$. The initial value of an $\langle \text{ASS} \rangle$ is $()$. Similarly to Lisp, property functions are defined as $\text{put}[x;y;v] = \text{assign}[\text{ass}[\text{tup}[x;y]]]$ and $\text{remprop}[x;y] = \text{put}[x;y;()]$, where $x, y \in \underline{ID}$ and v is a datum of any type. These functions are implemented

by making use of hashing for speed:

LEMMA 3. The time complexities of "ass", "key", "assign", "value", "put", "get" and "remprop" are all $O(1)$.

Note in ordinary Lisps that properties are more restrictive: $x \in \text{SYM}$ and $y \in (\text{INT} \cup \text{SYM})$, and that in case m properties are used on a SYM the time complexity may increase as $O(m)$ due to list implementation of properties.

(2.4) For $x, y \in \text{ID}$, the predicate function $\text{eq}[x;y]$ checks the equality of x, y in accordance with the mathematical common sense. Namely, in case x and y are of different types, $\text{eq}[x;y]=()$; for $x, y \in \text{INT}$, $\text{eq}[x;y]=\text{T}$ iff x and y are numerically equal; for $x, y \in \text{SYM}$, $\text{eq}[x;y]=\text{T}$ iff x and y have the same spelling; for $x, y \in \text{ASS}$, $\text{eq}[x;y]=\text{T}$ iff $\text{key}[x]=\text{key}[y]$; for $x, y \in \text{TUP}$ or SET, $\text{eq}[x;y]=\text{T}$ iff x and y represent the same n -TUPLE or SET mathematically. E.g.,

$\text{eq}[(A,B);(B,A)]=()$, $\text{eq}[\{A,B\};\{B,A\}]=\text{T}$,
 $\text{eq}[\{A,B\};\{B,B,A\}]=\text{T}$, $\text{eq}[(.A);(.A)]=()$.

LEMMA 4. The time complexity of "eq" is $O(1)$.

Note that for the equality checking of Lisp data CONS, the time consuming function "equal" has to be used[3]. TUPLE's essentially differ from LIST's in this regard.

(2.5) Outline of an Implementation called HLISP (Hashed LISP).

Each HLISP CELL in the FSA (Free Storage Area) consists of three fields: CELL ::= [TAG, CAR field, CDR field]. Besides for GBC (Garbage Collection) marking, the TAG is used to specify the data type of the cell. Similarly to Lisp 1.5, a CONS CELL ::= [CONS, x^* , y^*] is created in the FSA as the result of $\text{cons}[x;y]$. The FSA itself is used as the (only one) hash table with the size being a prime p . For $\text{tup}[x;y]$, a hash search (insert iff absent) is made for a TUP CELL ::= [TUP, x^* , y^*], using Knuth's algorithm D[4, p521], thereby ensuring uniqueness of the resultant TUP. For $\text{ass}[x]$, a hash search is made for an ASS CELL ::= [ASS, "don't care", x^*], using Knuth's algorithm U2[4, p539]. The value of the ASS is placed in the CAR field, which is not used as the key of the hash search. A Short INT is represented as a pointer (placed in CAR or CDR field) to a non existing memory address. An n -precision INT is uniquely represented like a TUP of Short INT's (i_1, i_2, \dots, i_n) with the head cell being changed to an INT CELL ::= [INT, i_1, t], where t is a TUPLE, (i_2, \dots, i_n). A SYM CELL, corresponding to an atom header cell of Lisp 1.5, is the same as an INT CELL, except the head cell SYM CELL ::= [SYM, i_1, t] with Short INT's i_1, \dots, i_n being an unique encoding of the character string which identifies the SYM. For $\text{setup}[t]$, $t=(e_1^*, \dots, e_m^*)$, a SYS1 CELL ::= [SYS1, "don't care", "don't care"] is made first, where SYS1 is a system data tag. Secondly, a TUP $t'=(e_1^*, \dots, e_m^*)$, free of duplicating elements is made from t by using hash searches for SYS2 CELLS ::= [SYS2, "pointer to the SYS1 cell", e_i^*] for removing duplications with time complexity $O(1)$ per element of t . Thirdly, using a symmetric (in respect to permutation of arguments) hash sequence $h_i(e_1^*, \dots, e_n^*)$ $i=1, 2, 3, \dots$ (e.g., $h_1=\text{mod}(e_1^* + \dots + e_n^*, p-1)+1$, $h_i=\text{mod}(i^*h_1, p)$ with time complexity $O(n+i)$; Algorithm U2[ibid], [5]), hash search is made for a cell $s=[\text{SET}, h_1, \text{"don't care"}]$. If unsuccessful, a new

SET CELL, SET CELL ::= [SET, $h_1, [\text{SYS1}, |s|, t']$], is created. If successful, $s = \text{setup}[t]$ (redefined SET) or \neq (hash conflicting SET's) is checked by utilizing the SYS2 CELL's of t . (Time complexity $O(|t'|)$ at the most.) The hash search is resumed in the latter case. The load factor α of the FSA is limit to $\alpha \leq \alpha_M < 1$ (e.g., $\alpha_M=80\%$). When $\alpha \geq \alpha_M$ the GBC is called. A trioccupancy ("occupied" (i.e., a cell in use), "deleted" (not in use but in hash conflict) and "empty" (neither in use nor in conflict) scheme is used to reclaim the garbage CELL's without cell relocations and without using secondary storage. (A detailed analysis is given in [6]; McCarthy [7], proposed a scheme essentially the same as the present uniquely represented n -TUPLES. However, he stated a difficulty in GBC: the necessity of the use of secondary storage.) If the result of GBC does not satisfy $\alpha < \alpha_M$ (e.g., $\alpha_M=60\%$), GGBC (Grand GBC; more details are given in IV) is called. If $\alpha < \alpha_M$ is still not satisfied the job is terminated because of insufficient storage. Note that the condition $\alpha_M < \alpha_M < 1$ ensures the time complexities as claimed in LEMMAS 1-4. If $\alpha_M = \alpha_M = 1$ were used, the FSA would be usable up to the very last one cell, but the LEMMAS would not be valid.

III. Application to Formula Manipulation.

Let IP be the set of polynomials with integer coefficients and positive integer exponents.

(3.1) The Sum of Product Normal Form.

Polynomials of IP can be expressed as sum of products (terms), e.g.,

$$p_1 = 2UV^2 + 3X^3Y^4, \quad p_2 = 3Y^4X^3 + VUV + UV^2.$$

These expressions represent the same polynomial, and they can be faithfully represented in terms of TUP's as follows:

SP* form ::= ((TERM ID), COEFFICIENT), ..., and
TERM ID ::= ((VARIABLE ID), EXPONENT), ..., where COEF \in INT0, VAR ID \in SYM and EXP \in INT+. E.g.,
 $\text{sp}^*(p_1) = (((V,2), (U,1)), 2), (((X,3), (Y,4)), 3)$
 $\text{sp}^*(p_2) = (((Y,4), (X,3)), 3), (((V,1), (U,1), (V,1)), 1),$
 $((U,1), (V,2)), 1)$.

These SP* forms can be transformed into a unique S₂ normal form in the following way (a program is given later): (1) Combine duplicating VAR ID's in a TERM ID as in $VUV=V^2U$. (2) Absorb the commutative nature of multiplications into a SET: TERM ID ::= {(VAR ID), EXP}, ..., E.g., $V^2U=UV^2$ is absorbed as $\{(V,2), (U,1)\} = \{(U,1), (V,2)\}$. (3) Combine duplicating TERM ID's as in $V^2U+V^2U=2V^2U$. (4) Absorb the commutative nature of additions into a SET: SP ::= {(TERM ID), COEF}, ..., E.g.,
 $\text{sp}[p_1] = \text{sp}[p_2] = \{(\{(V,2), (U,1)\}, 2), (\{(X,3), (Y,4)\}, 3)\}$.

We now define two data structures, in order to formalize the definition of the SP form: A CLUB is a SET of 2-TUPLE's of ID's (informally, CLUB ::= {..., (m_i, g_i), ...}) such that all of the first element, to be called the (club-) "member", of the 2-TUPLE's are distinct ($m_i \neq m_j$ for $i \neq j$). The second elements (g_i 's) of the 2-TUPLE's are called "grade"s. A MULTISET is a special CLUB of which the grades are restricted to positive integers. (This agrees with the "multiset" of Knuth[4] by regarding the "multiplicity" as the grade.) Thus, we can now state: "An SP is a CLUB of TERM ID's with non-zero integer grades, called COEF"; a TERM ID is a MULTISET of SYM's, called VAR ID's; specially, for the null and constant polynomials,

$$\text{sp}(0) = \{\}, \quad \text{sp}(n) = \{(\{\}, n)\}, \quad \text{where } n \in \text{INT0}."$$

Since the SP form obviously represents IP polynomials uniquely, i.e., for $p, q \in \underline{IP}$,
 $sp(p) = sp(q)$ (set equality) iff $p = q \equiv 0$,
 by LEMMA 4 we obtain:

PROPOSITION 1. Given two IP polynomials in the SP form, the time complexity for identity checking of the two is $O(1)$.

(3.2) Polynomial Manipulation in The SP Form:
 A Property Adding Auxiliary Function:

```
addprop[g;x;v;r] = prog[[y];y:=get[g;x];
  [null[y] → prog2[put[g;x;v];r:=tcons[x;r]];
  T → put[g;x;v+y]];return[r] ].
```

Given $g, x \in \underline{ID}$, $v \in \underline{INT}$ and $r \in \underline{TUP}$, if the G-property (i.e., the value of $get[g;x]$) is $()$, "addprop" puts v on the property and appends x to r in the result, otherwise, v is added into the property. By LEMMAS 1 and 3, the time complexity is $O(1)$. Similarly, we define:

```
subprop[g;x;v;r]=addprop[g;x;-v;r].
```

A Property into Club-Grade Function:

```
clubprop0[g;r]=prog[[c;y;w];w:=r;
  A [null[w] → return[setup[c]];y:=get[g;car[w]];
  [y≠0 → c:=tcons[tcons[car[w];tcons[y;()]];c] ;
  remprop[g;car[w]];w:=cdr[w];go[A] ].
```

Given $g \in \underline{ID}$ and r , a TUPLE of distinct IDs, "clubprop0" yields a club of the IDs with making the respective G-properties into grades and excluding 0-grade members. By LEMMAS 1, 2 and 3 and since loop A is executed $|r|$ times, the time complexity is $O(|r| + 1)$. 1 is added to account the time $O(1)$ needed in case $|r| = 0$, i.e., $r = ()$.

A Club Union and Grade-Adding Function:

```
addclub[p;q]=prog[[g;r;w];g:=gensym[];w:=tupset[p];
  A [null[w] → prog2[w:=tupset[q];go[B]];
  r:=addprop[g;caar[w];cadar[w];r];w:=cdr[w];go[A];
  B [null[w] → return[clubprop0[g;r]];
  r:=addprop[g;caar[w];cadar[w];r];w:=cdr[w];go[B]].
```

Given clubs p, q with numerical grades, "addclub" yields a club of the union of members of p and q with the grades of common members being added in and 0-grade members being excluded from the result. A "gensym" (i.e., a unique SYM generated by the system) is used to avoid possible confusions of properties in the auxiliary functions. Similarly, $subclub[p;q]$ is defined by replacing the "addprop" in the last line only by "subprop". Since loop A is repeated $|p|$ times and loop B, $|q|$ times and by LEMMAS 1, 2 and 3, the time complexity is $O(|p|+|q|+1)$. In case $p, q \in \underline{SP}$ "addclub" adds the two and gives the result in the SP normal form. Hence,

PROPOSITION 2. The time complexity of adding two polynomials p and q in the SP form is $O(|p|+|q|+1)$. (Multivariateness has no effect.)

A Polynomial Multiplier Function:

```
mulsp[p;q]=prog[[g;r;u;v];g:=gensym[];u:=tupset[p];
  A [null[u] → return[clubprop0[r]];v:=tupset[q];
  B [null[v] → prog2[u:=cdr[u];go[A]];
  r:=addprop[g;addclub[caar[u];caar[v]];
  cadar[u]*cadar[v];r];v:=cdr[v];go[B]].
```

Given $p, q \in \underline{SP}$, "mulsp" yields the product in the SP form. Note that "addclub" is used to multiply two TERM ID's as in $addclub\{(A,1), (B,2)\}; \{(B,3), (C,4)\} = \{(A,1), (B,5), (C,4)\}$. For $s \in \underline{SP}$, let $T(s) = |s| +$ (total number of elements in TERM ID's of s). The dominating term ($clubprop0[r]$

is $O(|p|+|q|)$ at the most) in the time complexity of "mulsp" is easily seen to be $O(|q|T(p)+|p|T(q))$, which arises from repeating the "addclub" on TERM ID's for $|p|+|q|$ times in the nested loops A and B. Hence, we obtain:

PROPOSITION 3. The time complexity of multiplying $p, q \in \underline{SP}$ is $O(|q|T(p)+|p|T(q))$; specially in case each term is K -variate at the most, it is $O(|p|+|q|(K+1))$ and in the univariate case it is $O(|p|+|q|)$. (Factors such as $\log_2 |p|$ or $\log_2 |q|$ are absent. Sparseness of the result has no effect.)

An SP* into SP Transformation Function:

$intosp[p]=mulsp[p; \{((),1)\}]$, where $\{((),1)\}=sp(1)$. This works correctly because of the "coercion rule" in (2.2). Let $T^*(p)=|p|+(\text{total number of elements in } \langle \underline{TERM ID} \rangle \text{ of } p \in \underline{SP}^*)$. We obtain:

PROPOSITION 4. The time complexity of transforming an IP polynomial p in an SP* form into the SP normal form is $O(T^*(p))$; specially in case the length of each term of p is K at the most, it is $O(|p|(K+1))$. (If TERM ID's and SP* were sorted into a sorted normal form, the time complexity would be $O(|p|(\log_2 |p|)(K+1)\log_2(K+1))$.)

(3.3) The Signed Absolute SP form:

Let $s=sp(p)$ be the SP form of a polynomial $p \in \underline{IP}$. As a SET, s can be partitioned uniquely as $s = s^+ \cup s^-$, wherein all grades of s^+ are positive and those of s^- , negative. Let $-s^-$ be the SP obtained by reversing all signs of grades of s^- .

Definition. The Absolute SP form $asp(p)$ of $p \in \underline{IP}$ is a SET: $asp(p) = \{s^+, -s^-\}$; specially $asp(0) = \{\}$.

PROPOSITION 5. For $p, q \in \underline{IP}$,

$$asp(p) = asp(q) \text{ iff } (p \equiv q \vee p \equiv -q).$$

Definition. The SASP normal form $sasp(p)$ of p is a 2-TUPLE: $sasp(p) = (asp(p), \text{sign}(p))$, where $\text{sign}(p) = +1$ in case the canonical order of the SET, $asp(p)$ is $tupset[asp(p)] = \{s^+, -s^-\}$, otherwise $\text{sign}(p) = -1$ (c.f., (2.2)); specially, $sasp(0) = ()$.

PROPOSITION 6. For $p, q \in \underline{IP}$,

$$sasp(p) = sasp(q) \text{ iff } p \equiv q.$$

(3.4) Unique Normal Forms for Rationals:

Let \underline{Q} be the set of (all) rational numbers. Hereinafter, for $q \in \underline{Q}$, we use the following obviously unique representation; if $q \in \underline{INT} \subset \underline{Q}$ use the integer q itself; otherwise use the 2-TUPLE, (a^*, b^*) such that $a, b \in \underline{INT}$, $b \geq 2$, $q = a/b$ and a, b are relative primes.

SP, ASP and SASP forms can be easily generalized to QF, polynomials with rational coefficients and positive integer exponents by changing the condition COEF $\in \underline{INTO}$ for IP's into COEF $\in (\underline{Q} - \{0\})$.

Let QF be the set of rational functions with rational coefficients and integer exponents, i.e., $\underline{QF} = \{x/y \mid x \in \underline{QF}, y \in (\underline{QF} - \{0\})\}$. Any function $r \in (\underline{QF} - \{0\})$ is known to be uniquely factorizable, except the arbitrariness of signs on the factors, as follows:

$$r = q p_1^{e_1} \dots p_i^{e_i} \dots p_k^{e_k},$$

wherein $q \in (\underline{Q} - \{0\})$, $e_i \in \underline{INTO}$ and $p_i \in (\underline{IP} - \underline{INT})$ such that p_i is not factorizable into elements of $(\underline{IP} - \{-1, 1\})$.

Definition. The Factorized SASP form $fsasp(r)$ of $r \in (\underline{QF} - \{0\})$ is a 2-TUPLE:

$fsasp(r) = (\{ \dots, (asp(p_i), e_i), \dots \}, \underline{+q})$, where
 $\underline{+q} = (\text{sign}(p_1))^{e_1} \dots (\text{sign}(p_i))^{e_i} \dots (\text{sign}(p_k))^{e_k} \dots q$;
 specially, $fsasp(0) = ()$.

PROPOSITION 7. For $x, y \in \underline{QF}$,
 $fsasp(x) = fsasp(y)$ iff $x \equiv y$.

PROPOSITION 8. For $x, y \in (\underline{QF} - \{0\})$,
 $\text{car}[fsasp(x)] = \text{car}[fsasp(y)]$ iff $x/y \in \underline{Q}$.

Proofs of PROPOSITIONS 5 to 8 have been omitted but they would be easy.

A Multiplier for $x, y \in (\underline{FSASP} - \{()\})$:

$\text{mulfsasp}[x;y] = \text{tcons}[\text{addclub}[\text{car}[x];\text{car}[y]];$
 $\text{tcons}[\text{mulq}[\text{cadr}[x];\text{cadr}[y]];()\]]$,

where "mulq" is a multiplication function of rational numbers. For a divider "divfsasp", replace "addclub" by "subclub" and "mulq" by a rational number divider "divq".

(3.5) Poisson series is a function as:

$$p = \sum_i a_i \cos(u_i) + \sum_j b_j \sin(v_j),$$

where $a_i, u_i, b_j, v_j \in \underline{QF}$.

A unique normal form POIS for this series can be obtained by absorbing the arbitrariness caused by $\cos(u) = \cos(-u)$ and $\sin(v) = -\sin(-v)$ into ASP forms: $\langle \text{POIS} \rangle ::= \langle \text{POIS COS} \rangle, \langle \text{POIS SIN} \rangle$, wherein $\langle \text{POIS COS} \rangle$ and $\langle \text{POIS SIN} \rangle$ are clubs:

$\langle \text{POIS COS} \rangle ::= \{ (\text{asp}(u), \text{sp}(a)), \dots \}$, and
 $\langle \text{POIS SIN} \rangle ::= \{ (\text{asp}(v), \text{sp}(\text{sign}(v)b)), \dots \}$.

with $u \in \underline{QF}$ and $a, b, v \in (\underline{QF} - \{0\})$. It would be a matter of exercise to define Lisp functions to perform addition, subtraction and multiplication on POIS normal forms.

(3.6) The $\langle \text{Associator List SP} \rangle$ Form:

So far stress has been laid on unique normal forms and on time complexities. However, for improvements in actual speed of computation, constant factors neglected in time complexities must be taken into account. Although time complexities of $\text{cons}[x;y]$ and $\text{tcons}[x;y]$ are both $O(1)$, "cons" would actually work faster than "tcons" because of extra hashing overhead time needed in "tcons" to ensure uniqueness. Similarly, "value", "key" and "assign" would be faster than "ass" (c.f., (2.3)). The same would hold for the $O(n)$ complexity for $\text{list}[x; \dots; xn]$ and $\text{setup}[t]$ with $|t|=n$. It would be a reasonable strategy to use unique normal forms only where they are essentially needed. For example, in the manipulation (add, sub and multiply) of $\langle \text{IP} \rangle$'s in the SP form, use of the unique normal forms for $\langle \text{TERM ID} \rangle$'s is essential but use of a $\langle \text{SET} \rangle$ for sum of terms is not. Use of the following $\underline{\text{ALSP}}$ form would be better for the sake of speed: $\langle \text{ALSP} \rangle ::= (\lambda (g^*, \langle \text{TERM ID} \rangle)) \dots \lambda$. For $p \in \underline{\text{IP}}$, $\text{alsp}(p)$ is a $\langle \text{LIST} \rangle$ of $\langle \text{ASSOCIATOR} \rangle$'s of 2- $\langle \text{TUPLE} \rangle$'s of a "gensym", g^* and a $\langle \text{TERM ID} \rangle$. $\langle \text{COEF} \rangle$'s of the $\text{sp}(p)$ are given as G-properties (i.e., $\text{get}[g^*; \text{i-th TERM ID}] = \text{i-th COEF}$). Rewriting functions for SP forms in (3.2) into those for $\underline{\text{ALSP}}$ forms would be a matter of exercise. The similar applies to Poisson series: Use ASP forms for u's and v's and $\underline{\text{ALSP}}$ forms for a's and b's.

IV. Computing Schemes with Reclaimable Hash Tables

The choice between tabulation and recomputation is a basic problem in programming. While (hashed) tabulation provides the best time complexity of $O(1)$ in many cases, extra storage space is needed to keep the tables.

In HLISP two features called tabulative and associative computing are provided, which enable users to utilize the full advantages of computing with hash tables. Moreover, in order to make a compromise between the space and time requirements automatically, a two staged garbage collection scheme, GBC and GGBC of (2.5), is employed. The $\langle \text{CELL} \rangle$'s used for hash table entries in "tab-" and "assoc-comp" schemes are reclaimed by GGBC but not by GBC. Hence, these entries are termed "reclaimable". After having been reclaimed, the table entries are reconstructed on demand.

(4.1) "Tabcomp" is applied to $\text{member}[x;s] = (x \in s)$ for $x \in \underline{\text{ID}}$, $s \in \underline{\text{SET}}$ and to n-way switching and selecting functions: $\text{tab}\alpha\beta[x;a;e^*]$ with $\alpha \in \{a, d, q, g\}$ and $\beta \in \{q, g\}$. The value of a must be an n- $\langle \text{TUPLE} \rangle$ of the form $a = (\dots, (m_i^*, g_i^*), \dots)$ and e^* must be a constant $\langle \text{ID} \rangle$ datum. If x matches with m_i ($\in \underline{\text{ID}}$), the resultant value is respectively $\text{cadr}[m_i^*, g_i^*] = g_i^*$, $\text{cdr}[(m_i^*, g_i^*)] = (g_i^*)$ or (m_i^*, g_i^*) for $\alpha = a, d$ or q ; for $\alpha = g$ the result is "GO TO g_i^* ". If no match, for $\beta = q$ the resultant value is e^* and for $\beta = g$ the result is "GO TO e^* ".

(4.2) "Assocomp" effectively avoids the recomputation of the same function for the same argument(s) by inserting the results of the previous computation in the reclaimable hash table entries. Evaluation of a function is made in the "assocomp" mode by so specifying to the compiler or interpreter. By "assocomp", the time complexity of recursive algorithms such as follows can be improved automatically without rewriting.

$\text{factorial}[n] = \text{fc}[n] = [n=0 \rightarrow 1; T \rightarrow n * \text{fc}[n-1]]$,
 $\text{fibonacci}[n] = \text{fb}[n] = [n \leq 1 \rightarrow n; T \rightarrow \text{fb}[n-1] + \text{fb}[n-2]]$,
 $\text{C} = \text{c}[n;m] = [m=0 \vee m=n \rightarrow 1; T \rightarrow \text{c}[n-1;m] + \text{c}[n-1;m-1]]$.
 $n \ m$

(4.3) LEMMA 5. Time Complexities of Tab- and Assoc-comp features are as in the following table:

Function	WITHOUT Tab- and Assoc-comp features.	WITH Tab- or Assoc-comp		
	TIME	INITIAL TIME	REPEATED TIME	EXTRA CELLS
$\text{member}[x,s]$	$O(s)$	$O(s)$	$O(1)$	$ s $
$\text{tab}\alpha\beta[x,a,e^*]$	$O(a)$	$O(a)$	$O(1)$	$ a +1$
$\text{factorial}[n]$	$O(n)$	$O(n)$	$O(1)$	$2n+3$
$\text{fibonacci}[n]$	$O(1.618^n)$	$O(n)$	$O(1)$	$2n+3$
$\text{C} = \text{c}[n,m]$	$O(\frac{C}{n \ m})$	$O(n^2)$	$O(1)$	$3n^2/2$

The initial time means the time complexity immediately after a GGBC call. Extra cells are the number of $\langle \text{CELL} \rangle$'s needed for reclaimable hash entries. E.g., repeated evaluation of $\text{fb}[21] = 10946$ runs 30,000 times faster in HLISP by merely feeding a card "ASSOCCOMP ((FB))". $\text{clubmember}[x;c] = \text{tabq}[x;\text{tupset}[c];()]$ checks whether x is a member of the $\langle \text{CLUB} \rangle$, c . The time complexity of $O(|s|+|t|)$ in the pure Lisp algorithms [3] for $s \cup t$ and $s \cap t$ of sets s, t is greatly improved by applying "tabcomp" to "member" (even immediately after a GGBC call):

LEMMA 6. Time complexity of $s \cup t$ and $s \cap t$ for $s, t \in \underline{\text{SET}}$ is $O(|s|+|t|)$.

(4.4) Outline of an HLISP Implementation:
 For "member" $\langle \text{SYS2 CELL} \rangle$'s of (2.5) are utilized. When $\langle \text{SYS2 CELL} \rangle$'s are reclaimed by GGBC, the $\langle \text{SYS1 CELL} \rangle$ is switched to a $\langle \text{SYS1* CELL} \rangle$ to indicate the necessity of reconstruction of the $\langle \text{SYS2 CELL} \rangle$'s. For "tab $\alpha\beta$ ", initially (i.e., after GGBC) a $\langle \text{SYS3 CELL} \rangle ::= [\text{SYS3}, a^*, e^*]$ is hash inserted (as a result of an unsuccessful search) and then

<SYS4 CELLS> ::= [SYS4, (mi*, gi*), [SYS3, a*, e*]]. are hash inserted by using a hash sequences determined by mi's (not the <TUP> (mi, gi)) and the pointer to the <SYS3 CELL>. Hash retrieval is made by utilizing these <SYS3 CELL> and <SYS4 CELL>'s, which are all reclaimed by GGBC. In the assoccomp mode, a function fb[n], say, is evaluated as: First, make a hash search for <SYS5 CELL> ::= [SYS5, "don't care", t], with t=tcons[n;FB], and if unsuccessful insert a <CELL>, [SYS5, l*, t], where l* is a <SYSTEM SYMBOL>, then compute fb[n] and replace l* by fb[n] for future retrieval of fb[n]. Else if successful retrieve the value from the <CAR field>. Specially, in case the <CAR field> contains l*, there must have been a vicious circle in the algorithm such as fb[n]=fnsl + n; T + fb[n]+ fb[n-1]]. Thus a message "CIRCULAR DEFINITION ERROR IN FB ..." is printed. GGBC reclaims <SYS5 CELL>'s except those containing l*. Hence,

LEMMA 7. "Assoccomp" effectively checks circular definitions at runtime.

(4.5) For fc[n], fb[n], c[n,m] etc., "assoccomp" is more convenient than "tabcomp" since the range of argument(s) is generally not known in advance. Conversely, if "assoccomp" were used for member[x;s], say, a great number of wasteful hash entries for x / s would be created. Thus, "tab- and assoc-comp" are complementary and each has its own raison d'être.

V. Concluding Remarks

The first version of HLISP without the SET feature has been in operation for two years[8], but with the TUP feature alone little advantage in formula manipulation could be found. The combination of SETs and TUPs is believed to have provided a really powerful tool for formula manipulation as indicated in III. Tab- and assoc-comp features would also be useful. Since the implementation of efficient hashing and garbage collection algorithms is a very specialized art, it would be better to

separate them from the general users. Therefore, external specifications of such algorithms have been given as LEMMAS in this paper.

The following improvements are now in progress to make the schemes presented in this paper into truly useful tools for symbolic and algebraic computations:

- (1) Writing of an efficient HLISP compiler[9].
 - (2) Implementation of a language system called "FLATS" which would enable us to absorb any existing algorithm written in Fortran, Lisp or Algol 60; and to write new algorithms with Tuples and Sets added to any of the three languages F, L or A, whichever the user may prefer (HLISP = FLATS).
 - (3) Design of hashing, GBC and runtime type check hardware to improve the ultimate speed of "FLATS".
- The authors acknowledge Messrs. M. Terashima[10] and F. Motoyoshi[9] for their valuable contributions in implementing HLISP.

VI. References Inf. Proc. Letters 5('76)pp.31-34

- [1] S.C. Johnson, SIGSAM Bulletin, 8, 3, p.63, '73.
- [2] E. Horowitz, J. ACM, 22, 4, p.450, 1975.
- [3] J. McCarthy, et al., LISP 1.5 Programmer's Manual, MIT press.
- [4] D.E. Knuth, The Art of Computer Programming, vol. 3, Addison-Wesley, Reading, Mass., '73.
- [5] M. Sassa and E. Goto, A Hashing Method for Fast Set Operations, submitted for publication.
- [6] T. Gunji, Tech. Rep. 76-03, ISD (Information Science Department, the University of Tokyo), 1976.
- [7] J. McCarthy, Page 151 of Symbol Manipulation Languages and Technique, D. Bobrow, ed., North-Holland, 1971.
- [8] Y. Kanada, Tech. Rep. 75-01, ISD, 1975.
- [9] F. Motoyoshi, Tech. Rep. 76-05, ISD, 1976.
- [10] M. Terashima, Tech. Rep. 75-03, ISD, 1975.
- [11] A.C. Hearn, REDUCE2 User's Manual, 2nd. ed., Salt Lake City, Utah., 1973.

APPENDIX. Actual Timing Data for Polynomial and Poisson Series Manipulations.

- REMARKS: (1) The machine used is HITAC 8800/8700 at the Computer Centre of the University of Tokyo.
 (2) The same HLISP interpreter system was used as the host system for REDUCE 2[11]. The free storage area was 75K cells in which 25K cells were reserved for <ID> objects.
 (3) The data for polynomial multiplication were obtained to observe the dependence of time on n (number of terms in polynomials) and multiplicity, K. Observed times were normalized by n²(K+1) as PROPOSITION 3 predicate. Unit of time is in msec. '*' means 'not measured'.
 (4) The FORTRAN data of univariate case were taken by a program with explicit code for hashing. The program is similar to the algorithm by Gustavson and Yun to be given at this SYMSAC '76. The hash area was selected to 5011 (a prime) and the hash probe sequence was given by Algorithm U2 of Knuth[4, p539].
 (5) The programs in HLISP were written for the ALSP and ASP forms of (3.6).

Formulas t=resultant # of terms	n				n				n			
	4	8	16	32	4	8	16	32	4	8	16	32
$(\sum_{i=1}^n A^i) * (\sum_{j=1}^n A^j)$ t=2n-1	1.71	1.69	1.60	1.60	1.85	1.73	1.71	1.67	1.82	1.74	1.74	1.77
	4.42	2.95	3.97	5.45	3.67	3.50	4.43	7.20	4.65	4.04	5.54	9.10
	.025	.024	.020	.016								
$(\sum_{i=1}^n A^i) * (\sum_{j=1}^n A^{jn+1})$ t=n*n	1.76	1.74	1.72	1.73	1.98	1.78	1.76	1.80	1.81	1.80	1.79	1.84
	5.50	6.08	15.4	51.3	4.33	7.37	21.6	*	4.40	8.48	*	*
	.025	.028	.020	.018								
$(\sum_{i=1}^n A^{-2+3i}) * (\sum_{j=1}^n A^{-3+4j})$ t=7n-12	1.96	1.71	1.68	1.63	1.88	1.82	1.73	1.74	1.84	1.83	1.79	1.77
	5.35	5.85	8.20	14.3	5.42	6.53	10.6	*	5.16	7.64	12.2	*
	.028	.025	.020	.016								
K-variate	1-variate (A=X)				2-variate (A=XY)				4-variate (A=XYZU)			

Timing Data for Poisson Series Manipulation: HLISP REDUCE
 (A1*COS(WT)+A3*COS(3*WT)+B1*SIN(WT)+B3*SIN(3*WT))**3 1587 msec 8077 msec