

32ビットアドレス向きLISP処理系の実現

橋本ユキ子 内田誠二 小川雄司 高橋哲夫
日本電気株式会社 基本ソフトウェア開発本部

東京大学和田研究室で、MC68010上に作成された高速なLISP処理系Utilisp
(University of Tokyo Interactive LISt Processor)
を32ビットアドレスをもつMC68020上に移植した処理方式について述べる。

MC68010では、アドレスバスが24ビットであるため、この上に作成されたUtilispは、ポ
インタデータ内にアドレス表現とは別にタグを設けることができる。このUtilispを4Gバイトのア
ドレス空間を扱えるようにするため、ポインタデータ内のアドレス表現の一部を直接タグとして使用する方
式を考案し、処理系作成、評価を行った。

IMPLEMENTATION OF LISP ON 32-BIT-ADDRESSING CHIPS

Yukiko HASHIMOTO, Seiji UCHIDA,
Yuji OGAWA and Tetsuo TAKAHASHI

Basic Software Development Division, NEC Corporation
1-10, Nisshin-cho, Fuchu-shi, Tokyo 183, JAPAN

This paper describes an implementation of Utilisp* on MC68020 using 32 bit addressing.
This version originates the MC68010 version of Utilisp designed and made at
the University of Tokyo, Wada Laboratory.

Because MC68010 has 24 bit address-bass, it is possible to use the left-most 8 bits of a pointer
data as a tag field. In case of MC68020, we can not provide a tag field in a pointer data in the
same way as MC68010. The authors will propose the way to recognize every lisp object with some of
middle bits of the address expression.

* University of Tokyo Interactive LISt Processor

1. はじめに

近年32ビットアドレスを持ちアドレッシング範囲が4Gバイトのアーキテクチャを持つコンピュータが増えている。このようなマシン上にLISP処理系を実現しようとした場合、必ず問題となるのがタグの扱いである。MC68010のようにアドレッシングの範囲が24ビットである場合には、ポインタデータ内にアドレス表現とは別にタグを設けることが多く、そうして実現されたLISP処理系では4Gバイトの空間を扱うことは不可能となる。32ビットアドレスで表現できるアドレス空間を扱えるようなLISP処理系の実現方法としては、ポインタデータとは別にタグ領域を設ける方法やポインタデータのさすオブジェクトデータ側へタグを持たせる方法など幾つかの方法が考えられる。本稿ではポインタデータ内にアドレス表現とは別にタグを含める方式と同等の効率を維持することを目的として、東京大学和田研究室でMC68010上に作成された Utilisp を32ビットアドレッシングをもつMC68020上へ移植した際の処理方式について報告する。以下従来方式としてはMC68010上の Utilisp を基に説明を行う。

2. 32ビットアドレッシングのための技法

2.1 ポインタ形式とタグ

32ビットアドレッシングが可能なLISP処理系を実現する場合、MC68010上の Utilisp で実現されているようにポインタデータ内にアドレス表現とは別にタグを設定することはできない。この方式でLISP処理系を32ビットアドレス空間上に実現した場合はLISP処理系がタグ領域を1バイトとすると16Mバイトより大きなアドレスをもつ場所へローディングされない保証が必要であり、またアドレッシングのときには必ずタグクリアが必要となってしまう。そこで本処理系ではアドレス表現自身の一部を直接タグとして使用する方式を採用した（以下本方式をアドレスタグ方式と呼ぶ）。

この場合アドレス表現のどの部分でもタグとして使用することができるがMC68020ではワード（16ビット）を扱う命令が豊富であり実行時間もロングワード（32ビット）を扱う場合より速い。そこでこれを有効に利用するため図1に示すようにポインタデータの第11ビットから第15ビットまでの5ビットを基本的にタグとして使用している。

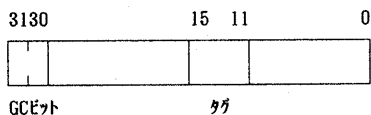


図1. ポインタデータの形式

GCはmark&sweep法を使用しており、GCビットは2ビット必要なため第30、31ビットを使用する。本システムではユーザ空間が1Gバイトに制限されて

いるためGCビットとして第30、31ビットを使用しているが、LISP処理系を4バイト境界に整列して作成しておくことにより第0、1ビットを使用することができ事実上4Gバイトのアドレッシングが可能となる。

Utilisp本来の高速性を損なわないように設計を行うためタグの決定についても従来の24ビットアドレス空間上の Utilisp と同じ効率で各タグが扱えるように考慮する必要がある。タグ決定のポイントは各データタイプ中出現頻度の高い順に領域の割合を定めることとアクセス頻度の高いデータをより高速に扱えるように工夫をすことである。

データタイプとタグの関係を表1に示す。

データタイプ	タグ (16進)
list	00 ~ 78
symbol	80 ~ 88
flonum	90 ~ 98
vector	A0 ~ A8
code piece	B0 ~ B8
string	C0 ~ D8
reference	E0
bignum (用)	E8
extra	F0
stream	F2
undefined	F4
fixnum	F8

表1. データタイプとタグの関係

listは各データタイプ中で最も大量に割り付けられるオブジェクトであり、listかatomかの判定回数も非常に多い。そのためlistかどうかの判定は従来通り正負の判定で可能とし、リスト領域もヒープ領域全体の1/2を占めるようになってきている。また、symbol、fixnumのアクセス頻度も高いためやはり高速に判定が行えるように定められている。stream、undefined、extraの各タイプについては必要最小限の領域が確保できればよいので更に細かく領域を分けポインタデータの第7ビットから第10ビットもタグとして使用することにする。

MC68010版とMC68020版におけるタグ判定法の代表例を表2に紹介する。

データタイプ	68020 判定法	68010 判定法
list	mov.w R1,R2 bpl yes	mov.l R1,R2 bmi yes
symbol	cmp.w Rn,SY bgt yes (SY=0x8fff)	cmp.l N,Rn bge yes
string	mov.w R1,R2 and.w 0xe000,R2 cmp.w 0xc000,R2 beq yes	mov.l R1,R2 swap R2 and.w 0xff00,R2 cmp.w 0x4000,R2 beq yes
fixnum	cmp.w 0xf800,Rn bcc yes	cmp.l maxfix,Rn bls yes

表2. タグ判定法

2.2 メモリ構成

2.1図1で示したようにポインタデータ内の第15ビットから第11ビットをタグとして使用した場合ヒープ領域は64Kバイト毎のセグメントに分割され更に各セグメントがタグの種類に従ってLISPの全データタイプに分割されることになる。

MC68020版Utilispのメモリ構成を図2に示す。(次頁)

本方式では当然のことながら predefined objectも各々決まったタグをアドレスとしてもつように配置されなくてはならないためメモリロードされた時点で第2図のように展開されるように、また無駄な領域が作られることのないよう調整を行っている。

本方式における問題点として例えばLISPで記述したアプリケーションがfloating データは使用していないがシンボルの数が非常に多いというように1セグメントのデータタイプ毎の分割の割合がアプリケーションのオブジェクト使用量に合っていない場合が考えられる。仮想記憶管理を備えているOS上では全く使用されていない領域は実メモリに割り付けられることはないでメモリ使用効率はあまり問題にはならないが、それでも実メモリが少ない場合は仮想記憶空間を圧迫することにはなる。それをできるだけ解消し最適な構成のLISP処理系を作成するためにタグの再編成方式を考案した。(詳細は「4. タグ再編成方式」参照)

2.3 データオブジェクトの特長

本節ではアドレスタグ方式の採用により元のUtilisp処理系とは異なった形式で実現しているデータオブジェクトについて説明する。変更が加えられたデータは fixnum,reference,extra である。

(1) fixnum

従来fixnumはオブジェクトを持たずポインタデータの代わりに直接値をタグとともに格納している。アドレスタグ方式を実現するとポインタデータがさすデータ側にfixnum値を格納しなければならなくなるがそれではfixnumをアクセスする際常に間接参照となり明らかに効率が低下する。またfixnum領域も大量に必要と

なってしまう。そこで従来と同様にポインタデータの代わりに値を格納する手段としてfixnumは次の構造をとっている。

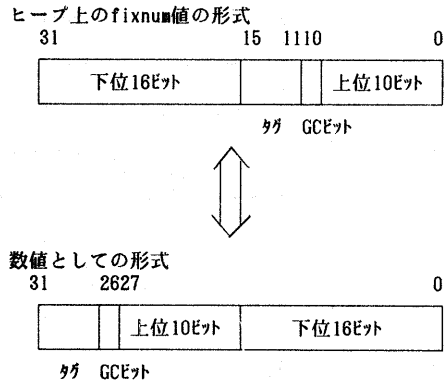


図3. fixnum形式

ヒープ上のfixnum値の形式と数値としての形式の間の変換はMC68020上のSWAP命令1命令で行うことができる。(間接参照よりSWAP命令を使用したほうが高速である。)

メモリ領域には図2に示す通り各セグメントにfixnum領域2Kバイトが確保されている。この領域はfixnumオブジェクトが存在していないため実際に使用されることはない。そこでこの領域はI/Oバッファとして現在利用している。

(2) extra

図2のようなメモリ構成をとっているため可変長オブジェクト スtring, ベクタは各々のサイズを越えてアロケートすることができない。そこでこれを越える大きさのString, ベクタを作成したい場合 extraオブジェクトと呼ぶ別のデータをアロケートしString, ベクタの実体へのポインタを格納する。

例としてXが10KバイトのString "ABC..."を値とするシンボルであるときのエクストラオブジェクトの構成を図4に示す。

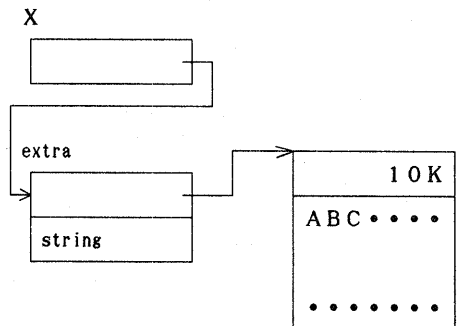


図4. extraオブジェクトの構成

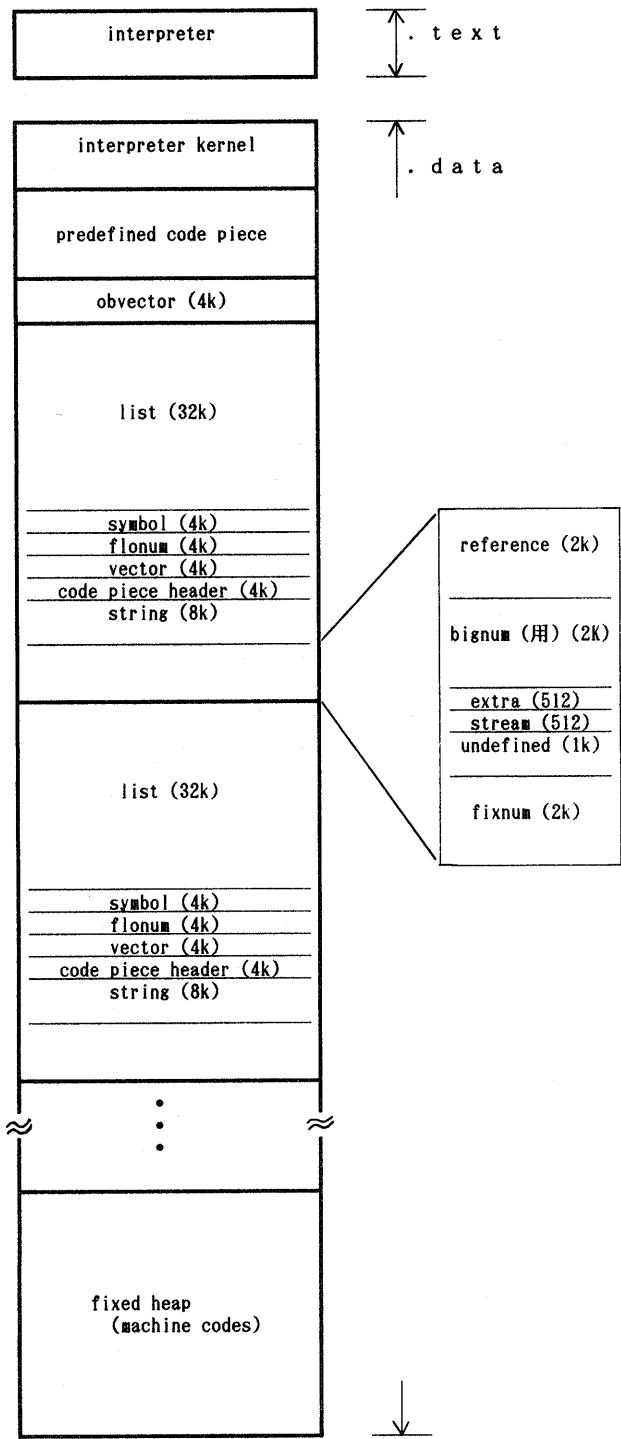


図2. Utilisp のメモリ構成

(3) reference

従来 referenceはオブジェクトを持たずポインタデータは直接ベクタ要素をさしていた。本方式においては reference オブジェクトが存在しベクタの要素をさすポインタとベクタの先頭をさすポインタから構成されている。

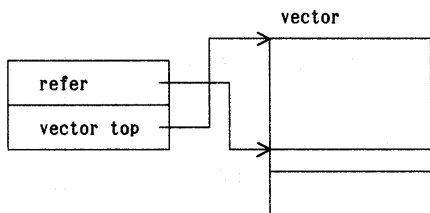


図5. referenceオブジェクトの構造

3. ガーベジコレクション

本処理系ではヒープ領域は図2に示すようにデータタイプ別に分類されているため、各領域の管理方式は固定長データ、可変長データにより異なっている。可変長オブジェクトは各領域の先頭から順に必要な大きさを割り付けているが固定長オブジェクトは未使用領域をフリーチェーンによって管理している。各データのセル自身をフリーチェーンとして使用しており当該セルから次のセルへの相対距離を格納している。

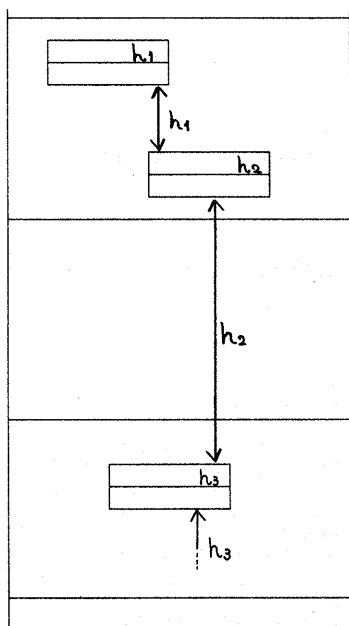


図6. 固定長オブジェクトの管理方式

GC (ガーベジコレクション)はこのヒープ領域管理に従い固定長オブジェクトと可変長オブジェクトで処理が異なる。固定長オブジェクトのGCはゴミとなったセルをフリーチェーンに繋いでいく処理を行い、可変長オブジェクトのGCはmark&sweep法により使用中のオブジェクトを各ヒープ領域の先頭に寄せ集める。

GCによって使用するビット位置は次のようになっている。fixnumの場合は第30、31ビットはfixnum値を表しGCビットとしては使用できないため第10ビットを使用している。

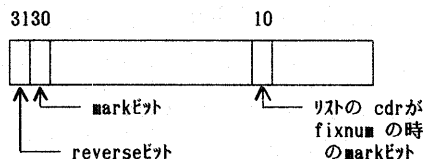


図7. GCビット

GC処理手順は次の通りである。

(1) マーキング処理

現在使用中のオブジェクトのマークビットをオンにする。ただしリストの cdrがfixnumのときは cdr 部の第10ビットがGCビットであり0のときがオン、1のときがオフである。(fixnumではタグとGCビットが全て1であれば負の数のときにタグの設定が不要になる。)またシンボルオブジェクトは pname部にマークをしその他のオブジェクトはヘッダ部にマークをする。

(2) ポインタの逆転処理

既定義のオブジェクト、スタック等から参照されているヒープ領域中の可変長オブジェクト(ストリング、ベクタ)へのポインタを逆転する。

(3) 未使用オブジェクトの回収

a. 固定長オブジェクト

各オブジェクトのヒープ領域を走査しフリーチェーンを作成する。使用中オブジェクト内にストリングまたはベクタをさすポインタがあれば逆転する。

b. 可変長オブジェクト

各オブジェクトのヒープ領域を走査し使用中オブジェクトは領域の先頭へ寄せ集める。このとき逆転ポインタをオブジェクトの移動先へ戻す。

本方式のGCは固定長オブジェクトはフリーチェーンで管理し移動しないため従来のmark&sweep法と比べると逆転ポインタ処理がかなり省略されておりGCは高速になっている。実際に逆転ポインタ処理を行うのはストリングとベクタだけである。また従来のmark&sweep法ではヒープ領域に対し順方向と逆方向の2度の走査処理が必要であった。本方式ではベクタ自身がポインタをもっているためコンパクションをするには2回の走査処理が必要となるが、他のオブジェクトは1度の走査で済ませることができる。

本GCでは固定長オブジェクトはヒープ上に寄せ集められず分散してしまうが アドレスタグ方式によりヒープ領域は64Kバイトのセグメントに分割されており本来ローカリティはあまり上がらない。従って高速にGCが行えるほうがメリットが高い。

4. タグ再編成方式

本節では、ヒープ領域内のメモリ構成を変更する方式について述べる。

4. 1 背景

第2節で述べたように、本処理系ではアドレスデータの内部の数ビットをタグとして使用している。したがってヒープ領域は各64kバイトの領域に分割され、図2に見るようにそのそれぞれの領域がさらにタグの値に対応する小領域に分割される。すなわちタグの編成がそのままヒープ領域の使用区分となっている。最も標準的と思われる編成をもって既定値としているが、特定の種類のデータを多用するアプリケーションプログラムを効率よく動作させるためには、多用される種類のオブジェクトを格納する領域がより大きくとれることが必要であり、そのためにはタグの編成を変更することが不可欠である。

4. 2 再編成の対象となるデータ領域

本処理系において、list,fixnum,symbolの3種のデータは特に高速に判定されるタグに対応する領域に格納されており、セグメント内の位置を変更することは困難である。一方、たとえばstreamオブジェクト用の領域は、オブジェクトの総量に上限があり、既定値のタグ編成において全てを使用可能としており、特に再編成の対象とする必要はない。再編成の対象としないデータ属性およびその理由は次の通りである。

- (1) list
タグの符号のみで判定が行われており、これが最高速のタグ判定であることによる
- (2) fixnum
判定の高速化のため、各セグメントの最終位置をあらわすタグがついており、最終位置であることを変更することはできない。また、実際にこのタグに対応する領域は2kバイト長のI/Oバッファとして使用されているため、サイズも変更することができないことによる。
- (3) streamおよびundefined
増減がないため。
- (4) extra
さほど多用されず、現行のままで十分であるとされるため。

上記以外のデータに対するヒープ領域、すなわち図8の領域に対して再編成処理を行う。

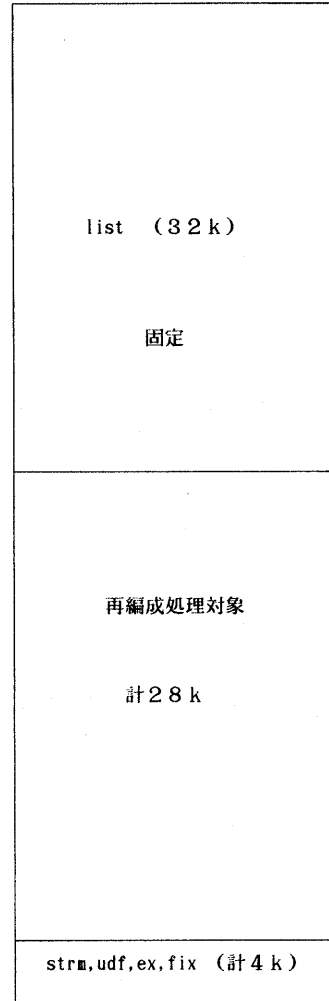


図8 再編成処理の対象となる領域

4. 3 再編成を行う契機

utilisp コマンドの通常の実行形態は、コマンドのオプションパラメータとしてヒープ領域の大きさ等が入力され、十分な大きさの領域が割り付けられた後、インタプリタ本体に制御がわたり、その初めの処理としてヒープ領域の初期化すなわちメモリ管理のためにヒープ上にチェーンを作成する処理を行い、システムファイルをロードした後、read-eval-print ループに入るといった形態をとっている。ヒープ領域等の拡張割り付け自体はタグの編成とは無関係であり、また、ヒープ領域の初期化処理はヒープが拡張されているか否かにはよらず、常にタグの編成に依存する定数を用いて行われる。これらから、タグの再編成処理の契機は、utilisp コマンドのオプションパラメータによって与えられる。契機が与えられなかった場合には

そのまま通常の実行形態をとる。 タグの再編成を行ったutilispの立ち上げは4.5で述べるが、このことによりインタプリタ自身は自分が既定値のutilispであるのかそれとも再編成されたutilispなのかを知らずに動作する。 すなわち、インタプリタ本体に制御が渡った後は、タグの再編成の有無による処理速度の変化は全くない。

4.4 処理概要

4.2で述べた再編成処理対象領域は、既定値のタグ編成において2k、4kあるいは8kバイトの小領域7つからなり、合計28kバイトである。 そのうち、symbolのタグ判定においてはレジスタとの比較が1回行われ、その他のタグ判定においては即値とのANDおよび即値比較が1回行われる。

即値を書きかえるのみでタグチェックを更新するためには、編成後の各領域の大きさは2k、4kあるいは8kバイトのいずれかでなくてはならない。 合計28kバイトの領域を7種のデータ領域に分割する可能な組み合わせは、大きさの順に、kバイト単位で

- (1) 8, 8, 4, 2, 2, 2, 2
- (2) 8, 4, 4, 4, 4, 2, 2
- (3) 4, 4, 4, 4, 4, 4, 4

の3種であり、データ属性の組み合わせにより、211種類の異なる編成が作成可能である。

タグの再編成処理は、

- (1) 新しいタグ編成の決定
 - (2) タグ値に依存する大域データの更新
 - (3) タグチェックに関するマスク値、タグ値の即値データの更新
 - (4) タグ値による多分岐用テーブルの更新
 - (5) predefinedオブジェクトに対する参照の再解決
 - (6) 再編成済みutilispをファイルに出力
- の6つの処理から構成されている。

新しいタグ編成の決定は対話的に行われる。 利用者は各データ属性およびそれらの領域サイズを入力する。 システムはそれらの値に基づき可能なタグ編成を決定し、利用者の確認を得る。 続く4つの更新処理に対しては、書きかえるべき即値のあるアドレス等をあらかじめ持っている必要がある。 これはLAP展開において、書きかえられる可能性のある箇所にラベルを出力し、同時に別のファイルにラベル名からなる表を作成し、utilisp作成と同時にアセンブル/リンクすることにより得る。 この表を格納した領域はutilisp本体が動作するときには不要であり、クリアされた後、ヒープの一部として使用される。 新しいタグ編成が決定されれば、上記の表を用いて処理(2) - (4)は容易に行われる。

処理(5)に対しては、再編成対象となるヒープ領域を2kバイトごとのブロックで管理し、各ブロックへの参照をLAP展開時に表にすることにより、参照の再解決を行う。

最後に処理(6)において、タグの再編成されたutilispをファイルに出力することにより、タグの再編成処理を終了する。

4.5 再編成されたutilispの立ち上げ

4.3で述べたように、utilispコマンドが起動されると、既定値のutilispがロードされる。 タグの再編成を行ったutilispを使用する旨のオプションパラメータと、前項で出力されたファイル名をコマンドの引数に指定することにより、そのファイルの内容が既定値のutilispに上書きされ、utilisp本体の初期化ルーチンに制御が渡る。 以後、タグの再編成されたutilispが動作する。

5. 性能評価

本処理系はEWS4800 (UNIX SYSTEMV) 上に作成した。 表3にLISPコンテストの幾つかのベンチマークの測定結果を示す。 従来の24ビットアドレスをもつUtilisp処理系と比べほぼ同等の性能を上げておりUtilispの高速性は維持されているといえる。

	(MS)	
ベンチマーク	インタプリタ	コパイラ
TARAI-4	2108.3	196.6
LIST-TARAI-4	3733	1033
STRING-TARAI-4	4434	2384
SEQ-100	91.51	48
BITA-5	25.83	5.16
TPU-7	2150	416
BOYER	141366	18232

表3. 第3回LISPコンテスト測定結果

アドレスタグ方式を採用しているためデータのローカリティが上がらないので大きなアプリケーションを実行した場合仮想記憶領域と実メモリの間のスラッシングの問題が非常に心配であったが、使用されていない領域は実メモリに割り付けられることはない、ワークステーションで実現したため実メモリがふんだんに使用可能となり、心配していた問題は起こらなかった。

6. おわりに

32ビットアドレス空間上に高速なLISP処理系 Utilispを実現した方法について述べた。ポインタデータのアドレスの一部をタグとして使用しているため幾つかの欠点もあるが、おおむね満足のいく処理系が作成できた。拡張機能として 他言語呼び出し機能 (C, FORTRAN)、日本語処理機能を追加している。今後の課題としては コンパイラの性能改善、プログラミング環境の充実などがある。

最後に、本処理系を作成するにあたりMC68010上の Utilispを提供して頂き、方式設計の際有益なご意見を賜った東京大学和田教授ならびに研究室の方々に深く感謝致します。

参考文献

- (1) 近山隆：“Utilisp システムの開発”，情報処理学会論文誌，Vol.24 No.5, Sep.1983
- (2) MOTOROLA：“MC68020 32-Bit Microprocessor User's Manual”，Second Edition
- (3) 近山隆：“UTILISP MANUAL”，TECHNICAL REPORTS, METR 81-6, September 1981