

Représentation des objets et des types

Introduction

Le contexte dans lequel on étudie les représentations des objets et de leur types est celui où:

- des objets sont créés dynamiquement et
- le typage est dynamique.

Par exemple, ce contexte correspond aux choix de représentation qu'il faut faire pour implanter Scheme.

Introduction

Définitions:

Une représentation est *directe* lorsque la valeur des objets représentés se trouve directement dans la représentation employée.

Par exemple, les données de types `char`, `int` et `float` en C sont représentées directement.

Une représentation est *indirecte* lorsque la représentation consiste en un pointeur (ou un index) et que la valeur des objets est accessible uniquement au travers du pointeur (ou de l'index).

Par exemple, les structures de données allouées dans le tas en C ou les objets en Java.

Introduction

Première tentative de représentation des objets de Scheme:

Représentation directe avec structures, unions et scalaires C.

Cette représentation est forcément inadéquate car un objet peut en contenir un autre.

Par exemple, pour la donnée (1 2), l'objet principal est une paire, disons de type DATA, laquelle doit contenir deux champs pour décrire le contenu de la paire. Chacun de ces champs doit être capable d'accommoder un autre objet, de type DATA. En particulier, le deuxième champ contient justement une autre paire. Ce qui nous mène à l'inéquation insatisfaisable:

$$\text{sizeof}(\text{DATA}) \geq 2 \times \text{sizeof}(\text{DATA})$$

Donc, il faut absolument employer l'allocation dynamique et les pointeurs.

Introduction

Deuxième tentative de représentation des objets Scheme:

Représentation indirecte avec des pointeurs sur tous les “contenus” possibles.

```
struct Data;

typedef struct Data DATA;

union Contenu
{
    int    entier;
    float reel;
    struct { DATA car; DATA cdr; } paire;
    ...
};

struct Data
{
    int type;
    union Contenu *ctn;
};
```

Les problèmes avec cette représentation sont:

- qu’il y a un gaspillage d’espace car l’union doit être aussi grande que le plus grand de ses membres;
- qu’on ne peut représenter les objets à taille variable comme les vector de Scheme.

Introduction

Approche générale:

- Par souci d'efficacité et d'économie de mémoire, on préfère laisser tomber les types C.
- On gère *soi-même* la représentation.
- On suppose que le tas est une longue plage de mémoire composée de mots machine; typiquement, des `int` de C.
- Une allocation dans notre tas consiste essentiellement à réserver quelques mots consécutifs et à conserver l'adresse du premier de ces mots.

Représentation par pointeurs d'objets

Description: Tout objet est représenté par un pointeur vers le tas. La partie du tas qui est pointée consiste en un champ pour indiquer le type de l'objet suivi de 0 ou plus champs en fonction du type.

Par exemple:

Une paire: →

PAIR	ref	ref
------	-----	-----

Un vecteur: →

VECTOR	len	ref	...
--------	-----	-----	-----

Un réel: →

REAL	float
------	-------

Une liste vide: →

EMPTY

Un caractère: →

CHAR	char
------	------

Un booléen: →

BOOL	int
------	-----

Un petit entier: →

FIXNUM	int
--------	-----

Un grand entier: →

BIGNUM	len	int	...
--------	-----	-----	-----

Une fermeture: →

CLOS	C-ptr	env
------	-------	-----

...

Représentation par pointeurs d'objets

Avantages et inconvénients:

- + Régularité et simplicité
- Inefficacité
 - Tout objet doit être créé (même un booléen!)
 - Tout test de type ou tout accès à la valeur d'un objet requiert un accès en mémoire

Représentation par pointeurs d'objets

Améliorations possibles:

- Assignation de types spécialisés pour les caractères et les booléens.

Par exemple:

Un caractère: \longrightarrow CHAR_i où $i \in \{T_C, \dots, T_C+255\}$

- Pré-allocation des objets de certains types. Plus de création d'objets de ces types. Comparaison par pointeur possible.

Par exemple, la liste vide (OK), les booléens (OK), les caractères ASCII (OK), les caractères Unicode (???)

Représentation par pointeurs taggés

Idée:

Certains bits peuvent être constants dans les adresses vers le tas:

- ceux du bas (lsb), à cause de l'alignement;
- ceux du haut (msb), à cause de l'étendue du tas.

On peut profiter de ces bits constants pour:

- encoder de l'information de typage;
- avoir une représentation directe de certains types d'objets.

Représentation par pointeurs taggés

Exemple:

Supposons que les adresses vers le tas sont alignées sur 8 octets et que les mots sont de 32 bits. Donc, les 3 lsb sont connus et sont égaux à 0.

xx...xx000	→	fixnums	(repr. directe)	(sur 29 bits)
xx...xx001	→	caractères	(repr. directe)	
xx...xx010	→	booléens	(repr. directe)	
xx...xx011	→	liste vide	(repr. directe)	
xx...xx100	→	paires	(repr. indirecte)	
xx...xx101	→	fermetures	(repr. indirecte)	
etc.				

Représentation par pointeurs taggés

Définitions:

L'opération d'*emballage* (*wrapping*) d'une adresse ou d'une donnée brute consiste à la convertir sous le format des pointeurs taggés.

L'opération de *déballage* (*unwrapping*) est l'opération inverse.

Il faut éviter de confondre ces opérations avec la création des objets et l'extraction des valeurs contenues par ceux-ci.

Exemples référant à la représentation de la page précédente.

- La création d'un petit entier consiste seulement en un emballage.
- La création d'une paire consiste en l'allocation et l'initialisation de quelques mots dans le tas et en l'emballage de l'adresse.
- L'extraction de la longueur d'un vecteur consiste en le déballage du pointeur taggé, en la lecture du deuxième champ du vecteur et en l'emballage de la valeur brute lue.
- Dans le cas de la représentation par pointeurs d'objets, les opérations d'emballage et de déballage sont inexistantes.

Représentation par pointeurs taggés

Exemple d'implantations des opérations liées au “tagging” en supposant que les 3 lsb servent au tagging:

Test de type: Pour tester si on a affaire à une fermeture:

$$\text{pred_clos}(p) \equiv (p \ \& \ 7) == 5$$

Emballage: Pour emballer un caractère brut:

$$\text{wrap}(c) \equiv (c \ \ll \ 3) \ | \ 1$$

Déballage: Pour déballer un booléen:

$$\text{unwrap}(b) \equiv b \ \gg \ 3$$

Exercice: refaire des exemples similaires en considérant l'utilisation des 3 msb pour le tagging.

Représentation par pointeurs taggés

L'utilisation de cette représentation combinée avec l'emploi des bits les moins significatifs pour le tagging permet d'avoir une arithmétique des petits entiers efficace.

Supposons que n bits de tagging sont employés par mot de 32 bits et que la valeur en décimal du tag des fixnums soit t . Ainsi, la représentation du nombre x est $2^n x + t \bmod 2^{32}$.

Soient x et y deux petits entiers et soient a et b leurs représentations respectives.

Naïvement, la représentation du résultat d'une opération \otimes entre x et y à partir de a et b est obtenue ainsi:

$$\text{wrap}(x \otimes y) \equiv \text{wrap}(\text{unwrap}(a) \otimes \text{unwrap}(b))$$

mais on peut faire mieux:

$$\begin{aligned} \text{wrap}(x + y) &\equiv a + b - t \\ \text{wrap}(x - y) &\equiv a - b + t \\ \text{wrap}(x * y) &\equiv (a - t) * (b \gg n) + t \\ \text{wrap}(x / y) &\equiv \text{wrap}((a - t) / (b - t)) \end{aligned}$$

Noter que ces implantations deviennent particulièrement efficaces si on choisit $t = 0$.

Représentation par pointeurs taggés

Pour mieux utiliser les bits de tagging, on a intérêt à utiliser des tags de longueurs diverses.

En effet, dans un exemple de tagging présenté précédemment, on constate que plusieurs bits sont gaspillés par la représentation (notamment, chez les caractères, les booléens et la liste vide).

De plus, typiquement, les bits de tagging sont rares.

Or, on peut souvent assigner plusieurs types aux mêmes tags.

Exemple:

xx...xxxxx000	→	fixnums	(repr. directe)
xx...xx000001	→	caractères	(repr. directe)
xx...xx001001	→	booléens	(repr. directe)
xx...xx010001	→	liste vide	(repr. directe)
etc.			
xx...xxxxx010	→	paires	(repr. indirecte)
xx...xxxxx011	→	fermetures	(repr. indirecte)
etc.			

Représentation par pointeurs taggés

L'utilisation de $m > 0$ msb et de $n > 0$ lsb à la fois entraîne des coûts supplémentaires.

Par exemple, pour effectuer un test de type sur un objet p , on peut faire quelque chose comme:

$$((p \ll m) \mid (p \gg (32 - m))) \& (2^{m+n} - 1) == tag$$

Toutefois, une implantation directement en assembleur peut profiter d'instructions de rotation qui ne sont pas disponibles en C.

$$(p \text{ rotateleft } m) \& (2^{m+n} - 1) == tag$$

Représentation par partition des codes

L'intervalle $[0, 2^{32} - 1]$ est subdivisé en sous-intervalles, un pour chaque type.

Par exemple:

$[0x00000000, 0x20000000)$ \longrightarrow fixnums

$[0x20000000, 0x60000000)$ \longrightarrow paires

etc.

Les frontières entre les sous-intervalles ne sont pas nécessairement conçues pour mener à des opérations efficaces.

Représentation par partition des codes

Exemples d'implantations des trois opérations liées à l'emballage.

Test de type: Pour tester si on a affaire à une fermeture:

$$\text{pred_clos}(p) \equiv (BASE_{\text{clos}} \leq p) \ \&\& \ (p < BASE_{\text{next}})$$

Emballage: Pour emballer un caractère brut:

$$\text{wrap}(c) \equiv c + BASE_{\text{char}}$$

Déballage: Pour déballer un booléen:

$$\text{unwrap}(b) \equiv b - BASE_{\text{bool}}$$

Représentation par partition des codes

Emplacement des objets alloués: 1 tas versus plusieurs tas.

- S'il y a un seul tas, le sous-intervalle de chacun des types indirects doit être assez grand pour fournir un code pour chaque adresse légale dans le tas.
- S'il y a plusieurs tas, il est possible d'assigner un sous-intervalle plus court aux types moins fréquents et de laisser plus de codes aux types directs. Toutefois, la gestion de plusieurs tas augmente la complexité de l'implantation.

Représentation par partition des codes

L'implantation du test de type peut être optimisée grâce au truc suivant:

$$\begin{aligned} & (BASE_{\text{type}} \leq p) \ \&\& \ (p < BASE_{\text{type}'}) \\ & \quad \quad \quad \equiv \\ & ((\text{unsigned}) \ p - BASE_{\text{type}}) < ((\text{unsigned}) \ (BASE_{\text{type}'} - BASE_{\text{type}})) \end{aligned}$$

qui permet de transformer un test comportant deux sauts conditionnels en un test comportant une soustraction et un saut conditionnel.

Représentation par emplacements typés

Description: le type est donné non pas par la référence à l'objet mais par l'endroit où on a trouvé cette référence.

Cette méthode de typage est normalement applicable chez les langages typés statiquement.

Malgré tout, cette forme de typage existe en Scheme. Les chaînes de caractères ne peuvent contenir que des caractères. On peut donc stocker ceux-ci de façon brute dans les chaînes.

Représentation par large wrappers

Description: la référence à un objet est conservée dans plus d'un champ. Par exemple: un champ contient le type et le champ suivant, la valeur ou l'adresse brute.

Avantages et inconvénients:

- + Fait disparaître les coûts d'emballage et de déballage.
- + Sur les architectures comportant de nombreux registres, le passage de la référence à un objet se fait via deux registres.
- Valable surtout pour des données manipulées directement par le processeur et non celles stockées en mémoire où elles causeraient un gaspillage de mémoire important.

Représentation hybride

Par exemple, un représentation hybride utilisant les pointeurs taggés sur 2 bits et les pointeurs d'objets:

xx...xxxx00	→	fixnums	(repr. directe)	(sur 30 bits)
xx...xxxx01	→	fermetures	(repr. indirecte)	
xx...xxxx10	→	autres types /tas	(repr. indirecte)	(sous-type /tas)
xx...xxx011	→	caractères	(repr. directe)	
xx...xx0111	→	booléens	(repr. directe)	
etc.				

Opérations optimisées diverses

Changement de casse pour des caractères représentés directement:

En minuscules: $c \mid (1 \ll j)$
En majuscules: $c \& (\sim(1 \ll j))$
Alternance maj./min.: $c \hat{=} (1 \ll j)$

Conversions entre les caractères et les entiers:

De caractère à entier: $(c \gg k) \hat{=} PATCH_{CE}$
D'entier à caractère: $(i \ll k) \hat{=} PATCH_{EC}$

Utilité de l'implantation efficace d'un branchement à plusieurs voies sur le type, i.e. sorte de `switch` sur le type:

<code>type_switch (exp)</code>		
<code> case fixnum:</code>		<code>if (pred_fixnum(exp))</code>
<code> { ... }</code>		<code> { ... }</code>
<code> case pair:</code>		<code>else if (pred_pair(exp))</code>
<code> { ... }</code>	qui soit plus rapide que:	<code> { ... }</code>
<code> ...</code>		<code> ...</code>
<code> case string:</code>		<code>else</code>
<code> { ... }</code>		<code> { ... }</code>

Applicabilité des techniques de représentations

Les techniques de représentation des objets et des types ne servent pas qu'aux langages typés dynamiquement. Ils peuvent aussi servir à l'implantation des langages typés statiquement.

Raisons:

- Présence de petits et grands entiers et conversion implicite entre les deux représentations.
- Types algébriques.

Soit le type algébrique suivant:

```
arbre23 a = Empty
          | Deux (arbre23 a) a (arbre23 a)
          | Trois (arbre23 a) a (arbre23 a) a (arbre23 a)
```

Si, par exemple, la variable 'v' est de type `arbre23 int`, on est garanti que sa valeur est une des trois possibilités énumérées. Toutefois, c'est seulement à l'exécution qu'on sait lequel des trois constructeurs est utilisé.

Il faut donc utiliser une des techniques de représentation pour encoder le type et les données contenues dans un objet de type `arbre23 int`.

Choix d'une représentation

Comment choisir une représentation? Il faut considérer:

- l'importance de la vitesse à l'exécution,
- l'importance de l'espace consommé,
- le nombre de bits constants disponibles (bits *taggables*),
- taille des références (ou mots),
- l'ensemble des types à encoder,
- l'importance de la clarté de l'implantation (???) ,
- l'importance de l'extensibilité de l'implantation (reste-t-il beaucoup de bits *taggables?*),
- la fréquence d'utilisation des objets des différents types,
- l'interaction avec le gestionnaire de mémoire,
- etc.