

projet de compilation

Mini Ada

version 3 — 29 décembre 2016

*I believe myself to possess a most singular combination
of qualities exactly fitted to make me pre-eminently
a discoverer of the hidden realities of nature.*
Ada.

L'objectif de ce projet est de réaliser un compilateur pour un fragment du langage Ada, appelé **Mini Ada** par la suite, produisant du code x86-64. Il s'agit d'un fragment relativement petit du langage Ada, avec parfois même quelques petites incompatibilités. Néanmoins, votre compilateur ne sera jamais testé sur des programmes incorrects au sens de **Mini Ada** mais corrects au sens de Ada. Le présent sujet décrit la syntaxe et le typage de **Mini Ada**, ainsi que la nature du travail demandé.

1 Syntaxe

Dans la suite, nous utilisons les notations suivantes dans les grammaires :

$\langle \text{r\`egle} \rangle^*$	répétition de la règle $\langle \text{r\`egle} \rangle$ un nombre quelconque de fois (y compris aucune)
$\langle \text{r\`egle} \rangle_t^*$	répétition de la règle $\langle \text{r\`egle} \rangle$ un nombre quelconque de fois (y compris aucune), les occurrences étant séparées par le terminal t
$\langle \text{r\`egle} \rangle^+$	répétition de la règle $\langle \text{r\`egle} \rangle$ au moins une fois
$\langle \text{r\`egle} \rangle_t^+$	répétition de la règle $\langle \text{r\`egle} \rangle$ au moins une fois, les occurrences étant séparées par le terminal t
$\langle \text{r\`egle} \rangle?$	utilisation optionnelle de la règle $\langle \text{r\`egle} \rangle$ (<i>i.e.</i> 0 ou 1 fois)
(\dots)	parenthésage
$\dots \mid \dots$	alternative

Attention à ne pas confondre « $*$ » et « $^+$ » avec « $*$ » et « $+$ » qui sont des symboles du langage Ada. De même, attention à ne pas confondre les parenthèses avec les terminaux $($ et $)$.

1.1 Conventions lexicales

Espaces, tabulations et retour-chariots constituent les blancs. Un commentaire débute par `--` et s'étend jusqu'à la fin de la ligne. Les identificateurs obéissent à l'expression régulière $\langle \text{ident} \rangle$ suivante :

$$\begin{aligned} \langle \text{chiffre} \rangle &::= 0-9 \\ \langle \text{alpha} \rangle &::= \text{a-z} \mid \text{A-Z} \\ \langle \text{ident} \rangle &::= \langle \text{alpha} \rangle (\langle \text{alpha} \rangle \mid \langle \text{chiffre} \rangle \mid _)^* \end{aligned}$$

Les identificateurs suivants sont des mots clés :

access	and	begin	else	elsif	end
false	for	function	if	in	is
loop	new	not	null	or	out
procedure	record	rem	return	reverse	then
true	type	use	while	with	

Les identificateurs et les mots clés sont insensibles à la casse. Les constantes obéissent aux expressions régulières $\langle \text{entier} \rangle$ et $\langle \text{caractère} \rangle$ suivantes :

$$\begin{aligned} \langle \text{entier} \rangle &::= \langle \text{chiffre} \rangle^+ \\ \langle \text{caractère} \rangle &::= ' \langle \text{tout caractère ASCII 7 bits} \rangle ' \end{aligned}$$

Les constantes entières doivent être comprises entre 0 et 2^{31} , inclus.

1.2 Syntaxe

La grammaire des fichiers sources considérée est donnée figures 1 et 2. Le point d'entrée est le non-terminal $\langle \text{fichier} \rangle$.

$$\begin{aligned} \langle \text{fichier} \rangle &::= \text{with Ada.Text_IO; use Ada.Text_IO;} \\ &\quad \text{procedure } \langle \text{ident} \rangle \text{ is } \langle \text{decl} \rangle^* \\ &\quad \text{begin } \langle \text{instr} \rangle^+ \text{ end } \langle \text{ident} \rangle? ; \text{ EOF} \\ \langle \text{decl} \rangle &::= \text{type } \langle \text{ident} \rangle ; \\ &\quad | \text{type } \langle \text{ident} \rangle \text{ is access } \langle \text{ident} \rangle ; \\ &\quad | \text{type } \langle \text{ident} \rangle \text{ is record } \langle \text{champs} \rangle^+ \text{ end record ;} \\ &\quad | \langle \text{ident} \rangle^+ : \langle \text{type} \rangle (:= \langle \text{expr} \rangle)? ; \\ &\quad | \text{procedure } \langle \text{ident} \rangle \langle \text{params} \rangle? \text{ is } \langle \text{decl} \rangle^* \\ &\quad \quad \text{begin } \langle \text{instr} \rangle^+ \text{ end } \langle \text{ident} \rangle? ; \\ &\quad | \text{function } \langle \text{ident} \rangle \langle \text{params} \rangle? \text{ return } \langle \text{type} \rangle \text{ is } \langle \text{decl} \rangle^* \\ &\quad \quad \text{begin } \langle \text{instr} \rangle^+ \text{ end } \langle \text{ident} \rangle? ; \\ \langle \text{champs} \rangle &::= \langle \text{ident} \rangle^+ : \langle \text{type} \rangle ; \\ \langle \text{type} \rangle &::= \langle \text{ident} \rangle \\ &\quad | \text{access } \langle \text{ident} \rangle \\ \langle \text{params} \rangle &::= (\langle \text{param} \rangle^+ ;) \\ \langle \text{param} \rangle &::= \langle \text{ident} \rangle^+ : \langle \text{mode} \rangle? \langle \text{type} \rangle \\ \langle \text{mode} \rangle &::= \text{in} \mid \text{in out} \end{aligned}$$

FIGURE 1 – Grammaire des fichiers de Mini Ada.

```

⟨expr⟩      ::=  ⟨entier⟩ | ⟨caractère⟩ | true | false | null
              |  ( ⟨expr⟩ )
              |  ⟨accès⟩
              |  ⟨expr⟩ ⟨opérateur⟩ ⟨expr⟩
              |  not ⟨expr⟩ | - ⟨expr⟩
              |  new ⟨ident⟩
              |  ⟨ident⟩ ( ⟨expr⟩+ )
              |  character ' val ( ⟨expr⟩ )
⟨instr⟩     ::=  ⟨accès⟩ := ⟨expr⟩ ;
              |  ⟨ident⟩ ;
              |  ⟨ident⟩ ( ⟨expr⟩+ ) ;
              |  return ⟨expr⟩? ;
              |  begin ⟨instr⟩+ end ;
              |  if ⟨expr⟩ then ⟨instr⟩+ (elsif ⟨expr⟩ then ⟨instr⟩+)*
              |  (else ⟨instr⟩+)? end if ;
              |  for ⟨ident⟩ in reverse? ⟨expr⟩ .. ⟨expr⟩
              |  loop ⟨instr⟩+ end loop ;
              |  while ⟨expr⟩ loop ⟨instr⟩+ end loop ;
⟨opérateur⟩ ::=  = | /= | < | <= | > | >=
              |  + | - | * | / | rem
              |  and | and then | or | or else
⟨accès⟩     ::=  ⟨ident⟩ | ⟨expr⟩ . ⟨ident⟩

```

FIGURE 2 – Grammaire des instructions et expressions de Mini Ada.

Les associativités et précédences des diverses constructions sont données par la table suivante, de la plus faible à la plus forte précedence :

opérateur ou construction	associativité
or, or else	gauche
and, and then	gauche
not	
=, /=	
>, >=, <, <=	
+, -	gauche
- (unaire)	
*, /, rem	gauche
.	gauche

V3

Vérification supplémentaire. Dans une déclaration de procédure ou de fonction, si un identificateur suit le mot-clé **end**, alors il doit être identique au nom de la procédure ou de la fonction déclarée.

2 Typage statique

Une fois l'analyse syntaxique effectuée avec succès, on vérifie la conformité du fichier source. Dans tout ce qui suit, les types sont de la forme suivante :

$$\tau ::= \text{integer} \mid \text{character} \mid \text{boolean} \mid R \mid \text{access } R \mid \text{typenull}$$

où R désigne un type enregistrement. Il s'agit là d'une notation pour la syntaxe *abstraite* des expressions de types. On note $x : \tau \in R$ pour signifier que le type enregistrement R possède un champ x de type τ .

On introduit la relation \equiv sur les types comme la plus petite relation réflexive et symétrique telle que

$$\text{typenull} \equiv \text{access } R$$

Un environnement Γ contient des déclarations de variables, de types enregistrements, de fonctions et de procédures. On notera $d \in \Gamma$ pour spécifier qu'une certaine déclaration d appartient à Γ . Dans un environnement Γ , un type enregistrement R peut être uniquement *déclaré* (avec une déclaration de la forme **type** R ;) ou bien être *défini* (avec une déclaration de la forme **type** R **is record** ... **end record**;). On notera $R \in \Gamma$ pour signifier que R est déclaré ou défini dans Γ .

2.1 Typage des expressions

On introduit le jugement $\Gamma \vdash e : \tau$ signifiant « dans l'environnement Γ , l'expression e est bien typée de type τ ». Ce jugement est défini par les règles d'inférence suivantes :

$$\begin{array}{c} \frac{}{\Gamma \vdash \text{null} : \text{typenull}} \quad \frac{c \text{ constante de type } \tau}{\Gamma \vdash c : \tau} \quad \frac{x : \tau \in \Gamma}{\Gamma \vdash x : \tau} \\[10pt] \frac{\Gamma \vdash e : R \quad x : \tau \in R}{\Gamma \vdash e.x : \tau} \quad \frac{\Gamma \vdash e : \text{access } R \quad x : \tau \in R}{\Gamma \vdash e.x : \tau} \\[10pt] \frac{\Gamma \vdash e_1 : \text{integer} \quad \Gamma \vdash e_2 : \text{integer} \quad op \in \{+, -, *, /, \text{rem}\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{integer}} \\[10pt] \frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \tau_2 \quad op \in \{=, /=\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{boolean}} \\[10pt] \frac{\Gamma \vdash e_1 : \text{integer} \quad \Gamma \vdash e_2 : \text{integer} \quad op \in \{<, <=, >, >=\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{boolean}} \\[10pt] \frac{\Gamma \vdash e_1 : \text{boolean} \quad \Gamma \vdash e_2 : \text{boolean} \quad op \in \{\text{and}, \text{and then}, \text{or}, \text{or else}\}}{\Gamma \vdash e_1 \text{ op } e_2 : \text{boolean}} \\[10pt] \frac{\Gamma \vdash e : \text{integer}}{\Gamma \vdash - e : \text{integer}} \quad \frac{\Gamma \vdash e : \text{boolean}}{\Gamma \vdash \text{not } e : \text{boolean}} \\[10pt] \frac{\text{function } F(x_1 : \tau'_1, \dots, x_n : \tau'_n) \text{ return } \tau \in \Gamma \quad \forall i, \Gamma \vdash e_i : \tau_i \wedge \tau_i \equiv \tau'_i}{\Gamma \vdash F(e_1, \dots, e_n) : \tau} \\[10pt] \frac{\Gamma \vdash e : \text{integer}}{\Gamma \vdash \text{character'val}(e) : \text{character}} \\[10pt] \frac{R \in \Gamma}{\Gamma \vdash \text{new } R : \text{access } R} \end{array}$$

2.2 Typage des instructions

On introduit le jugement $\Gamma \vdash^r i$ signifiant « dans l'environnement Γ , l'instruction i est bien typée, pour un type de retour r ». Un type de retour est soit \perp , pour signifier que l'instruction se trouve dans une procédure, soit un type τ , pour signifier que l'instruction se trouve dans une fonction dont le type de retour est τ . Ce jugement est établi par les règles d'inférence suivantes :

$$\begin{array}{c}
\frac{x : \tau \in \Gamma \quad \Gamma \vdash e' : \tau' \quad \tau \equiv \tau'}{\Gamma \vdash^r x := e} \\
\\
\frac{\Gamma \vdash e : R \quad \textcolor{red}{e \text{ valeur gauche}} \quad x : \tau \in R \quad \Gamma \vdash e' : \tau' \quad \tau \equiv \tau'}{\Gamma \vdash^r e.x := e'} \quad \textcolor{red}{V2} \\
\\
\frac{\Gamma \vdash e : \text{access } R \quad x : \tau \in R \quad \Gamma \vdash e' : \tau' \quad \tau \equiv \tau'}{\Gamma \vdash^r e.x := e'} \\
\\
\frac{\Gamma \vdash e : \text{boolean} \quad \Gamma \vdash^r i_1 \quad \Gamma \vdash^r i_2}{\Gamma \vdash^r \text{if } e \text{ then } i_1 \text{ else } i_2} \\
\\
\frac{\Gamma \vdash e : \text{boolean} \quad \Gamma \vdash^r i}{\Gamma \vdash^r \text{while } e \text{ loop } i \text{ end loop}} \\
\\
\frac{\Gamma \vdash e_1 : \text{integer} \quad \Gamma \vdash e_2 : \text{integer} \quad \Gamma + \textcolor{red}{x : integer} \vdash^r i}{\Gamma \vdash^r \text{for } x \text{ in (reverse) } e_1..e_2 \text{ loop } i \text{ end loop}} \quad \textcolor{red}{V2} \\
\\
\frac{\text{procedure } P(x_1 : \tau'_1, \dots, x_n : \tau'_n) \in \Gamma \quad \forall i, \Gamma \vdash e_i : \tau_i \wedge \tau_i \equiv \tau'_i}{\Gamma \vdash^r P(e_1, \dots, e_n)} \\
\\
\frac{}{\Gamma \vdash^\perp \text{return}} \quad \frac{\Gamma \vdash e : \tau' \quad \tau' \equiv \tau}{\Gamma \vdash^\tau \text{return } e} \\
\\
\frac{\forall j, \Gamma \vdash^r i_j}{\Gamma \vdash^r i_1 \dots i_n} \quad \frac{}{\Gamma \vdash^r \text{new_line}} \quad \frac{\Gamma \vdash e : \text{character}}{\Gamma \vdash^r \text{put}(e)}
\end{array}$$

2.3 Typage des déclarations

On note $\Gamma \vdash \tau$ qui signifie « le type τ est bien formé dans l'environnement Γ ». Ce jugement est défini par les règles suivantes :

$$\begin{array}{c}
\overline{\Gamma \vdash \text{integer}} \quad \overline{\Gamma \vdash \text{character}} \quad \overline{\Gamma \vdash \text{boolean}} \\
\\
\frac{R \text{ est défini dans } \Gamma}{\Gamma \vdash R} \quad \frac{R \in \Gamma}{\Gamma \vdash \text{access } R}
\end{array}$$

On introduit le jugement $\Gamma \vdash d$ qui signifie « dans l'environnement Γ , la déclaration d est bien formée ». Si \vec{d} est une séquence de déclarations d_1, \dots, d_n , on note $\Gamma_0 \vdash \vec{d}$ pour indiquer que $\Gamma_0 \vdash d_1, \Gamma_0 + d_1 \vdash d_2, \dots, \Gamma_0 + d_1 + \dots + d_{n-1} \vdash d_n$. On a les règles suivantes :

$$\begin{array}{c}
\frac{\Gamma \vdash \tau}{\Gamma \vdash x : \tau} \quad \frac{\Gamma \vdash \tau \quad \Gamma \vdash e : \tau' \quad \tau' \equiv \tau}{\Gamma \vdash x : \tau := e} \\
\\
\frac{\forall i, \Gamma \vdash \tau'_i \quad \Gamma' := \Gamma + \textcolor{red}{P} + x_1 : \tau'_1 + \dots + x_n : \tau'_n \quad \Gamma' \vdash \vec{d} \quad \Gamma' + \vec{d} \vdash^\perp i}{\Gamma \vdash \text{procedure } P(x_1 : \tau'_1, \dots, x_n : \tau'_n) \text{ is } \vec{d} \text{ begin } i \text{ end}} \quad \textcolor{red}{V2}
\end{array}$$

$$\begin{array}{c}
\Gamma \vdash \tau \quad \forall i, \Gamma \vdash \tau'_i \\
\hline
\Gamma' := \Gamma + \textcolor{red}{F} + x_1 : \tau'_1 + \dots + x_n : \tau'_n \quad \Gamma' \vdash \vec{d} \quad \Gamma' + \vec{d} \vdash^\tau i \\
\hline
\Gamma \vdash \textbf{function } \textcolor{red}{F}(x_1 : \tau'_1, \dots, x_n : \tau'_n) \textbf{ return } \tau \textbf{ is } \vec{d} \textbf{ begin } i \textbf{ end} \\
\\
\hline
\Gamma \vdash \textbf{type } \overline{R} \\
R \in \Gamma \\
\hline
\Gamma \vdash \textbf{type } x \textbf{ is access } \overline{R} \\
\\
\hline
\forall i, \Gamma + \textbf{type } \textcolor{red}{R} \vdash \tau_i \\
\hline
\Gamma \vdash \textbf{type } R \textbf{ is record } x_1 : \tau_1, \dots, x_n : \tau_n \textbf{ end record}
\end{array}$$

V2

V2

2.4 Vérifications supplémentaires

Enfin, les conditions suivantes doivent être vérifiées :

- Dans un appel de fonction ou de procédure, si un paramètre formel est déclaré **in out**, alors le paramètre effectif correspond doit être une valeur gauche, c'est-à-dire,
 - soit une variable,
 - soit une expression $e.f$ avec e une **valeur gauche** de type R ,
 - soit une expression $e.f$ avec e une expression de type **access** R .
- Dans une fonction ou une procédure, si un paramètre formel X est déclaré **in** (explicitement ou par défaut), alors sa valeur ne peut être modifiée avec une affectation $X := e$ et, si X est de type R , un champ de X ne peut être modifié avec une affectation $X.f := e$.
De même, sa valeur ou la valeur de son champ ne peut être modifiée en étant passée à une fonction ou procédure attendant un paramètre in out.
- La variable d'une boucle **for** ne peut être affectée.
- L'exécution d'une fonction doit nécessairement terminer sur une instruction **return**.
- Toutes les déclarations d'un même niveau doivent porter des noms différents, la seule exception étant celle d'un type enregistrement déclaré puis défini plus loin.
- Un type enregistrement déclaré doit être défini avant la fin des déclarations de même niveau et avant toute introduction d'un niveau de déclarations supérieur.
- Tous les champs d'un même enregistrement doivent porter des noms différents.
- Les noms **put** et **new_line** sont réservés.

V2

V2

2.5 Anticipation

Dans la phase suivante (production de code), certaines informations provenant du typage peuvent être nécessaires. Il vous est conseillé d'anticiper ces besoins en programmant des fonctions de typage qui ne se contentent pas de parcourir les arbres de syntaxe abstraite issus de l'analyse syntaxique mais en renvoient de nouveaux, contenant plus d'information lorsque c'est nécessaire.

3 Production de code

L'objectif est de réaliser un compilateur simple mais correct. En particulier, on ne cherche pas à faire d'allocation de registres mais on se contente d'utiliser la pile pour stocker les éventuels calculs intermédiaires. Bien entendu, il est possible d'utiliser localement les registres. On ne cherche pas à libérer la mémoire.

En cas de doute sur un point de sémantique du langage Ada, on pourra utiliser un compilateur Ada existant (par exemple **gnatmake** en salle info) comme référence.

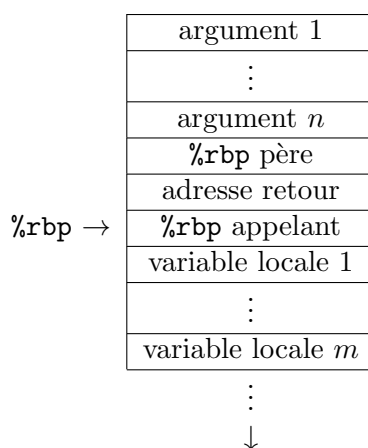
3.1 Représentation des valeurs

Types primitifs. Les entiers sont représentés de manière immédiate par des entiers machine 32 bits signés. Les booléens et les caractères sont représentés par des entiers. En particulier, 0 représente `false` et 1 représente `true`.

La valeur null. Elle est représentée par l'entier 0. En particulier, elle est différente de toute adresse d'un objet alloué sur le tas.

3.2 Schéma de compilation

Les arguments sont passés sur la pile. L'appelant place aussi une référence vers le tableau d'activation de la procédure père. La valeur de retour est passée dans le registre `%rax`. Le tableau d'activation a donc la structure suivante :



4 Travail demandé

Le projet est à faire seul ou en binôme. Il doit être remis par email à `filliatr@lri.fr`. Votre projet doit se présenter sous forme d'une archive `tar` compressée (option "`z`" de `tar`), appelée `vos_noms.tgz` qui doit contenir un répertoire appelé `vos_noms` (exemple : `dupont-durand.tgz`). Dans ce répertoire doivent se trouver les *sources* de votre compilateur (inutile d'inclure les fichiers compilés). Quand on se place dans ce répertoire, la commande `make` doit créer votre compilateur, qui sera appelé `adac`. La commande `make clean` doit effacer tous les fichiers que `make` a engendrés et ne laisser dans le répertoire que les fichiers sources.

L'archive doit également contenir un rapport expliquant les différents choix techniques qui ont été faits et, le cas échéant, les difficultés rencontrées ou les éléments non réalisés. Ce rapport peut être fourni dans un format ASCII, PostScript ou PDF. Ce rapport n'a pas besoin d'être long, mais il doit être informatif.

Première partie, à rendre avant le dimanche 11 décembre 18h. Votre compilateur `adac` doit accepter sur sa ligne de commande les options `--parse-only` et `--type-only` et exactement un fichier Mini Ada portant l'extension `.adb`. Il doit alors réaliser l'analyse syntaxique du fichier. En cas d'erreur lexicale ou syntaxique, celle-ci doit être signalée le plus précisément possible, par sa nature et sa localisation dans le fichier source. On adoptera le format suivant pour cette signalisation :

```
File "test.adb", line 4, characters 5-6:  
syntax error
```

L'anglicisme de la première ligne est nécessaire pour que la fonction `next-error` d'Emacs puisse interpréter la localisation et placer ainsi automatiquement le curseur sur l'emplacement de l'erreur. En revanche, le message d'erreur proprement dit pourra être écrit en français si vous le souhaitez. En cas d'erreur, le compilateur doit terminer avec le code de sortie 1 (`exit 1`).

Si le fichier est syntaxiquement correct, votre compilateur doit terminer avec le code de sortie 0 si l'option `--parse-only` a été passée sur la ligne de commande. Sinon, il doit poursuivre avec le typage du fichier source. Lorsqu'une erreur de typage est détectée par votre compilateur, elle doit être signalée le plus précisément possible, de la manière suivante :

```
File "test.adb", line 4, characters 5-6:  
this expression has type Integer but is expected to have type Boolean
```

Là encore, la nature du message est laissée à votre discrétion, mais la forme de la localisation est imposée. Le compilateur doit alors terminer avec le code de sortie 1 (`exit 1`). Si en revanche il n'y a pas d'erreur de typage, le compilateur doit terminer avec le code de sortie 0. En cas d'erreur du compilateur lui-même, le compilateur doit terminer avec le code de sortie 2 (`exit 2`).

Seconde partie, à rendre avant le dimanche 15 janvier 18h. Étendre votre compilateur pour stopper la compilation après l'analyse sémantique (typage) si l'option de ligne de commande `--type-only` est spécifiée. En l'absence de cette option, et si le fichier d'entrée est conforme à la syntaxe et au typage décrits dans ce sujet, votre compilateur doit produire du code x86-64 et terminer avec le code de sortie 0 (`exit 0` explicite ou terminaison normale du programme), sans rien afficher. Si le fichier d'entrée est `file.adb`, le code assembleur doit être produit dans le fichier `file.s` (même nom que le fichier source mais suffixe `.s` au lieu de `.adb`). Ce fichier x86-64 doit pouvoir être exécuté avec la commande

```
gcc file.s -o file  
./file
```

Le résultat affiché sur la sortie standard doit être identique à celui donné par l'exécution du fichier Ada `file.adb` avec

```
gnatmake file.adb -o file  
./file
```

Remarque importante. La correction du projet est réalisée en partie automatiquement, à l'aide d'un jeu de petits programmes réalisant des affichages avec les procédures `Put` et `New Line`, qui sont compilés avec votre compilateur et dont la sortie est comparée à la sortie attendue. Il est donc très important de correctement compiler les appels à ces deux procédures.

Au delà de ce sujet. Si vous trouvez ce projet trop facile, vous êtes libres d'ajouter à votre compilateur d'autres fragments du langage Ada (en veillant à rester compatible avec les tests fournis). Voici quelques suggestions, non exhaustives : surcharge, tableaux, types dérivés, vérifications dynamiques, exceptions, paquetages. Si vous réalisez une telle extension, merci de le mentionner explicitement dans votre rapport.

Remerciements. Merci à Nathanaël Courant d'avoir relevé des coquilles dans la première version de ce sujet.