

## projet de compilation

**Petit Koka**

version 2 — 29 octobre 2024

L’objectif de ce projet est de réaliser un compilateur pour un fragment de Koka<sup>1</sup>, appelé **Petit Koka** par la suite, produisant du code x86-64. Il s’agit d’un fragment 100% compatible avec Koka, au sens où tout programme de **Petit Koka** est aussi un programme Koka correct. Ceci permettra notamment d’utiliser ce dernier comme référence. Le présent sujet décrit précisément **Petit Koka**, 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)
$( \langle \text{r\`egle} \rangle )$	parenthésage ; attention à ne pas confondre ces parenthèses avec les terminaux ( et )

### 1.1 Conventions lexicales

Espaces et tabulations constituent des blancs. Les commentaires peuvent prendre deux formes :

- débutant par `/*` et s’étendant jusqu’à `*/`, et ne pouvant pas être imbriqués ;
- débutant par `//` et s’étendant jusqu’à la fin de la ligne.

Les identificateurs obéissent à l’expression régulière  $\langle \text{ident} \rangle$  suivante :

$$\begin{aligned}
 \langle \text{digit} \rangle &::= 0-9 \\
 \langle \text{lower} \rangle &::= \text{a-z} \mid \_ \\
 \langle \text{upper} \rangle &::= \text{A-Z} \\
 \langle \text{other} \rangle &::= \langle \text{lower} \rangle \mid \langle \text{upper} \rangle \mid \langle \text{digit} \rangle \mid - \\
 \langle \text{ident} \rangle &::= \langle \text{lower} \rangle \langle \text{other} \rangle^* \text{'*}
 \end{aligned}$$

Par ailleurs, un tiret (`-`) dans un identificateur doit être précédé d’une lettre ou d’un chiffre et suivi d’une lettre s’il n’est pas en dernière position.

Les identificateurs suivants sont des mots clés :

<code>elif</code>	<code>else</code>	<code>fn</code>	<code>fun</code>	<code>if</code>
<code>return</code>	<code>then</code>	<code>val</code>	<code>var</code>	

---

1. <https://koka-lang.github.io/>

Les constantes entières obéissent à l'expression régulière  $\langle integer \rangle$  suivante :

$$\langle integer \rangle ::= -? (0 \mid 1-9 \langle digit \rangle^*)$$

Une chaîne de caractères  $\langle string \rangle$  s'écrit entre guillemets ("). Il y a quatre séquences d'échappement : `\` pour le caractère `"`, `\\` pour le caractère `\`, `\t` pour le caractère de tabulation et `\n` pour un retour chariot. Une chaîne ne peut pas contenir de retour chariot.

## 1.2 Indentation significative

Le langage Koka utilise l'alignement dans le texte source pour définir la structure des blocs. Plus précisément, des lexèmes `{` (début de bloc), `}` (fin de bloc) et `;` (séparation des éléments d'un bloc) sont parfois ajoutés implicitement par l'analyse lexicale. L'analyse syntaxique est ensuite réalisée de manière traditionnelle (voir section suivante). Voici un exemple, avec à gauche le texte source et à droite la suite de lexèmes qui est produite par l'analyseur lexical :

texte source	lexèmes produits
<pre>fun main()   repeat(10)     println("a")     println("b")     println("c")</pre>	<pre> ; fun main() {   repeat(10) {     println("a");     println("b"); } ;     println("c"); } ;</pre>

Explications :

- Le premier retour chariot, à la fin de la première ligne qui est vide, a produit un lexème `;` (et on pourra vérifier que la grammaire s'accommode de points-virgules en début de fichier).
- Le deuxième retour chariot a produit un lexème `{` car le lexème `repeat` est sur une colonne plus grande que le lexème `fun`. De même pour le troisième retour chariot.
- Le quatrième retour chariot (après `println("a")`) a produit un lexème `;` car l'indentation est identique à celle de la ligne précédente.
- Le cinquième retour chariot a produit d'une part les lexèmes `;` et `}`, car l'indentation diminue, puis le lexème `;` car elle retrouve l'indentation de `repeat`.
- Enfin, la fin de fichier a un effet similaire à retrouver une indentation sur la colonne 0, ce qui produit d'abord les lexèmes `;` et `}` puis le lexème `;`, comme au point précédent.

**Définitions.** On introduit deux classes de lexèmes pour les besoins de l'algorithme ci-dessous. On dit qu'un lexème est *une fin de continuation* s'il fait partie de cette liste de lexèmes :

`+ - * / % ++ < <= > >= == != && || ( { ,`

De même, on dit qu'un lexème est *un début de continuation* s'il fait partie de cette liste de lexèmes :

`+ - * / % ++ < <= > >= == != && || then else elif ) } , -> { = . :=`

Ces deux notions sont utilisées pour déterminer si une action doit être effectuée lorsqu'on rencontre un retour chariot.

**Algorithme.** On utilise une *pile* qui contient des entiers. Ce sont les colonnes des blocs d'indentation actuellement en cours. Ces entiers sont strictement croissants, le plus grand étant au sommet de la pile. Initialement, la pile contient l'entier 0.

On lit les lexèmes un par un, avec un outil de son choix, et on retient en permanence le dernier lexème qui a été produit (noté **last** dans la suite). Lorsqu'on rencontre un retour chariot, on lit le prochain lexème **next**, on note sa colonne  $c$  et on effectue les actions suivantes :

- si  $c > m$ , où  $m$  est le sommet de la pile, alors
  1. si **last** n'est pas une fin de continuation et **next** n'est pas un début de continuation, alors émettre { ;
  2. si le dernier lexème émis est {, alors empiler  $c$  ;
  3. émettre **next**.
- si  $c \leq m$ , où  $m$  est le sommet de la pile, alors
  1. tant que  $c < m$ ,
    - (a) dépiler  $m$  et donner à  $m$  la nouvelle valeur du sommet de pile,
    - (b) émettre } si **next**  $\neq$  } ;
  2. si  $c > m$ , échouer ;
  3. sinon, c'est que  $c = m$ , et
    - (a) si **last** n'est pas une fin de continuation et **next** n'est pas un début de continuation, alors émettre ;,
    - (b) émettre **next**.

Pour traiter la fin de fichier correctement, On peut considérer que le lexème EOF est sur la colonne 0 et le traiter dans le cas général (second point ci-dessus). Enfin, toute émission du lexème } est précédée de l'émission d'un lexème ;.

**Suggestion.** Commencer par écrire un analyseur lexical traditionnel (avec un outil de type **lex**), puis construire par dessus un second analyseur lexical qui insère à la volée les lexèmes supplémentaires. On pourra par exemple utiliser une *file* pour stocker les lexèmes en attente d'être envoyés à l'analyseur syntaxique, afin de retrouver le type attendu d'un analyseur lexical qui envoie les lexèmes un par un. De nombreux tests, négatifs et positifs, sont fournis pour aider à comprendre ce mécanisme et à mettre au point l'analyseur lexical.

### 1.3 Syntaxe

La grammaire des fichiers sources considérée est donnée figure 1. Le point d'entrée est le non terminal  $\langle file \rangle$ . 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
<b>if</b>	—
<b>  </b>	gauche
<b>&amp;&amp;</b>	gauche
<b>:=, ==, !=, &gt;, &gt;=, &lt;, &lt;=</b>	—
<b>+, -, ++</b>	gauche
<b>*, /, %</b>	gauche
<b>~, !</b>	—
<b>., {} (postfixe), fn (postfixe)</b>	—

$\langle \text{file} \rangle$	$::=$	<code>;* (⟨decl⟩ ;<sup>+</sup>)* EOF</code>
$\langle \text{decl} \rangle$	$::=$	<code>fun ⟨ident⟩ ⟨funbody⟩</code>
$\langle \text{funbody} \rangle$	$::=$	<code>( ⟨param⟩<sup>*</sup> ) ⟨annot⟩? ⟨expr⟩</code>
$\langle \text{param} \rangle$	$::=$	<code>⟨ident⟩ : ⟨type⟩</code>
$\langle \text{annot} \rangle$	$::=$	<code>: ⟨result⟩</code>
$\langle \text{result} \rangle$	$::=$	<code>(&lt; ⟨ident⟩<sup>*</sup> , &gt;)? ⟨type⟩</code>
$\langle \text{type} \rangle$	$::=$	<code>⟨atype⟩</code>
		<code>  ⟨atype⟩ -&gt; ⟨result⟩</code>
		<code>  ( ⟨type⟩<sup>*</sup> , ) -&gt; ⟨result⟩</code>
$\langle \text{atype} \rangle$	$::=$	<code>⟨ident⟩ (&lt; ⟨type⟩ &gt;)?</code>
		<code>  ( ⟨type⟩ )</code>
		<code>  ( )</code>
$\langle \text{atom} \rangle$	$::=$	<code>True   False   ⟨integer⟩   ⟨string⟩   ( )</code>
		<code>  ⟨ident⟩</code>
		<code>  ( ⟨expr⟩ )</code>
		<code>  ⟨atom⟩ ( ⟨expr⟩<sup>*</sup> , )</code>
		<code>  ⟨atom⟩ . ⟨ident⟩</code>
		<code>  ⟨atom⟩ fn ⟨funbody⟩</code>
		<code>  ⟨atom⟩ ⟨block⟩</code>
		<code>  [ ⟨expr⟩<sup>*</sup> , ]</code>
$\langle \text{expr} \rangle$	$::=$	<code>⟨block⟩</code>
		<code>  ⟨bexpr⟩</code>
$\langle \text{bexpr} \rangle$	$::=$	<code>⟨atom⟩</code>
		<code>  ~ ⟨bexpr⟩</code>
		<code>  ! ⟨bexpr⟩</code>
		<code>  ⟨bexpr⟩ ⟨binop⟩ ⟨bexpr⟩</code>
		<code>  ⟨ident⟩ := ⟨bexpr⟩</code>
		<code>  if ⟨bexpr⟩ then ⟨expr⟩ (elif ⟨bexpr⟩ then ⟨expr⟩)* (else ⟨expr⟩)?</code>
		<code>  if ⟨bexpr⟩ return ⟨expr⟩</code>
		<code>  fn ⟨funbody⟩</code>
		<code>  return ⟨expr⟩</code>
$\langle \text{block} \rangle$	$::=$	<code>{ ;* (⟨stmt⟩ ;<sup>+</sup>)* }</code>
$\langle \text{stmt} \rangle$	$::=$	<code>⟨bexpr⟩</code>
		<code>  val ⟨ident⟩ = ⟨expr⟩</code>
		<code>  var ⟨ident⟩ := ⟨expr⟩</code>
$\langle \text{binop} \rangle$	$::=$	<code>==   !=   &lt;   &lt;=   &gt;   &gt;=   +   -   *   /   %   ++   &amp;&amp;     </code>

FIGURE 1 – Grammaire de Petit Koka.

**Blocs.** Dans un bloc (non terminal  $\langle block \rangle$ ), le dernier élément doit être une expression ( $\langle expr \rangle$ ) et non une déclaration **val** ou **var**.

**Sucre syntaxique.** Certaines constructions n'existent que dans la syntaxe :

syntaxe	traduction
$e . x$	$x(e)$
$e(e_1, \dots, e_n) \text{ fn } f$	$e(e_1, \dots, e_n, \text{fn } f)$
$e \text{ fn } f$	$e(\text{fn } f)$ si $e$ n'est pas un appel
$e(e_1, \dots, e_n) \{ b \}$	$e(e_1, \dots, e_n, \text{fn } () \{ b \})$
$e \{ b \}$	$e(\text{fn } () \{ b \})$ si $e$ n'est pas un appel
$\dots \text{elif} \dots$	$\dots \text{else if} \dots$
$\text{if } e_1 \text{ then } e_2$	$\text{if } e_1 \text{ then } e_2 \text{ else } \{ \}$
$\text{if } e_1 \text{ return } e_2$	$\text{if } e_1 \text{ then return } e_2 \text{ else } \{ \}$

## 2 Typage statique

Une fois l'analyse syntaxique effectuée avec succès, on vérifie la conformité du fichier source. Les types sont de la forme suivante :

$$\begin{aligned} \tau &::= \text{unit} \mid \text{bool} \mid \text{int} \mid \text{string} \mid \text{list}(\tau) \mid \text{maybe}(\tau) \mid (\tau, \dots, \tau) \rightarrow \kappa \\ \kappa &::= \varepsilon / \tau \\ \varepsilon &\subseteq \{ \text{div}, \text{console} \} \end{aligned}$$

Un type  $\tau$  représente un type de valeur et un type  $\kappa$  représente un type de *calcul*, qui combine un type de valeur et un effet  $\varepsilon$ . Ce dernier indique une possible non terminaison (**div**, pour divergence) et/ou une possible écriture sur la sortie standard (**console**).

**Environnement de typage.** L'environnement de typage local, noté  $\Gamma$ , contient une suite ordonnée de déclarations de variables de la forme **val**  $x : \tau$  ou **var**  $x : \tau$ . On note  $\Gamma + \text{val } x : \tau$  l'environnement  $\Gamma$  étendu avec une nouvelle déclaration de variable (même chose avec **var**). Une précédente déclaration de  $x$  dans  $\Gamma$  est, le cas échéant, remplacée par cette nouvelle déclaration.

### 2.1 Typage des expressions

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

$$\begin{array}{c} \frac{c \text{ constante de type } \tau}{\Gamma \vdash c : \emptyset / \tau} \quad \frac{\text{val } x : \tau \in \Gamma}{\Gamma \vdash x : \emptyset / \tau} \quad \frac{\text{var } x : \tau \in \Gamma}{\Gamma \vdash x : \emptyset / \tau} \quad \frac{\text{var } x : \tau \in \Gamma \quad \Gamma \vdash e : \varepsilon / \tau}{\Gamma \vdash x := e : \varepsilon / \text{unit}} \\[10pt] \frac{\Gamma \vdash e : \varepsilon / \text{int}}{\Gamma \vdash \sim e : \varepsilon / \text{int}} \quad \frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{int} \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{int} \quad op \in \{+, -, *, /, \%\}}{\Gamma \vdash e_1 \text{ op } e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{int}} \\[10pt] \frac{\Gamma \vdash e : \varepsilon / \text{bool}}{\Gamma \vdash ! e : \varepsilon / \text{bool}} \quad \frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{int} \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{int} \quad op \in \{<, <=, >, >=\}}{\Gamma \vdash e_1 \text{ op } e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{bool}} \\[10pt] \frac{\Gamma \vdash e_1 : \varepsilon_1 / \tau \quad \Gamma \vdash e_2 : \varepsilon_2 / \tau \quad \tau \in \{\text{bool}, \text{int}, \text{string}\} \quad op \in \{==, !=\}}{\Gamma \vdash e_1 \text{ op } e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{bool}} \\[10pt] \frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{bool} \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{bool} \quad op \in \{\&\&, \mid\mid\}}{\Gamma \vdash e_1 \text{ op } e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{bool}} \end{array}$$

$$\begin{array}{c}
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{string} \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{string}}{\Gamma \vdash e_1 ++ e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{string}} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{list}(\tau) \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{list}(\tau)}{\Gamma \vdash e_1 ++ e_2 : \varepsilon_1 \cup \varepsilon_2 / \text{list}(\tau)} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{bool} \quad \Gamma \vdash e_2 : \varepsilon_2 / \tau \quad \Gamma \vdash e_3 : \varepsilon_3 / \tau}{\Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : \varepsilon_1 \cup \varepsilon_2 \cup \varepsilon_3 / \tau} \\
\frac{\forall i, \Gamma \vdash e_i : \varepsilon_i / \tau}{\Gamma \vdash [e_1, \dots, e_n] : \bigcup \varepsilon_i / \text{list}(\tau)} \\
\frac{}{\Gamma \vdash \{ \} : \emptyset / \text{unit}} \quad \frac{\Gamma \vdash e : \kappa}{\Gamma \vdash \{ e \} : \kappa} \quad \frac{\Gamma \vdash e_1 : \varepsilon_1 / \tau_1 \quad \Gamma \vdash \{ e_2; \dots \} : \varepsilon_2 / \tau_2}{\Gamma \vdash \{ e_1; e_2; \dots \} : \varepsilon_1 \cup \varepsilon_2 / \tau_2} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \tau_1 \quad \Gamma \vdash \text{val } x : \tau_1 \vdash \{ e_2; \dots \} : \varepsilon_2 / \tau_2}{\Gamma \vdash \{ \text{val } x = e_1; e_2 \dots \} : \varepsilon_1 \cup \varepsilon_2 / \tau_2} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \tau_1 \quad \Gamma \vdash \text{var } x : \tau_1 \vdash \{ e_2; \dots \} : \varepsilon_2 / \tau_2}{\Gamma \vdash \{ \text{var } x := e_1; e_2 \dots \} : \varepsilon_1 \cup \varepsilon_2 / \tau_2} \\
\frac{\Gamma \vdash \text{val } x_1 : \tau_1 + \dots + \text{val } x_n : \tau_n \vdash e : \kappa}{\Gamma \vdash \text{fn } (x_1 : \tau_1, \dots, x_n : \tau_n) : \kappa \ e : \emptyset / (x_1 : \tau_1, \dots, x_n : \tau_n) \rightarrow \kappa} \\
\frac{\Gamma \vdash \text{val } x_1 : \tau_1 + \dots + \text{val } x_n : \tau_n \vdash e : \kappa}{\Gamma \vdash \text{fn } (x_1 : \tau_1, \dots, x_n : \tau_n) \ e : \emptyset / (x_1 : \tau_1, \dots, x_n : \tau_n) \rightarrow \kappa} \\
\frac{\Gamma \vdash e : \varepsilon / ((\tau_1, \dots, \tau_n)) \rightarrow \varepsilon_f / \tau \quad \forall i, \Gamma \vdash e_i : \varepsilon_i / \tau_i}{\Gamma \vdash e(e_1, \dots, e_n) : \varepsilon \cup \varepsilon_f \cup \bigcup \varepsilon_i / \tau} \quad \frac{\Gamma \vdash e_1 : \varepsilon_1 / \tau_1}{\Gamma \vdash \text{return } e_1 : \varepsilon_1 / \tau}
\end{array}$$

### Typage des fonctions prédéfinies.

$$\begin{array}{c}
\frac{\Gamma \vdash e : \varepsilon / \tau \quad \tau \in \{\text{unit}, \text{bool}, \text{int}, \text{string}\}}{\Gamma \vdash \text{println}(e) : \varepsilon \cup \{\text{console}\} / \text{unit}} \quad \frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{maybe}(\tau) \quad \Gamma \vdash e_2 : \varepsilon_2 / \tau}{\Gamma \vdash \text{default}(e_1, e_2) : \varepsilon_1 \cup \varepsilon_2 / \tau} \\
\frac{\Gamma \vdash e : \varepsilon / \text{list}(\tau)}{\Gamma \vdash \text{head}(e) : \varepsilon / \text{maybe}(\tau)} \quad \frac{\Gamma \vdash e : \varepsilon / \text{list}(\tau)}{\Gamma \vdash \text{tail}(e) : \varepsilon / \text{list}(\tau)} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{int} \quad \Gamma \vdash e_2 : \varepsilon_2 / \text{int} \quad \Gamma \vdash e_3 : \varepsilon_3 / ((\text{int}) \rightarrow \varepsilon / \text{unit})}{\Gamma \vdash \text{for}(e_1, e_2, e_3) : \varepsilon_1 \cup \varepsilon_2 \cup \varepsilon_3 \cup \varepsilon / \text{unit}} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / \text{int} \quad \Gamma \vdash e_2 : \varepsilon_2 / (() \rightarrow \varepsilon / \text{unit})}{\Gamma \vdash \text{repeat}(e_1, e_2) : \varepsilon_1 \cup \varepsilon_2 \cup \varepsilon / \text{unit}} \\
\frac{\Gamma \vdash e_1 : \varepsilon_1 / (() \rightarrow \varepsilon_3 / \text{bool}) \quad \Gamma \vdash e_2 : \varepsilon_2 / (() \rightarrow \varepsilon_4 / \text{unit})}{\Gamma \vdash \text{while}(e_1, e_2) : \varepsilon_1 \cup \varepsilon_2 \cup \varepsilon_3 \cup \varepsilon_4 \cup \{\text{div}\} / \text{unit}}
\end{array}$$

## 2.2 Typage d'un fichier

Les déclarations constituant un fichier sont traitées dans l'ordre d'apparition. Les fonctions sont ajoutées à l'environnement au fur et à mesure de leur analyse. Une fonction ne peut faire référence qu'à des fonctions antérieures (ou à elle-même). Chaque fonction ne peut être définie qu'une seule fois.

**Déclaration de fonction.** Une fonction  $f$  est introduite par une déclaration de la forme suivante

$$\text{fun } f(x_1 : \tau_1, \dots, x_n : \tau_n) (: \kappa)? e$$

avec un type de retour  $\kappa$  optionnel et un corps de fonction qui est une expression  $e$ . Les paramètres  $x_i$  doivent porter des noms distincts. On vérifie que le corps  $e$  est bien typé dans l'environnement  $\Gamma$  des fonctions précédentes étendu avec  $f$  et ses paramètres, c'est-à-dire

$$\Gamma + \text{val } f : ((\tau_1, \dots, \tau_n) \rightarrow \kappa) + \text{val } x_1 : \tau_1 + \dots + \text{val } x_n : \tau_n \vdash e : \kappa$$

où  $\kappa$  est le type donné par l'utilisateur, le cas échéant, ou sinon un type de calcul d'effet minimal (au sens de l'inclusion). Si le corps  $e$  de la fonction fait référence à la fonction  $f$  (dans un appel récursif ou même en passant  $f$  à une autre fonction), alors  $\kappa$  doit inclure l'effet `div`. Enfin, il convient de vérifier que toute expression `return`  $e_1$  apparaissant dans le corps de la fonction renvoie bien une valeur du type attendu par  $\kappa$ .

**Fonction main.** Le fichier source doit contenir une fonction `main` sans paramètre. Son type de retour est libre.

## 2.3 Limitations par rapport à Koka

Le langage Petit Koka souffre d'un certain nombre de limitations par rapport à Koka. En particulier,

- Petit Koka possède moins de mots clés que Koka ;
- Koka autorise à ne pas indiquer le type d'un paramètre de fonction (et l'infère).

Cependant, votre compilateur ne sera jamais testé sur des programmes incorrects au sens de Petit Koka mais corrects au sens de Koka.

## 2.4 Indications

Il est fortement conseillé de procéder construction par construction, en compilant et testant systématiquement son projet à chaque étape. De nombreux tests sont fournis sur la page du cours, avec un script pour lancer votre compilateur sur ces tests.

En cas de doute concernant un point de sémantique, vous pouvez utiliser le compilateur Koka comme référence. Vous pouvez d'ailleurs vous inspirer de ses messages d'erreur pour votre compilateur (en les traduisant ou non en français).

Dans la phase suivante (production de code), certaines informations provenant du typage seront 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 notamment des types. Ces nouveaux arbres de syntaxe abstraite peuvent différer plus ou moins des arbres issus de l'analyse syntaxique.

### 3 Production de code

L'objectif est de réaliser un compilateur simple mais correct. En particulier, on ne cherchera pas à faire d'allocation de registres mais on se contentera d'utiliser la pile pour stocker les éventuels calculs intermédiaires. Bien entendu, il est possible, et même souhaitable, d'utiliser localement les registres de x86-64. On ne cherchera pas à libérer la mémoire. On propose ici un schéma de compilation à titre indicatif. Vous êtes libres de faire tout autre choix.

#### 3.1 Représentation des valeurs

Les entiers sont représentés de manière immédiate par des entiers 64 bits signés (même si le type `int` de Koka est un type d'entiers de précision arbitraire). La valeur `()` est représentée par l'entier 0. Les booléens sont représentés par des entiers, avec la même convention que x86-64 : 0 représente `False` et 1 représente `True`. Une chaîne est une adresse vers un bloc alloué sur le tas contenant une chaîne terminée par un caractère 0 (comme en C). Une liste est soit l'entier 0 (la liste vide), soit une adresse vers un bloc de deux mots alloué sur le tas contenant la valeur en tête de liste et le reste de la liste. Une valeur de type `maybe< $\tau$ >` est représentée soit par l'entier 0 (rien), soit par un pointeur vers un bloc d'un mot alloué sur le tas contenant une valeur de type  $\tau$ . Enfin, une valeur de type fonction est représentée par une fermeture, c'est-à-dire une adresse vers un bloc de  $n + 1$  mots alloué sur le tas de la forme

code	$v_1$	$v_2$	$\dots$	$v_n$
------	-------	-------	---------	-------

où le premier mot contient un pointeur de code et les mots suivants contiennent **les variables  $v_1, \dots, v_n$  capturées dans la fermeture.**

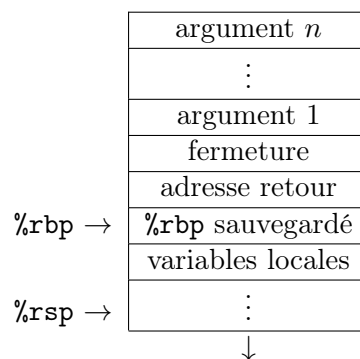
#### 3.2 Schéma de compilation

Chaque fonction de Petit Koka est compilée vers une fonction assembleur qui prend tous ses arguments sur la pile et place sa valeur de retour dans le registre `%rax`. Si la fonction est une fermeture, celle-ci est passée comme un tout premier paramètre. Cela simplifie les choses de considérer toute fonction comme une fermeture, y compris les fonctions globales.

On conseille fortement de procéder en deux temps, en commençant par une explicitation des fermetures (cf cours 8). Cela étant, il reste possible de commencer par le fragment sans aucune fonction anonyme, dans un schéma de compilation plus simple (cf cours 7) et de ne rajouter les fonctions anonymes que dans un second temps. Attention toutefois aux constructions `for/repeat/while` qui sont en Koka des fonctions d'ordre supérieur. Il faut alors soit les traiter de façon particulière, soit attendre d'avoir ajouté le traitement des fermetures.

On pourra écrire quelques fonctions directement en assembleur (concaténation de chaînes et de listes, instances de `println`, etc.) et ajouter cela au code produit par le compilateur. Attention, avant d'appeler des fonctions de bibliothèque C comme `printf` ou `malloc`, il convient d'aligner la pile, c'est-à-dire de garantir que `%rsp` est un multiple de 16 avant de faire `call`. Le plus simple pour cela consiste à définir de petites fonctions assembleur qui appellent les fonctions de bibliothèque après avoir correctement aligné la pile, par exemple comme ceci :

```
my_malloc:
    pushq    %rbp
```





```

movq    %rsp, %rbp
andq    $-16, %rsp    # alignement de la pile
movq    24(%rbp), %rdi # argument de malloc ici passé sur la pile
call    malloc
movq    %rbp, %rsp
popq    %rbp
ret

```

Si on adopte cette approche, il n'est en revanche pas nécessaire d'aligner la pile pour ses propres fonctions.

## 4 Travail demandé

Le projet est à faire seul ou en binôme. Il doit être remis par email<sup>2</sup>, sous la forme d'une archive compressée (avec `tar` ou `zip`), appelée `vos_noms.tgz` ou `vos_noms.zip` qui doit contenir un répertoire appelé `vos_noms` (exemple : `dupont-durand`). Dans ce répertoire doivent se trouver les *sources* du compilateur (inutile d'inclure les fichiers compilés). Quand on se place dans ce répertoire, la commande `make` doit créer le compilateur, qui sera appelé `kokac`. La commande `make clean` doit effacer tous les fichiers que `make` a produits et ne laisser dans le répertoire que les fichiers sources. Bien entendu, la compilation du projet peut être réalisée avec d'autres outils que `make` (par exemple `dune` si le projet est réalisé en OCaml) et le `Makefile` se réduit alors à quelques lignes seulement pour appeler ces outils.

L'archive doit également contenir un court rapport expliquant les différents choix techniques qui ont été faits et, le cas échéant, les difficultés rencontrées ou les éléments non réalisés. Ce rapport pourra être fourni dans un format texte (par exemple Markdown) ou PDF.

**Partie 1 (à rendre pour le dimanche 15 décembre 18:00).** Dans cette première partie du projet, le compilateur `kokac` doit accepter sur sa ligne de commande une option éventuelle (parmi `--parse-only` et `--type-only`) et exactement un fichier Petit Koka portant l'extension `.koka`. 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.koka", line 4, characters 5-6:
syntax error

```

L'anglicisme de la première ligne est nécessaire pour que la fonction `next-error` d'Emacs<sup>3</sup> 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, le 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 le compilateur, elle doit être signalée le plus précisément possible, de la manière suivante :

```

File "test.koka", line 4, characters 5-6:
this expression has type int but is expected to have type string

```

---

2. à [Jean-Christophe.Filliatre@cnrs.fr](mailto:Jean-Christophe.Filliatre@cnrs.fr)

3. Le correcteur est susceptible d'utiliser cet éditeur.

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`). L’option `--type-only` indique de stopper la compilation après l’analyse sémantique (typage). Elle est donc sans effet dans cette première partie, mais elle doit être honorée néanmoins.

**Partie 2 (à rendre pour le dimanche 19 janvier 18:00).** Si le fichier d’entrée est conforme à la syntaxe et au typage de ce sujet, votre compilateur doit produire du code x86-64 et terminer avec le code de sortie 0, sans rien afficher. Si le fichier d’entrée est `file.koka`, le code assembleur doit être produit dans le fichier `file.s` (même nom que le fichier source mais suffixe `.s` au lieu de `.koka`). Ce fichier x86-64 doit pouvoir être exécuté avec la commande

```
gcc -no-pie file.s -o file
./file
```

Le résultat affiché sur la sortie standard doit être identique à celui donné par la commande suivante :

```
koka --console=raw -v0 -e file.koka
```

**Remarque importante.** La correction du projet est réalisée en partie automatiquement, à l’aide d’un jeu de petits programmes réalisant des affichages avec l’instruction `println`, qui sont compilés avec votre compilateur et dont la sortie est comparée à la sortie attendue. (Voir les tests fournis dans le sous-répertoire `exec/`.) Il est donc très important de correctement compiler les appels à `println`.

**Conseils.** Il est fortement conseillé de procéder construction par construction que ce soit pour le typage ou pour la production de code, dans cet ordre : affichage (`println`), arithmétique, variables locales, `if`, fonctions, fonctions anonymes.