

# テキスト処理言語における文字列のための正規表現型

田 淵 直<sup>†</sup> 住井 英二郎<sup>†</sup> 米 澤 明 憲<sup>†</sup>

Perl, Ruby, Python などに代表されるスクリプト言語は, CGI のような文字列操作が中心となるアプリケーションの開発に広く使用されており, また開発効率を高めるために有用であることが経験的に知られている. しかしながら, これらの言語はその柔軟すぎる記述力のため, プログラムの挙動が実際の実行時まで分からないことが多く, デバッグ・保守が困難になりがちであるという問題をかかえている. 一方, プログラムの性質を静的に検証するという目的では, 型システムを利用することが有効であるが, Java・ML など既存の言語の型システムは, スクリプト言語のプログラミングスタイルと必ずしも相性が良くないという問題がある. 本研究では, これらのスクリプト言語のプログラミングスタイルとして「正規表現を用いた文字列操作」が多用されるという点に着目し, 正規表現を一種の「型」と見なして文字列に対する型付けを行う型システムを提案する. 我々の手法は, 文字列の内容や, 文字列操作の結果を静的に検査し, プログラムの信頼性を向上させることができる. また, 正規表現を用いることにより, 古典的な型システムでは達成できない記述力と柔軟性を得ることができる.

## Regular Expression Types for Strings in a Text Processing Language

NAOSHI TABUCHI,<sup>†</sup> EIJIRO SUMII<sup>†</sup> and AKINORI YONEZAWA<sup>†</sup>

Scripting language such as Perl, Ruby, Python, etc. are now widely used in developing string-manipulating software including CGIs. Although they are known to be useful for rapid application development, their excessive flexibility often makes it difficult to debug and maintain programs. On the other hand, type systems are known to be a good tool for static verification of programs. However, existing type systems as in Java, ML, etc. are not necessarily compatible with scripting languages. To overcome this problem, we propose a new type system in which “types” consist of regular expressions, observing that an important feature of these scripting languages is to manipulate string values by means of regular expressions. Our method enables programmers to statically verify string values and operations on them, and thereby improve reliability of software. Using regular expressions, our approach achieves higher degree of expressiveness and flexibility than traditional type systems.

### 1. はじめに

#### 背 景

Perl<sup>13)</sup>, Ruby<sup>9),12)</sup>, Python<sup>8)</sup> などのスクリプト言語は, いわゆる RAD (Rapid Application Development) のための強力なツールとして, 特に web プログラミングのような文字列・テキスト処理が中心となる開発ドメインにおいて広く受け入れられている. しかしながら, 確かにこれらの言語ではプログラムを書くことそのものは簡単になるものの, 一度書いたプログラムのデバッグ・保守という観点から見ると, いくつかの問題点をかかえている. スクリプト言語を使っ

た開発の典型的な工程では (1) スクリプト言語の記述力を活かした高速なコーディングの後 (2) プログラムを実行させると実行時エラーに遭遇し (運が悪ければ何の出力も得られないこともある) (3) 苦痛のともなうデバッグ作業に多大な時間を費やす, というステップをたどることになる. さらに悪いことには, テストの段階で見落とされた境界条件が原因となってプログラムの運用開始後に初めてバグが発見されるというケースも多い. Web のようなオープンな環境において, この種のバグは致命的なセキュリティホールとなりうる. よく知られた例では, HTML のメタ文字をエスケープし忘れることによる CGI プログラムのクロスサイトスクリプティング問題<sup>1)</sup>などをあげることができる.

一方, プログラムの実行時エラーを防止するための有効なアプローチとしては静的型システムによる型付

<sup>†</sup> 東京大学情報理工学系研究科コンピュータ科学専攻  
Department of Computer Science, Graduate School of  
Information Science and Technology, The University of  
Tokyo

けがある。しかし、MLに見られるような標準的な型システムは必ずしも上述の問題に対する有効な解決策とはならない。なぜなら、これらの型システムはルールが厳格にすぎたため古いプログラム（あるいは古いプログラマ）との相性が必ずしも良くないからである。さらに標準的な型システムでは、テキストデータは単に string という型のみが与えられ、それ以上の情報を得ることができないという点も重要である。入力テキストを解析して、より構造化されたデータの形にして操作するという方法ももちろん考えられるが、それではスクリプト言語の簡便さというメリットを損なってしまう。

#### 我々のアプローチ

以上のような手軽さと頑健さのジレンマを克服する手段として、我々は文字列の型を正規表現として表した文字列のための正規表現型を提案する。たとえば、 $a$ ,  $aa$ ,  $aaa$  といった文字列定数は皆  $a^*$  という型を付けることができる。他にも空文字列（長さ 0 の文字列）を  $\varepsilon$ 、文字列の連結操作を  $\wedge$  と書くことにすると、 $f(n) = \text{if } n \leq 0 \text{ then } \varepsilon \text{ else } a \wedge f(n-1) \wedge b$  という再帰関数には  $\text{int} \rightarrow a^*b^*$  という型を付けることができるだろう。プログラムがより複雑になってくると、このような文字列の形式に関する情報はデバッグ・保守の有用な手掛かりとなることが期待される。

ここで、なぜ文脈自由文法など他のクラスの形式言語ではなく、正規表現を採用するのかという疑問が生じるかもしれない。正規表現の代わりに“文脈自由文法型”なるものを採用すれば、上の例での関数  $f$  の値域は  $\{a^n b^n \mid n = 0, 1, 2, \dots\}$  のようにより正確に表すことができる。我々が正規表現を選んだ 1 つの理由は、その構文・意味論的な基礎が確立されており、しかもプログラマにとって馴染みの深いものだからである。より重要な理由として、正規表現が交差・和・差・商など多くの基本的な演算に関して閉じているということがある。そして、等価性・包含など最も基本的な判定問題に関しても正規表現ならば決定可能である。これらの特質は我々の型システムにおいて、効率的な型チェックの実現のために不可欠である。他のクラスの形式言語はこれほど扱いやすい特徴を持っていない。たとえば広く知られているように文脈自由言語の等価性（それゆえ包含も）の判定は決定不能である<sup>4)</sup>。

#### 我々の貢献

上述のアイデアを形式的、かつ簡易に示すための必要最小限の計算体系として我々は言語  $\lambda^{re}$ 、およびその上での文字列の正規表現型を定義する。 $\lambda^{re}$  はきわめて小さな体系であるが、正規表現の持つ柔軟性によ

り、得られる結果の多くはより現実的な言語に拡張できると期待される。たとえば、文字列に評価されるあらゆる式は、最悪でも  $*$ （“すべての文字列”）として型付けできるし、多くの場合にはもっと詳細にできるだろう。

データの形やプログラムの挙動を表現するために正規表現、あるいは他の形式言語を用いるのはまったく新しいアイデアではない。我々の知る限り最も関連の深い研究として、XML 文書処理のための言語である XDuce<sup>5),6)</sup> がある。XDuce では XML 文書をラベル付きの木と見なし、それに対する正規表現型を与えている。理論的に見れば文字列は木の特別な場合であるから（Scheme でのリストの cons エンコーディングを考えればよい）、XDuce は本論文の内容を包括していると考えられるかもしれない。しかしながら、本論文は以下の点で独立の価値があると我々は考える。

- (1) 一般的な木構造ではなく文字列に焦点を当てているため、技術的な詳細が簡略化され、問題点を明確にできる
- (2) 我々の研究は XDuce の結果を一部、拡張している。詳細は後の章に譲るが、主要な拡張は (i) 副作用の正規表現による型付けと (ii) 末尾以外の場所での  $as$  パターンに対する型推論である。

本論文では型推論の完全なアルゴリズムや効率的な実装については言及していない（ただし 4 章でパターンマッチの変数束縛に関する部分的な型推論を、6 章で完全な型推論に向けての大まかなアイデアを紹介する）。プログラマが大きな負担を払わずに型システムの恩恵を受けるために、型推論は重要な要素ではあるが、それ自身で 1 つの大きなテーマとなる問題であるため、今回は将来の課題とした。

本論文の以降の構成は次のとおりである。2 章で  $\lambda^{re}$  の構文を、続いて 3 章で操作的意味論を定義する。4 章で型システムを定義し、5 章では  $\lambda^{re}$  によるプログラムの例を紹介する。6 章で関連研究と今後の課題についての考察を与え、結論とする。

なお、紙幅の都合により証明の詳細な部分については省略した。これに関しては WWW で <http://www.yl.is.s.u-tokyo.ac.jp/~tabee/xper1/> から参照することができる。

## 2. $\lambda^{re}$ の構文

$\lambda^{re}$  の項の構文は図 1 で定義される  $.x, y, z, f, g, h$  などの変数名であり、これらは加算無限個存在するとする。 $\text{fix}(f, x, M)$  は  $f(x) = M$  で定義される再帰関数  $f$  を表すものとし、これを含めて最初の 3 つの要素は標準的なものである。 $\lambda x. M$  を  $f$  が  $M$  の自由変数に表

$M$ (term)	::=	$x$	(variable)
		$\text{fix}(f, x, M)$	(recursive function)
		$M_1 M_2$	(function application)
		$s$	(constant string)
		$M_1 \hat{\ } M_2$	(string concatenation)
		$\text{match } M \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n$	(pattern matching)
		$\text{print } M$	(output)
$P$ (pattern)	::=	$s$	(constant string)
		$P_1 \mid P_2$	(choice)
		$x \text{ as } P$	(variable binding)
		$P^*$	(repetition)
		$P_1 P_2$	(sequence)

図 1 項とパターンの構文

Fig. 1 Syntax of terms and patterns.

れない場合の  $\text{fix}(f, x, M)$  の略記,  $\text{let } x = M_1 \text{ in } M_2$  を  $(\lambda x. M_2)M_1$  の略記と定める。ただし, 自由変数の定義についてはやはり標準的なものに従うこととし, ここでは省略する。

文字列を操作するためのプリミティブとして, 以下の4つの要素を導入する。(1)文字列そのものを表す文字列定数  $s$ , (2)2つの文字列を連結して新しい文字列を作る  $M_1 \hat{\ } M_2$ , (3)文字列に対して正規表現のパターンマッチを行う  $\text{match } M \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n$ , (4)そして文字列を出力して副作用を引き起こす  $\text{print } M$  である。4章で見ると, これが我々の正規表現 effect system の対象となる。

パターンの構文はほぼ標準的な正規表現に準ずるが, 変数束縛のための  $x \text{ as } P$  という構文を導入して拡張している。これは ML における  $\text{as}$  パターンと同様に, 最初に入力文字列  $s$  をパターン  $P$  にマッチさせ, 成功すれば  $x$  に  $s$  を束縛する。以降,  $P$  に現れるすべての変数の集合を表すのに  $\text{var}(P)$  と書くこととする。

### 3. $\lambda^{re}$ の操作的意味論

$\lambda^{re}$  の項の意味は図2にあるとおり, small-step の簡約意味論として定義される。Big-step ではなく small-step の意味論を用いるのは, 発散するプログラムについても副作用に関しては停止するプログラムと同等に扱えるようにするためである。意味論は  $M_1 \xrightarrow{s} M_2$  の形の関係で表され, これは“項  $M_1$  は項  $M_2$  に簡約され, その際文字列  $s$  を出力する”と読むことができる。副作用の扱いを簡単にするため, 評価順序は call-by-value かつ left-to-right で固定する。

規則 (R-App), (R-Cat) および (R-Ctx) はそれぞ

れ標準的な(再帰的かもしれない)関数適用, 文字列結合, および評価コンテキストを表している。(R-Print) は  $\text{print } s$  という項を評価すると文字列  $s$  を出力し, ダミーの値(空文字列  $\varepsilon$ )に簡約されることを表している。(R-Match)では  $\lambda^{re}$  におけるパターンマッチの意味を,  $s \triangleright P \Rightarrow \theta$  という関係を使って定義している。直観的には  $s \triangleright P \Rightarrow \theta$  という関係は, 文字列  $s$  をパターン  $P$  にマッチさせると, 失敗するか( $\theta = \perp$  の場合)変数から文字列への置換  $\theta$  を生じる( $\theta \neq \perp$  の場合)ことを意味する。すなわち, (R-Match)は  $\text{match } s \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n$  という項が  $\theta M_n$  に簡約され, ただし  $P_m$  は文字列  $s$  にマッチする(そして置換  $\theta$  を生む)最初のパターンであることをいっている。

パターンマッチの関係  $s \triangleright P \Rightarrow \theta$  は図3で定義されている。以降,  $s \triangleright P \Rightarrow \theta$  なる  $\theta (\neq \perp)$  が存在するとき単に  $s \triangleright P$  と書くこととする。直観的には, 図3の規則はパターンマッチの(素朴な)アルゴリズムをボトムアップの形で記述したものに相当し, 基本的にはパターン  $P$  の構造に関して帰納的に定義されている。(M-Const-Succ)と(M-Const-Fail)はパターンが文字列定数  $s$  の場合の規則, (M-Choice-Fst)と(M-Choice-Snd)はパターンが  $P_1 \mid P_2$  という選択の形の場合の規則, (M-Bind)は変数束縛  $x \text{ as } P$  のための規則, そして(M-Rep)はパターンが繰返し  $P^*$  の場合の規則である。(M-Bind)において  $\theta \uplus \{x \mapsto s\}$  は  $\theta'(x) = s, \theta'(y) = \theta(y) (y \neq x)$  を満たす  $\theta'$  を表すこととする(ただし  $\perp \uplus \{x \mapsto s\} = \perp$  と定める)。また,  $\emptyset$  は恒等置換を表すとする。

注意すべき点として, (M-Choice-Snd)規則は(M-

$v$ (value)	$::= \text{fix}(f, x, M)$ $  s$
$C[]$ (evaluation context)	$::= []M$ $  v[]$ $  []^M$ $  v^[]$ $  \text{match } [] \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n$ $  \text{print } []$
$\frac{}{\text{fix}(f, x, M)v \xrightarrow{\varepsilon} [v/x][\text{fix}(f, x, M)/f]M} \text{(R-App)} \quad \frac{}{s_1 \wedge s_2 \xrightarrow{\varepsilon} s_1 + s_2} \text{(R-Cat)}$	
$\frac{}{\text{print } s \xrightarrow{s} \varepsilon} \text{(R-Print)} \quad \frac{\forall i < m. s \triangleright P_i \Rightarrow \perp \quad s \triangleright P_m \Rightarrow \theta \quad \theta \neq \perp}{\text{match } s \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n \xrightarrow{\varepsilon} \theta M_m} \text{(R-Match)}$	
$\frac{M_1 \xrightarrow{s} M_2}{C[M_1] \xrightarrow{s} C[M_2]} \text{(R-Ctx)}$	

図 2 項の操作的意味

Fig. 2 Operational semantics of terms.

$\frac{}{s \triangleright s \Rightarrow \emptyset} \text{(M-Const-Succ)} \quad \frac{s \neq s'}{s \triangleright s' \Rightarrow \perp} \text{(M-Const-Fail)}$	$\frac{s \triangleright P_1 \Rightarrow \theta \quad \theta \neq \perp}{s \triangleright P_1 \mid P_2 \Rightarrow \theta} \text{(M-Choice-Fst)} \quad \frac{s \triangleright P_1 \Rightarrow \perp \quad s \triangleright P_2 \Rightarrow \theta}{s \triangleright P_1 \mid P_2 \Rightarrow \theta} \text{(M-Choice-Snd)}$
$\frac{s \triangleright P \Rightarrow \theta}{s \triangleright x \text{ as } P \Rightarrow \theta \uplus \{x \mapsto s\}} \text{(M-Bind)}$	$\frac{s \triangleright PP^* \mid \varepsilon \Rightarrow \theta}{s \triangleright P^* \Rightarrow \theta} \text{(M-Rep)}$
$\frac{s_2^{-1}s_1 \triangleright P \Rightarrow \theta}{s_1 \triangleright s_2P \Rightarrow \theta} \text{(M-Seq-Const-Succ)}$	$\frac{s_2^{-1}s_1 \text{ not exist}}{s_1 \triangleright s_2P \Rightarrow \perp} \text{(M-Seq-Const-Fail)}$
$\frac{s \triangleright P_1P_3 \mid P_2P_3 \Rightarrow \theta}{s \triangleright (P_1 \mid P_2)P_3 \Rightarrow \theta} \text{(M-Seq-Choice)}$	
$\frac{y \notin \text{var}(P_1P_2) \quad s_1 \triangleright P_1(y \text{ as } P_2) \Rightarrow \theta \uplus \{y \mapsto s_2\}}{s_1 \triangleright (x \text{ as } P_1)P_2 \Rightarrow \theta \uplus \{x \mapsto s_1s_2^{-1}\}} \text{(M-Seq-Bind)}$	
$\frac{s \triangleright (P_1P_1^* \mid \varepsilon)P_2 \Rightarrow \theta}{s \triangleright P_1^*P_2 \Rightarrow \theta} \text{(M-Seq-Rep)}$	$\frac{s \triangleright P_1(P_2P_3) \Rightarrow \theta}{s \triangleright (P_1P_2)P_3 \Rightarrow \theta} \text{(M-Seq-Seq)}$

図 3 パターンの操作的意味

Fig. 3 Operational semantics of patterns.

Choice-Fst) 規則が適用されなかった場合のみ適用される。これによって、いわゆる *first-match* を実現できる。同様に (M-Rep) 規則は  $P^*$  を  $\varepsilon \mid PP^*$  ではなく  $PP^* \mid \varepsilon$  に展開する。これと *first-match* の原則

からいわゆる *longest-match* を実現できる。

パターンが  $P_1P_2$  という連接の形をしている場合、(M-Seq-...) 規則に見られるように  $P_1$  の形に従ってさらに 6 つの場合に分ける。入力文字列  $s$  を、単に

$\frac{}{\varepsilon \triangleright a(a^*(y \text{ as } a^*)) \Rightarrow \perp} \text{ (M-Seq-Const-Fail)}$ $\frac{}{\varepsilon \triangleright (aa^*)(y \text{ as } a^*) \Rightarrow \perp} \text{ (M-Seq-Seq)}$ $\frac{\varepsilon \triangleright (aa^*)(y \text{ as } a^*) \mid \varepsilon(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon}{\varepsilon \triangleright (aa^* \mid \varepsilon)(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Seq-Choice)}$ $\frac{\varepsilon \triangleright a^*(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon}{\varepsilon \triangleright a^*(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Seq-Rep)}$ $\frac{a \triangleright a(a^*(y \text{ as } a^*)) \Rightarrow y \mapsto \varepsilon}{a \triangleright a(a^*)(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Seq-Const-Succ)}$ $\frac{a \triangleright (aa^*)(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon}{a \triangleright (aa^*)(y \text{ as } a^*) \mid \varepsilon(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Choice-Fst)}$ $\frac{a \triangleright (aa^* \mid \varepsilon)(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon}{a \triangleright a^*(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Seq-Rep)}$ $\frac{a \triangleright a^*(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon}{a \triangleright (x \text{ as } a^*)a^* \Rightarrow x \mapsto a} \text{ (M-Seq-Bind)}$	$\frac{}{\varepsilon \triangleright aa^* \Rightarrow \perp} \text{ (M-Seq-Const-Fail)}$ $\frac{}{\varepsilon \triangleright \varepsilon \Rightarrow \emptyset} \text{ (M-Const-Succ)}$ $\frac{\varepsilon \triangleright aa^* \mid \varepsilon \Rightarrow \emptyset}{\varepsilon \triangleright a^* \Rightarrow \emptyset} \text{ (M-Rep)}$ $\frac{\varepsilon \triangleright y \text{ as } a^* \Rightarrow y \mapsto \varepsilon}{\varepsilon \triangleright \varepsilon(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Bind)}$ $\frac{}{\varepsilon \triangleright \varepsilon(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Seq-Const-Succ)}$ $\frac{}{\varepsilon \triangleright \varepsilon(y \text{ as } a^*) \Rightarrow y \mapsto \varepsilon} \text{ (M-Choice-Snd)}$
---	--

図 4 パターンマッチの例

Fig. 4 Example of pattern matching.

$s_1$  が  $P_1$  にマッチし,  $s_2$  が  $P_2$  にマッチする, という分割をしてしまうとパターンマッチに非決定性の挙動が含まれてしまうためである. 6 つの規則はそれぞれ  $P_1P_2$  を, 導出の前件部で  $P_1$  の形がより“小さく”なるように  $P_1$  の意味に従って書き換える(ただし繰返しの場合を除く). (M-Seq-Const-...) と (M-Seq-Bind) 規則に見られる  $s_1^{-1}s_2$ ,  $s_1s_2^{-1}$  という表記は,  $s_1s = s_2$  あるいは  $s_1 = ss_2$  を満たすような文字列  $s$  を(もしあれば)表す記法とする.

例 1 文字列  $a$  をパターン  $(x \text{ as } a^*)a^*$  にマッチさせると置換  $x \mapsto a$  を生じる. この導出の例を図 4 に示す. first-match, longest-match のポリシーがどのように働いているかに注意されたい.

ところで, このパターンマッチの導出が“無限ループ”に陥って停止しなくなってしまうケースがある. 空文字列の繰返しがあるパターン, すなわちパターンが  $P^*$  で  $P$  が  $\varepsilon$  にマッチするような場合がそれに相当する. もう少し厳密ないい方をすれば, 上述の(帰納的に定義された)規則で, たとえば  $\varepsilon \triangleright \varepsilon^* \Rightarrow \emptyset$  のような導出をすることはできない. 幸いにも, このような繰返しは前処理によってパターンの意味を変えることなく排除することができる. この技法についてはすでに研究がなされている. たとえば文献 3) の 4 章などを参照されたい. よって本論文では以降, このようなパターンは出現しないことを仮定する.

注意深い読者は,  $s \triangleright P \Rightarrow \perp$  を単に “ $s \triangleright P \Rightarrow \theta$  を満たす  $\theta$  が存在しない” ことで定義しないことに疑問を持たれるかもしれない. しかしながらこの定義では停止しないこととパターンマッチに失敗することが区別できないという問題がある. のみならず, このや

り方ではパターンマッチの関係の定義そのものが実は成立しない. というのも (M-Choice-Snd) 規則が前提の  $s \triangleright P_1 \Rightarrow \perp$  のために帰納的な定義でなくなってしまうからである.

#### 4. 型システム

$\lambda^{re}$  の型  $\tau$  は図 5 のとおりに定義される. 関数型  $\tau_1 \xrightarrow{T} \tau_2$  は  $\tau_1$  型の引数を取り,  $T$  型の文字列を出力し, 返り値があればその型は  $\tau_2$  となるような関数を表す. 文字列型  $T$  は  $\llbracket T \rrbracket$  と表記される文字列の集合を表す. これは文字列定数  $s$ , 選択  $T_1 \mid T_2$ , 繰返し  $T^*$ , 接続  $T_1T_2$  のような正規表現であるか, 交差  $T_1 \cap T_2$ , 差分  $T_1 \cap \overline{T_2}$ , left quotient  $T_1^{-1}T_2$ , right quotient  $T_1T_2^{-1}$  のような正規表現に対する演算の形をとる. 正則言語の集合はこれらの演算に関して閉じていることが知られている<sup>4)</sup>. しかしながら, このような演算は型の構文として定義する方が, 型付け規則を与えるうえで有用である. ところで, 型の構文は変数の束縛を除いてパターンの構文を含んでいる. そこで以降, パターン  $P$  からすべての変数を除去して得られる型を  $novar(P)$  と書くことにする. すなわち,  $novar(s) = s$ ,  $novar(P_1 \mid P_2) = novar(P_1) \mid novar(P_2)$ ,  $novar(x \text{ as } P) = novar(P)$ ,  $novar(P^*) = novar(P)^*$ ,  $novar(P_1P_2) = novar(P_1)novar(P_2)$  と定める. さらに, 型の交差  $\cap$  と選択  $\mid$  の定義を文字列型  $T$  から一般の型  $\tau$  に拡張する. 関数型  $\tau_1 \xrightarrow{T} \tau_2$  についてはそれぞれの演算を以下のように定義する.

実のところ XDuce の論文 5) のパターンマッチの定義はこの問題をかかえている.

$\tau$ (type)	$::= \tau_1 \xrightarrow{T} \tau_2$		
	$T$		
$T$ (string type)	$::= s$	$[[s]] = \{s\}$	
	$T_1   T_2$	$[[T_1   T_2]] = [[T_1]] \cup [[T_2]]$	
	$T^*$	$[[T^*]] = \{s_1 + \dots + s_n \mid s_i \in [[T]] (i = 1, \dots, n)\}$	
	$T_1 T_2$	$[[T_1 T_2]] = \{s_1 + s_2 \mid s_1 \in [[T_1]] \wedge s_2 \in [[T_2]]\}$	
	$T_1 \cap T_2$	$[[T_1 \cap T_2]] = [[T_1]] \cap [[T_2]]$	
	$T_1 \cap \overline{T_2}$	$[[T_1 \cap \overline{T_2}]] = [[T_1]] \setminus [[T_2]]$	
	$T_1^{-1} T_2$	$[[T_1^{-1} T_2]] = \{s' \mid s \in [[T_1]] \wedge s + s' \in [[T_2]]\}$	
	$T_1 T_2^{-1}$	$[[T_1 T_2^{-1}]] = \{s \mid s + s' \in [[T_1]] \wedge s' \in [[T_2]]\}$	

図 5 型の構文と意味

Fig. 5 Syntax and semantics of types.

$$(\tau_1 \xrightarrow{T_1} \tau'_1) \mid (\tau_2 \xrightarrow{T_2} \tau'_2) = (\tau_1 \cap \tau_2) \xrightarrow{T_1 \mid T_2} (\tau'_1 \mid \tau'_2)$$

$$(\tau_1 \xrightarrow{T_1} \tau'_1) \cap (\tau_2 \xrightarrow{T_2} \tau'_2) = (\tau_1 \mid \tau_2) \xrightarrow{T_1 \cap T_2} (\tau'_1 \cap \tau'_2)$$

このような定義の代わりに、 $\tau_1 \cap \tau_2$  や  $\tau_1 \mid \tau_2$  も型の構文として定義し、本物の intersection 型 (つまりオーバロード)、union 型を考えることもできる。同様のアイデアが文献 2) で研究されており、この成果を本研究に取り入れることも可能ではあるが、ここでは簡潔さのためその方法はとらない。

型付け規則は図 6 で与えられる。規則 (S-Str) にあるように、文字列型  $T_1$  と  $T_2$  の間の部分型関係  $T_1 \leq T_2$  は、型の表示の包含関係として定義する。型の表示は正則集合なのでその包含関係は決定可能である (ただし、効率的な実装に関しては本論文の範囲を超える)。文字列型の部分型関係は一般の型の間に  $\tau_1 \leq \tau_2$  として拡張される。関数型に関しては (S-Fun) 規則にあるように標準的な部分型の定義を用いる。

型判定  $\Gamma \vdash M : \tau, T$  は “型環境  $\Gamma$  の下で項  $M$  は型  $\tau$  と副作用  $T$  を持つ” と読むことができる。型付け規則は副作用とパターンマッチの部分を除いて標準的なものである。副作用に関しては、たとえば (T-Print) 規則を考えると  $\text{print } M$  という項は最初に  $M$  を評価して結果の文字列  $s$  を出力し、空文字列  $\varepsilon$  に評価されるので、 $M$  が文字列型  $T$  と副作用  $T'$  を持つならば  $\text{print } M$  は文字列型  $\varepsilon$  と副作用  $T' T$  で型付けされる。副作用の型付けは他の規則でも同様の推論に基づいている。

パターンマッチに関しては、(T-Match) 規則を参照されたい。まず、入力  $M$  の型は文字列型  $T_1$  でなければならない。 $T_1$  型の入力文字列に対し、最初のパターン  $P_1$  とのマッチを試みる。もしこのマッチが成功すれば、何らかの変数束縛が生成される。このとき

生じた束縛変数の型を計算するため、補助的な関係として  $T \rightsquigarrow P \Rightarrow \Gamma$  を導入する。この関係は “もし型  $T$  を持つ文字列がパターン  $P$  にマッチしたならば、各束縛変数  $x$  の型は  $\Gamma(x)$  に従う” と読むことができる。このようにして、型環境  $\Gamma_1$  が  $T_1$  と  $P_1$  から  $T_1 \rightsquigarrow P_1 \Rightarrow \Gamma_1$  によって計算され、 $P_1$  に対応する項  $M_1$  が新しい型環境の下で  $\Gamma, \Gamma_1 \vdash M_1 : \tau_1, T'_1$  のように型付けされる。ところで、もし  $T_1$  型のどんな文字列も  $P_1$  にマッチしない、すなわち、 $T_1$  型の文字列で  $P_1$  にマッチするものの集合  $T_1 \cap \text{novar}(P_1)$  が空ならば、このパターンマッチが成功することはありえないので  $P_1$  は冗長である。このことが型システムの健全性に直接影響を与えることはないが、我々はプログラマの利便のため、この種の冗長性に対して警告を与えることにする。ただし、冗長性をただちにエラーと見なすことはできない。このようにすると後の subject reduction の性質が成り立たなくなってしまうからである。さらに、我々は first-match のポリシーを採用したので、 $P_1$  にマッチする入力文字列は次のパターン  $P_2$  の入力とはならない。そこで、 $P_2$  に対する入力文字列の型  $T_2$  は、 $T_1$  と  $P_1$  (にマッチする文字列の型) の差分  $T_1 \cap \overline{\text{novar}(P_1)}$  とする。このようにして各パターン  $P_m$  と対応する項  $M_m$  ( $m = 1, \dots, n$ ) に対する型付けが行われる。最後に、パターンマッチが exhaustive、つまり入力がいずれかのパターンに必ずマッチすることを保証するため、入力文字列のうちのどのパターンにもマッチしないものの型  $T_{n+1}$  は空でなければならない。ここで空の型 (その表示が空集合となる型) は型  $\varepsilon$  (表示が 1 要素  $\varepsilon$  のみからなる集合の型) と区別するように注意されたい。以後、空の型のことを  $\emptyset$  と表記することとする。空の型としてどのような型を具体的に取るかは問題ではないが、たと

$$\begin{array}{c}
\frac{[[T_1]] \subseteq [[T_2]]}{T_1 \leq T_2} \text{(S-Str)} \quad \frac{\tau'_1 \leq \tau_1 \quad \tau_2 \leq \tau'_2 \quad T_1 \leq T_2}{\tau_1 \xrightarrow{T_1} \tau_2 \leq \tau'_1 \xrightarrow{T_2} \tau'_2} \text{(S-Fun)} \\
\\
\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau, \varepsilon} \text{(T-Var)} \quad \frac{\Gamma, f : \tau_1 \xrightarrow{T} \tau_2, x : \tau_1 \vdash M : \tau_2, T' \quad T' \leq T}{\Gamma \vdash \text{fix}(f, x, M) : \tau_1 \xrightarrow{T} \tau_2, \varepsilon} \text{(T-Fix)} \\
\\
\frac{\Gamma \vdash M_1 : \tau_2 \xrightarrow{T} \tau_1, T_1 \quad \Gamma \vdash M_2 : \tau'_2, T_2 \quad \tau'_2 \leq \tau_2}{\Gamma \vdash M_1 M_2 : \tau_1, T_1 T_2 T} \text{(T-App)} \quad \frac{}{\Gamma \vdash s : s, \varepsilon} \text{(T-Const)} \\
\\
\frac{\Gamma \vdash M_1 : T_1, T'_1 \quad \Gamma \vdash M_2 : T_2, T'_2}{\Gamma \vdash M_1 \wedge M_2 : T_1 T_2, T'_1 T'_2} \text{(T-Cat)} \quad \frac{\Gamma \vdash M : T, T'}{\Gamma \vdash \text{print } M : \varepsilon, T' T} \text{(T-Print)} \\
\\
\frac{\Gamma \vdash M : T_1, T' \quad T_m \rightsquigarrow P_m \Rightarrow \Gamma_m \quad \text{warn of redundancy if } T_m \cap \text{novar}(P_m) \leq \emptyset \quad \Gamma, \Gamma_m \vdash M_m : \tau_m, T'_m \quad T_{m+1} = T_m \cap \overline{\text{novar}(P_m)} \quad (m = 1, \dots, n) \quad T_{n+1} \leq \emptyset}{\Gamma \vdash \text{match } M \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n : \tau_1 \mid \dots \mid \tau_m, T'(T'_1 \mid \dots \mid T'_m)} \text{(T-Match)} \\
\\
\frac{([[T]], P) \in \Pi \quad \forall x \in \text{dom}(\Gamma) = \text{var}(P), \Gamma(x) = \emptyset}{\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma} \text{(P-Mem)} \quad \frac{}{\Pi \vdash T \rightsquigarrow s \Rightarrow \emptyset} \text{(P-Const)} \\
\\
\frac{\Pi \vdash T \rightsquigarrow P_1 \Rightarrow \Gamma_1 \quad \Pi \vdash T \cap \overline{\text{novar}(P_1)} \rightsquigarrow P_2 \Rightarrow \Gamma_2 \quad \forall x \in \text{dom}(\Gamma) = \text{dom}(\Gamma_1) = \text{dom}(\Gamma_2), \Gamma(x) = \Gamma_1(x) \mid \Gamma_2(x)}{\Pi \vdash T \rightsquigarrow P_1 \mid P_2 \Rightarrow \Gamma} \text{(P-Choice)} \\
\\
\frac{\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow x \text{ as } P \Rightarrow \Gamma \uplus \{x : \Gamma(P) \cap T\}} \text{(P-Bind)} \quad \frac{\Pi \uplus \{([T], P^*)\} \vdash T \rightsquigarrow P P^* \mid \varepsilon \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow P^* \Rightarrow \Gamma} \text{(P-Rep)} \\
\\
\frac{\Pi \vdash s^{-1} T \rightsquigarrow P \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow s P \Rightarrow \Gamma} \text{(P-Seq-Const)} \quad \frac{\Pi \vdash T \rightsquigarrow P_1 P_3 \mid P_2 P_3 \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow (P_1 \mid P_2) P_3 \Rightarrow \Gamma} \text{(P-Seq-Choice)} \\
\\
\frac{\forall ([T], P) \in \Pi, y \notin \text{var}(P) \quad y \notin \text{var}(P_1 P_2) \quad \Pi \vdash T_1 \rightsquigarrow P_1(y \text{ as } P_2) \Rightarrow \Gamma \uplus \{y : T_2\}}{\Pi \vdash T_1 \rightsquigarrow (x \text{ as } P_1) P_2 \Rightarrow \Gamma \uplus \{x : \Gamma(P_1) \cap T_1 T_2^{-1}\}} \text{(P-Seq-Bind)} \\
\\
\frac{\Pi \uplus \{([T], P_1^* P_2)\} \vdash T \rightsquigarrow (P_1 P_1^* \mid \varepsilon) P_2 \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow P_1^* P_2 \Rightarrow \Gamma} \text{(P-Seq-Rep)} \quad \frac{\Pi \vdash T \rightsquigarrow P_1 (P_2 P_3) \Rightarrow \Gamma}{\Pi \vdash T \rightsquigarrow (P_1 P_2) P_3 \Rightarrow \Gamma} \text{(P-Seq-Seq)}
\end{array}$$

図 6 型付け規則

Fig. 6 Typing rules.

例えば  $\varepsilon \cap \bar{\varepsilon}$  などがあるだろう。

パターンマッチの型付け  $T \rightsquigarrow P \Rightarrow \Gamma$  は操作的意味論  $s \triangleright P \Rightarrow \theta$  と似た形で定義される。ここでも各規則はボトムアップに、パターン  $P$  と型  $T$  をとって結果の型環境  $\Gamma$  を計算するアルゴリズムと見なせる。文字列型は文字列の ( 正則 ) 集合を表すので、パター

ンの型付け規則は実際、パターンの操作的意味論に非常に近い ( 図 3 参照 )。しかしながら、個別の文字列の代わりに文字列の型をとる、という明らかな違いに加えて、大きく 2 つの違いがある。

1 つは型推論の停止性を保証するために “memoization” の手法を利用している点である。具体的には、

導出の関係を  $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma$  の形に一般化する．ここで  $\Pi$  は文字列型  $T$  (の表示) とパターン  $P$  のペア  $([T], P)$  の集合である． $\Pi = \emptyset$  のとき単に  $T \rightsquigarrow P \Rightarrow \Gamma$  と書くことにする．直観的には、 $\Pi$  は型付けの過程で“すでに一度出現した”繰返しのパターンと型の組を記憶している．規則 (P-Req), (P-Seq-Rep) などを参照されたい．アルゴリズムの途中で再びこのような組に遭遇すると、規則 (P-Mem) にあるようにすべての変数に空の型を割り当てた型環境を返す．この規則がないと、繰返しを含むパターンに関する型付けはつねに停止しなくなってしまう．

もう 1 つの違いは、この型付けそのものはパターンマッチが成功することを保証しない点にある．型付けは単にマッチが成功したならばどのような文字列が変数に束縛されるかを推論するにすぎない．これが最も明確に現れるのは (P-Const) 規則で、たとえば  $\emptyset \vdash a^* \rightsquigarrow b \Rightarrow \emptyset$  のような導出が可能である．このことの意味は“ $a^*$  型の文字列を定数文字列のパターン  $b$  にマッチさせると変数に対する束縛は起こらない”ということであるが、当然このパターンマッチが成功することはありえない．このような定義の仕方は明快さを損ねるようにも思えるが、実は必然的なものである．特に  $P_1 \mid P_2$  のような選択のパターンで、入力がつねに  $P_1$  にマッチすることを (P-Choice) で要求してしまうと、制約が強くなりすぎてしまう．一例として、パターン  $*$  を  $\text{match } M \text{ with } P \Rightarrow M_1 \mid x \text{ as } * \Rightarrow M_2$  のように“デフォルト”として使うことを考えよう(ただし、はすべての文字の選択  $a \mid b \mid \dots$  の省略記法とする)． $*$  への入力文字列が任意の形をとりうるわけではないが、このパターンをデフォルトとして用いるのはきわめて自然な発想であろう．とはいえ、今の例では  $P$  にマッチする文字列は  $x$  に束縛されないことは確かなのだから、 $x$  の型にはその事実を何らかの形で反映させたい．それを実現するのが (P-Bind) 規則で、ここで入力の型  $T$  と型環境  $\Gamma$  の元でのパターン  $P$  の値域  $\Gamma(P)$  の交差をとって変数の型としている．ただし、ここで  $\Gamma(P)$  は  $\Gamma(s) = s$ ,  $\Gamma(P_1 \mid P_2) = \Gamma(P_1) \mid \Gamma(P_2)$ ,  $\Gamma(x \text{ as } P) = \Gamma(x) \cap \Gamma(P)$ ,  $\Gamma(P^*) = \Gamma(P)^*$ ,  $\Gamma(P_1 P_2) = \Gamma(P_1)\Gamma(P_2)$  で定義される型である．同じ発想で (P-Seq-Bind) 規則も定義されている．なお、型環境に新しい要素を追加する際には  $\boxplus$  を使って保守的な形でのみ拡張を許している．すなわち、追加し

ようとしている要素と同じ変数に関する要素がすでに型環境に含まれていた場合、拡張の結果は未定義としているので、パターン中の変数については線形性が要求される．

以降、簡単のため  $[T] = [T']$  のときに型  $T$  を別の型  $T'$  に暗黙に置き換えて使用することがある．これによって型付けの結果が影響を受けることはない．

例 2 型  $aa^*$  の文字列をパターン  $(x \text{ as } a^*)(y \text{ as } a^*)$  にマッチさせると、 $x$  の型は  $aa^*$ ,  $y$  の型は  $\varepsilon$  となる．すなわち、 $aa^* \rightsquigarrow (x \text{ as } a^*)(y \text{ as } a^*) \Rightarrow x : aa^*, y : \varepsilon$  という導出ができる．この導出の例を図 7 に示す．この型付けが first-match, longest-match のポリシーを反映していることに注意されたい．

上述の型システムの健全性は以下のように証明される．議論の多くは一般的なものであるが、パターンマッチの性質についていくつか言及している．証明の詳細については紙幅の都合上割愛するが、1 章で紹介した web サイトに掲載しているので、そちらを参照されたい．

定理 3 (型と副作用の健全性)  $\vdash M : \tau, T$  かつ  $M \xrightarrow{s_1} M_1 \xrightarrow{s_2} M_2 \xrightarrow{s_3} \dots \xrightarrow{s_n} M_n$  ならば、 $s_1 s_2 s_3 \dots s_n$  はある文字列  $s \in [T]$  の prefix である．さらに、 $M_n$  がこれ以上簡約できない ( $M_n \xrightarrow{s} M'$  なる  $s, M'$  が存在しない) ならば、 $M_n$  は値でその型は  $\tau$  の部分型であり (すなわち、ある  $v, \tau'$  について、 $M_n = v$ ,  $\vdash v : \tau', \varepsilon, \tau' \leq \tau$ )、同時に  $s_1 s_2 s_3 \dots s_n \in [T]$  が成り立つ．

証明．下の progress と subject reduction の補題よりただちに導かれる．  $\square$

補題 4 (Progress)  $\vdash M : \tau, T$  で  $M$  が値でないならば、 $M$  は簡約できる (すなわち、 $M \xrightarrow{s} M'$  なる  $s$  と  $M'$  が存在する)．

証明． $\vdash M : \tau, T$  の導出に関する帰納法．ただし (T-Match) の場合に以下の 2 つの補題を使う．  $\square$

補題 5 (パターンマッチの決定性・停止性)  $s \triangleright P \Rightarrow \theta_1$  かつ  $s \triangleright P \Rightarrow \theta_2$  ならば  $\theta_1 = \theta_2$ ．さらに、 $s \in [\text{novar}(P)]$  ならば  $s \triangleright P$ ．逆に  $s \notin [\text{novar}(P)]$  ならば  $s \triangleright P \Rightarrow \perp$ ．

証明概略．証明は見かけよりもやや複雑なものとなる．単純にパターンの構造に関して帰納的に証明する方法では、規則 (M-Rep), (M-Seq-Rep) のところでパターンの構造が結論部分よりも仮定部分において大きくなってしまいうまくいかない．我々の証明では、まず最初にパターンマッチのもう 1 つのアルゴリズムを定義し、新しいアルゴリズムの決定性・停止性と元々のアルゴリズムとの互換性を示す、という段階

これとは逆に、XDuce<sup>5)</sup> では末尾以外の場所に現れる  $\text{as}$  パターンについて正確な型を与えることができない．Regular tree と正規表現は理論的には近い位置にあるため、我々のアプローチは XDuce にも適用できるかもしれないと考えている．



$$\begin{array}{c}
\vdots \\
\frac{}{\Pi_1 \vdash \emptyset \rightsquigarrow a^* \Rightarrow \emptyset} \text{(P-Rep)} \\
\frac{}{\Pi_1 \vdash \emptyset \rightsquigarrow y \text{ as } a^* \Rightarrow y : \emptyset} \text{(P-Bind)} \\
\frac{}{\Pi_1 \vdash \emptyset \rightsquigarrow z \text{ as } (y \text{ as } a^*) \Rightarrow y : \emptyset, z : \emptyset} \text{(P-Bind)} \\
\frac{}{\Pi_1 \vdash \emptyset \rightsquigarrow \varepsilon(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \emptyset, z : \emptyset} \text{(P-Seq-Const)} \\
\frac{\Delta}{\Pi_1 \vdash aa^* \rightsquigarrow (aa^*)(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Seq)} \\
\frac{}{\Pi_1 \vdash aa^* \rightsquigarrow (aa^*)(z \text{ as } (y \text{ as } a^*)) \mid \varepsilon(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Choice)} \\
\frac{}{\Pi_1 \vdash aa^* \rightsquigarrow (aa^* \mid \varepsilon)(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Choice)} \\
\frac{}{\emptyset \vdash aa^* \rightsquigarrow a^*(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Rep)} \\
\frac{}{\emptyset \vdash aa^* \rightsquigarrow (x \text{ as } a^*)(y \text{ as } a^*) \Rightarrow x : aa^*, y : \varepsilon} \text{(P-Seq-Bind)} \\
\Delta = \\
\\
\vdots \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow a^* \Rightarrow \emptyset} \text{(P-Rep)} \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow y \text{ as } a^* \Rightarrow y : \varepsilon} \text{(P-Bind)} \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow z \text{ as } (y \text{ as } a^*) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Bind)} \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow \varepsilon(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Const)} \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow (aa^*)(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \emptyset, z : \emptyset} \text{(P-Seq-Seq)} \\
\frac{}{\Pi_2 \vdash \varepsilon \rightsquigarrow a(a^*(z \text{ as } (y \text{ as } a^*))) \Rightarrow y : \emptyset, z : \emptyset} \text{(P-Seq-Const)} \\
\frac{}{\Pi_2 \vdash a^* \rightsquigarrow a^*(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \emptyset, z : \emptyset} \text{(P-Mem)} \\
\frac{}{\Pi_2 \vdash a^* \rightsquigarrow (aa^*)(z \text{ as } (y \text{ as } a^*)) \mid \varepsilon(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Choice)} \\
\frac{}{\Pi_2 \vdash a^* \rightsquigarrow (aa^* \mid \varepsilon)(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Rep)} \\
\frac{}{\Pi_1 \vdash a^* \rightsquigarrow a^*(z \text{ as } (y \text{ as } a^*)) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Const)} \\
\frac{}{\Pi_1 \vdash aa^* \rightsquigarrow a(a^*(z \text{ as } (y \text{ as } a^*))) \Rightarrow y : \varepsilon, z : \varepsilon} \text{(P-Seq-Const)}
\end{array}$$

ただし  $\Pi_1 = \{([aa^*], a^*(z \text{ as } (y \text{ as } a^*)))\}$ ,  $\Pi_2 = \Pi_1 \cup \{([a^*], a^*(z \text{ as } (y \text{ as } a^*)))\}$ .

図 7 導出の例

Fig. 7 Example of derivation.

をふむ。新しいアルゴリズムでは first/longest-match を含むすべての可能なマッチの集合が計算される。このアルゴリズムは文献 14) のものに基づいているが、停止性を保証するためにパターンが 3 章で述べたように前処理<sup>3)</sup> されている必要がある、という点で文献のものとは多少異なっている。

#### 補題 6 (パターンマッチの Exhaustiveness)

$\Gamma \vdash \text{match } s \text{ with } P_1 \Rightarrow M_1 \mid \dots \mid P_n \Rightarrow M_n : \tau, T$  ならば、ある  $1 \leq i \leq n$  について  $s \triangleright P_i$ 。

補題 7 (Subject Reduction)  $\Gamma \vdash M : \tau, T$  かつ  $M \xrightarrow{s} M'$  ならば、ある  $\tau', T'$  について  $\Gamma \vdash M' : \tau', T'$  であり、 $\tau', T'$  は  $\tau' \leq \tau$  と  $sT' \leq T$  を満たす。証明。  $\Gamma \vdash M : \tau, T$  の導出に関する帰納法。(T-App) と (T-Match) の場合に以下の 2 つの補題をそれぞれ使う。□

補題 8 (Substitution)  $\Gamma, x : \tau' \vdash M : \tau, T$  かつ  $\Gamma \vdash v : \tau', \varepsilon$  ならば  $\Gamma \vdash [v/x]M : \tau, T$ 。

補題 9 (パターンマッチの型推論の健全性・完全性)  $T \rightsquigarrow P \Rightarrow \Gamma$  ならば、各々の  $x \in \text{var}(P) = \text{dom}(\Gamma)$  について  $[[\Gamma(x)]] = \{\theta(x) \mid s \in [T] \wedge s \triangleright P \Rightarrow \theta \wedge \theta \neq \perp\}$ 。

最後の補題は見かけよりもかなり巧妙な証明を要する(特に  $\supseteq$  の部分について)。たとえば、もしパターンに  $\varepsilon^*$  のような空(かもしれない)文字列の繰返しを認めると、この性質は成り立たなくなる。反例とし

ては  $a \rightsquigarrow \varepsilon^*(y \text{ as } a) \Rightarrow y : \emptyset$  などをあげることができる。このようなケースでは型システムの健全性が破壊される結果を生んでしまう(実際には 3 章で述べたように、この種のパターンは一般性を失わずあらかじめ排除できるので、起こりえないケースである)。以上のような事情があるため、証明は複雑なものとなる。少し具体的には、 $\Pi$  に関していくつかの補助的な関係を定義したうえで、成り立つ不変条件について言及するという方法をとる。

証明概略。まず、なぜ単純な方法での証明では巧みかないのかを考察する。今、我々が示したいことは“ $\emptyset \vdash T \rightsquigarrow P \Rightarrow \Gamma$  ならば  $[[\Gamma(x)]]$  は実際に  $x$  に束縛されうる文字列の集合  $\{\theta(x) \mid s \in [T] \wedge s \triangleright P \Rightarrow \theta \wedge \theta \neq \perp\}$  に一致する”ということである。この証明に、単純に  $\emptyset \vdash T \rightsquigarrow P \Rightarrow \Gamma$  の導出に関する帰納法などを適用することはできない。なぜならば、規則 (P-Rep) などでは  $\Pi$  は  $\emptyset$  ではなくってしまい、帰納法の仮定が適用できなくなるからである。次に考えられるのは補題全体の言明を“ $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma$  ならば...”のように強めるやり方であるが、この方法もうまくいかない。なぜならば、 $\Pi$  の要素として任意のものを許してしまうと、 $\Pi$  の中に導出に悪影響を与える「ごみ」が含まれうるためである。そのような「ごみ」によって正しくない結果を生じてしまう導出の例を図 8 に示す。ここで変数  $x$  の型は  $\emptyset$  であると推論

$$\frac{\frac{([\![ab]\!], a(y \text{ as } b)) \in \{([\![ab]\!], a(y \text{ as } b))\}}{\{([\![ab]\!], a(y \text{ as } b))\} \vdash ab \rightsquigarrow a(y \text{ as } b) \Rightarrow \{y : \emptyset\}} \text{(P-Mem)}}{\{([\![ab]\!], a(y \text{ as } b))\} \vdash ab \rightsquigarrow (x \text{ as } a)b \Rightarrow \{x : \emptyset\}} \text{(P-Seq-Bind)}$$

図 8 正しくない導出の例

Fig. 8 Example of wrong derivation.

されているが、一見して分かるように  $x$  の正しい型は  $a$  になるはずである。以上の事柄をふまえ、証明に先立って導出の各過程で  $\Pi$  の中にこのような「ごみ」が含まれないことを保証しなければならない。そのために新しい導出規則と  $\Pi$  に関する述語を導入する。定義の詳細は前出の web サイトに譲るが、以下にそれぞれを簡単に解説する。

#### 導出の拡張

まず、導出  $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma$  を  $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma; \Pi'$  の形に拡張する。ここで  $\Pi'$  は  $\Pi$  の要素のうち、(P-Mem) 規則によって実際に使われた要素を記録する集合である。

#### $\Pi$ の正当性

次に、導出  $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma; \Pi'$  において  $\Pi$  の中の要素が正当なものである、ということ定義するための関係を導入する。 $\Pi \stackrel{P}{\bowtie} T$  と書かれる関係によって、 $\Pi$  の中のどの要素も、少なくとも  $T, P$  に関して悪影響を持たないということを表す。この関係が  $\Pi$  の中に「ごみ」が含まれていない、ということの定義となる。

#### 証明の流れ

以上の定義を導入したうえで、証明は以下のようになされる。(1) 最初に、導出  $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma; \Pi'$  において  $\Pi \stackrel{P}{\bowtie} T$  が成り立っているならば、導出の仮定の段階でも同様の性質が成り立つことを示す。すなわち、導出の帰結において  $\Pi$  が「ごみ」を含まないならば、この性質は導出全体の不変条件となることを証明する。(2) 導出が (1) の性質を満たすならば、実際に使われた要素の集合  $\Pi'$  も特定の性質を満たすことを示す。(3) 最後に、導出において  $\Pi'$  が (2) の性質を満たすならば型推論の健全性・完全性が成り立つことを示して証明の完了となる。結局、最終的に証明されるのは以下の言明となる：“ $\Pi \vdash T \rightsquigarrow P \Rightarrow \Gamma; \Pi'$  かつ  $\Pi \stackrel{P}{\bowtie} T$  ならば健全性・完全性が ( $\Pi'$  に関して少し強められた形で) 満たされる。”  $\Pi = \emptyset$  のとき、(直観的には  $\emptyset$  に「ごみ」は含まれないので) 任意の  $T, P$  について  $\emptyset \stackrel{P}{\bowtie} T$  が成り立つ。このことから、補題 9 は上の言明の系として示される。

## 5. プログラム例

本章では  $\lambda^{re}$  によるプログラムの小さな例(とその型付け)を紹介する。

### ブール値

文字列  $t$  と  $f$  を true, false と見なせば、ブール代数は  $\lambda^{re}$  で表現できる。条件式  $\text{if } M_0 \text{ then } M_1 \text{ else } M_2$  はパターンマッチを使って  $\text{match } M_0 \text{ with } t \Rightarrow M_1 \mid f \Rightarrow M_2$  と表すことができる。型  $\text{bool}$  はそれゆえ  $t \mid f$  として定義できる。これと再帰関数を組み合わせれば、 $\lambda^{re}$  は Turing 完全である。

### リスト

型  $T$  を、文字、(カンマ) を中に含まない任意の文字列を表す型としよう。このとき、 $T$  型の文字列のリストは  $(T,)^*$  型の文字列として表現できる。その考え方は以下のとおりである：空リストは空文字列  $\varepsilon$  で表し、要素  $s$  とリスト  $\ell$  の cons は  $s^\wedge, \wedge \ell$  で表す。リスト  $\ell$  を分解するコードは、パターンマッチを使って  $\text{match } \ell \text{ with } \varepsilon \Rightarrow M_1 \mid (x \text{ as } T), (y \text{ as } *) \Rightarrow M_2$  のようにして実現できる。 $x$  は  $\ell$  の先頭要素(“car”)に、 $y$  は残りの要素(“cdr”)に対応する。 $y$  のパターンはプログラム上は  $*$  と書かれているが、 $\ell$  の型が  $(T,)^*$  であれば、 $y$  の型は  $(T,)^*$  と正確に、自動的に推論される。この例が末尾以外の場所の  $\text{as}$  パターンの例にもなっていることに注意されたい。これは先行研究<sup>5)</sup>と比較した我々の型システムの利点である。

### フィルタリング

(上の方法で表現された)リストに対しては、“grep”のような簡単なフィルタリングが以下のように明らかな方法で実現できる。ただしフィルタリングに使用する正規表現  $T'$  はコンパイル時に静的に与えられるものとする。

$$\begin{aligned} & \text{fix}(f, x, \text{match } x \text{ with} \\ & \quad \varepsilon \Rightarrow \varepsilon \\ & \quad \mid (y \text{ as } T'), (z \text{ as } *) \Rightarrow y^\wedge, \wedge f(z) \\ & \quad \mid T, (z \text{ as } *) \Rightarrow f(z) \end{aligned}$$

このフィルタ関数の型は期待どおり  $(T,)^* \xrightarrow{\varepsilon} ((T \cap T'),)^*$  となる。

## パーズと整形表示

自然数とそのうでの基本的な演算が(簡単のため厳密な定義抜きに)組み込み済みであるとすると,自然数とその文字列による表記の間の相互変換は下のプログラムで実現できる.

```
let dig2chr = λd.
  if d = 9 then 9 else
  if d = 8 then 8 else
  ...
  if d = 1 then 1 else
  0 in
fix(nat2str, n,
  if n < 10 then dig2chr(n) else
  nat2str(n div 10) ^ dig2chr(n mod 10))
```

および,

```
let chr2dig = λc.
  match c with 9 ⇒ 9
  | ...
  | 0 ⇒ 0 in
fix(str2int, s,
  match s with (d as .) ⇒ chr2dig(d)
  | (d as .)(t as .*) ⇒
  chr2dig(d) × 10 + str2int(t)
```

それぞれの型は  $\text{nat} \xrightarrow{\varepsilon} u^*$  と  $u^* \xrightarrow{\varepsilon} \text{nat}$  である. ただし  $\iota = 0 \mid \dots \mid 9$ . しかしながら, 我々の型システムは  $\text{nat2str}$  の値域の型を  $0 \mid (\iota \cap \bar{0})\iota^*$  と型付けできるほど強力ではない. 実際には  $00$  のような結果は起こりえないのであるが, それを型システムで表現するには  $\text{nat}$  に対する dependent type を必要とするだろう.

## 6. 結 論

我々はテキスト処理のための小さな言語として  $\lambda^{re}$  を定義し, 文字列と副作用に対して正規表現を用いた型システムを導入した.

XDuce のほかには Igarashi らによる *resource usage analysis*<sup>7)</sup> (および先行研究) がおそらく我々と最も関連の深い研究である. Resource usage analysis は, ファイル・メモリ領域など種々の計算機資源に対してプログラムがどのようなパターンでアクセスするかを静的に解析(そして検証)するための一般的な枠組みである. 我々の正規表現による effect system は彼らの枠組みの簡単な例と見なすことができる. しかしながら彼らの枠組みでは, アクセスパターンの記述

言語は再帰を含んでしまうため等価性と包含の判定は決定不能であり, 型チェックも一般には決定不能となってしまう. 一方, 我々の体系では文字列の型は正則であるため, fix の束縛変数に型の注釈が付いていれば型チェックは決定可能である.

本論文では正規表現による型システムの定義とその健全性という, 理論的な基礎となる部分について紹介したが, 実用化のためには他にも解決すべき課題が残されている. とりわけ, 型チェック, そしてより重要な型推論の効率的なアルゴリズムは実用にあたって本質的な位置を占めるが, 現在ははまだ用意されていない. 技術的に重大な問題として, 制約解消に基づく標準的な型推論の手法では, 制約の中に再帰を含んでしまうため正規表現よりも表現力の強い(少なくとも文脈自由文法程度の)解が必要になるという点があげられる. 1章で見た  $f(n) = \text{if } n \leq 0 \text{ then } \varepsilon \text{ else } a^{\varepsilon} f(n-1)^{\varepsilon} b$  という関数の例を思い出すと, 型推論の過程で  $\varepsilon \mid a\alpha b \leq \alpha$  のような制約が得られるが, これは正規表現の範囲では最小の解を持たない. なぜなら  $a^*b^* \geq (\varepsilon \mid aa^*b^*b) \geq (\varepsilon \mid ab \mid aaa^*b^*bb) \geq (\varepsilon \mid ab \mid aabb \mid aaaa^*b^*bbb) \geq \dots$  などがすべて有効な解となるからである. 我々は目下“適度な”近似解を与える型推論アルゴリズムの構築に向けて, 自然言語処理の分野から文献(10), (11)など類似の研究の成果を取り入れながら作業中である.

## 参 考 文 献

- 1) CERT Advisory CA-2000-02 (2000). <http://www.cert.org/advisories/CA-2000-02.html>
- 2) Frisch, A., Castagna, G. and Benzaken, V.: Semantic subtyping, *Logic in Computer Science*, Los Alamitos, CA, USA, pp.137-146 IEEE Computer Society (2002).
- 3) Harper, R.: Proof-Directed Debugging, *Journal of Functional Programming*, Vol.9, No.4, pp.463-469 (1999).
- 4) Hopcroft, J.E., Motwani, R. and Ullman, J.D.: *Introduction to Automata Theory, Languages, and Computation*, 2nd edition, Addison-Wesley Publishing (2000).
- 5) Hosoya, H. and Pierce, B.C.: Regular expression pattern matching for XML, *Journal of Functional Programming* (2002). To appear. Extended abstract in *Symposium on Principles of Programming Languages*, pp.67-80 (2001).
- 6) Hosoya, H., Vouillon, J. and Pierce, B.C.: Regular expression types for XML, *International Conference on Functional Programming*, pp.11-22 (2000).
- 7) Igarashi, A. and Kobayashi, N.: Resource Us-

age Analysis, *Symposium on Principles of Programming Languages* (2002).

- 8) Lutz, M.: *Programming Python*, 2nd edition, O'Reilly & Associates (2001).
- 9) Matsumoto, Y.M. and Ishituka, K.: *The Ruby Programming Language*, Addison-Wesley Publishing (2002).
- 10) Mohri, M. and Nederhof, M.-J.: Regular approximation of context-free grammars through transformation, *Robustness in Language and Speech Technology*, Junqua, J.-C. and Noord, G.V.(Eds.), Kluwer Academic Publishers (2001).
- 11) Nederhof, M.-J.: Regular Approximation of CFLs: A Grammatical View, *Advances in Probabilistic and Other Parsing Technologies*, Bunt, H.C. and Nijholt, A.(Eds.), Kluwer Academic Publishers (2000).
- 12) Thomas, D., Hunt, A. and Thomas, D.: *Programming Ruby: A Pragmatic Programmer's Guide*, Addison-Wesley Publishing (2000).
- 13) Wall, L., et al.: *Programming Perl*, 3rd edition, O'Reilly & Associates (2000).
- 14) 山本 篤：正規表現の関数としての定義とパターンマッチアルゴリズム (2001).  
<http://aglaia.c.u-tokyo.ac.jp/~yamamoto/regex.function/>

(平成 14 年 5 月 23 日受付)

(平成 14 年 8 月 9 日採録)



田淵 直

1977 年生 . 2001 年 3 月京都大学文学部心理学専攻卒業 . 同年 4 月東京大学情報理工学系研究科コンピュータ科学専攻修士課程進学 . ソフトウェアの自動検証等の領域にお

いて , 型システムの理論および応用に関する研究に従事 .



住井英二郎

1975 年生 . 2000 年 3 月東京大学大学院理学系研究科情報科学専攻修士課程修了 . 同年 4 月同専攻博士課程進学 . 同年より翌年にかけて , 米国ペンシルバニア大学訪問研究員 .

2001 年 4 月より東京大学大学院情報理工学系研究科コンピュータ科学専攻助手 ( 2003 年 3 月まで同大学院情報学環へ流動 ) . 関数型言語 , プロセス計算 , 部分評価 , 情報セキュリティ等の領域における , 先進的な型システムの理論と応用に関する研究に従事 . 日本ソフトウェア科学会員 . ACM 会員 .



米澤 明恵 ( 正会員 )

1947 年生 . 1977 年 Ph.D. in Computer Science ( MIT ) . 1989 年より東京大学理学部情報科学科教授 . 超並列・分散ソフトウェアアーキテクチャ等に興味を持つ . 共著書「モデルと表現」等 ( 岩波書店 ) , 編著書「ABCL」( MIT Press ) 等がある . 1992 年 ~ 1996 年ドイツ国立情報処理研究所 ( GMD ) 科学顧問 , ACM Transaction on Programming Languages and Systems 副編集長 , IEEE Parallel & Distributed Technology および Computer 編集委員等を歴任 , 元ソフトウェア科学会理事長 , ACM Fellow .