# Semantics of Regular Expressions

## 1 Operational Semantics

$$\frac{\vdash r_1 \text{ matches } s_1 \text{ leaving } s_2 \quad \vdash r_2 \text{ matches } s_2 \text{ leaving } s_3}{\vdash r_1 r_2 \text{ matches } s_1 \text{ leaving } s_3}$$

$$\frac{\vdash r_1 \text{ matches } s_1 \text{ leaving } s_2}{\vdash r_1 | r_2 \text{ matches } s_1 \text{ leaving } s_2}$$

$$\frac{\vdash r_2 \text{ matches } s_1 \text{ leaving } s_2}{\vdash r_1 | r_2 \text{ matches } s_1 \text{ leaving } s_2}$$

$$\frac{}{\vdash r_1 * \text{ matches } s_1 \text{ leaving } s_1}$$

$$\frac{\vdash r \text{ matches } s_1 \text{ leaving } s_2 \quad \vdash r * \text{ matches } s_2 \text{ leaving } s_3}{\vdash r_1 * \text{ matches } s_1 \text{ leaving } s_3}$$

## 2 Denotational Semantics

### 2.1 Disjunction

$$\mathcal{R}[\![\, r_1 | r_2 \,]\!](s) = \mathcal{R}[\![\, r_1 \,]\!](s) \ \cup \ \mathcal{R}[\![\, r_2 \,]\!](s)$$

or, equivalently:

$$\mathcal{R}[\![\, r_1 | r_2 \,]\!](s) = \{x \mid x \in \mathcal{R}[\![\, r_1 \,]\!](s) \vee x \in \mathcal{R}[\![\, r_2 \,]\!](s)\}$$

### 2.2 Concatenation

$$\mathcal{R}[\![\, r_1 r_2 \,]\!](s) = \{x \mid \exists y.\ y \in \mathcal{R}[\![\, r_1 \,]\!](s) \wedge x \in \mathcal{R}[\![\, r_2 \,]\!](y)\}$$

or, equivalently:

$$\mathcal{R}[\![\, r_1 r_2 \,]\!](s) = \bigcup_{y \in \mathcal{R}[\![\, r_1 \,]\!] s} \mathcal{R}[\![\, r_2 \,]\!](y)$$

### 2.3 Kleene Closure

Let $r^0 \equiv \mathsf{empty}$ and $r^n \equiv r_1 r_2 \ldots r_n$ (i.e., $r$ concatenated with itself $n$ times).

$$\mathcal{R}[\![\, r* \,]\!](s) = \bigcup_{k \in 0\ldots\infty} \mathcal{R}[\![\, r^k \,]\!](s)$$

or, equivalently:

Consider the unwinding equation $r* \equiv rr*$. We define a context $C$ (a regexp with a hole) so that $C \equiv r\bullet$. Note that $r* \equiv C[r*]$. The meaning of a context is a semantic function $F$ such that $F[\![\, C[r*] \,]\!] = F[\![\, r* \,]\!]$. The type of $F$ is:

$$F : (S \to \mathcal{P}(S)) \to (S \to \mathcal{P}(S))$$

We want the least fixed point of $F$, where *least* is interpreted with respect to set inclusion $\subseteq$. We assert that $F$ is monotonic and continuous. Let $F^0(W) = \mathcal{R}[\![\, \mathsf{empty} \,]\!] = \lambda s.\{s\}$. We define $F^{k+1}$ as follows:

$$F^{k+1}(W) = FF^k(W) = \lambda s. \bigcup_{y \in \mathcal{R}[\![\, r \,]\!](s)} F^k(y)$$

Then we want the least fixed point:

$$\mathcal{R}[\![\, r* \,]\!](s) = \bigsqcup_k F^k(\lambda s.\{s\}) = \bigcup_k F^k(\lambda s.\{s\})$$

## 3 Incorrect Answers

The following definition of Kleene star is *incorrect*:

$$\mathcal{R}[\![\, r* \,]\!](s) \neq \{s\} \cup \mathcal{R}[\![\, rr* \,]\!]$$

Using the rule for concatenation above, it is equivalent to the following also-*incorrect* definition:

$$\mathcal{R}[\![\, r* \,]\!](s) \neq \{s\} \cup \{x \mid \exists y.\ y \in \mathcal{R}[\![\, r \,]\!](s) \wedge x \in \mathcal{R}[\![\, r* \,]\!](y)\}$$

The definitions are *incorrect* because they define $\mathcal{R}[\![\, r* \,]\!]$ directly in terms of itself. Such circular definitions correspond to implementation code such as:

```
1   | Star(r) -> (* incorrect *)
2       matches (Or(Empty,Concat(r,Star(r)))) s
```

On regular expressions such as $r = \mathsf{empty}*$, this leads to an infinite loop (and usually a stack overflow).

There are two typical approaches for a correct implementation. The first chooses some large $k$ (say, based on the length of the input string $s$) and computes $\cup_{i=0..k} \mathcal{R}[\![\, r^k \,]\!](s)$. The second actually computes the fixed point (instead of picking $k$ in advance) by repeating the process until nothing new is added to the answer.

Regular expression matching is used almost everywhere. Note that understanding the denotational semantics actually helps one to write a real-world program correctly.