

Regular Expressions I

*Languages, Equality Types, and
the `regexp` type*

15-150 M21

Lecture 0707

07 July 2021

0 Decision Problems

Birth Year of Computation



1936

January							February							March						
Su	Mo	Tu	We	Th	Fr	Sa	Su	Mo	Tu	We	Th	Fr	Sa	Su	Mo	Tu	We	Th	Fr	Sa
			1	2	3	4							1	1	2	3	4	5	6	7
5	6	7	8	9	10	11	2	3	4	5	6	7	8	8	9	10	11	12	13	14
12	13	14	15	16	17	18	9	10	11	12	13	14	15	15	16	17	18	19	20	21
19	20	21	22	23	24	25	16	17	18	19	20	21	22	22	23	24	25	26	27	28
26	27	28	29	30	31		23	24	25	26	27	28	29	29	30	31				
1:○	8:○	16:○	24:●	30:○			7:○	15:○	22:●	29:○				8:○	16:○	22:●	29:○			

April							May							June						
Su	Mo	Tu	We	Th	Fr	Sa	Su	Mo	Tu	We	Th	Fr	Sa	Su	Mo	Tu	We	Th	Fr	Sa
			1	2	3	4						1	2	1	2	3	4	5	6	
5	6	7	8	9	10	11	3	4	5	6	7	8	9	7	8	9	10	11	12	13
12	13	14	15	16	17	18	10	11	12	13	14	15	16	14	15	16	17	18	19	20
19	20	21	22	23	24	25	17	18	19	20	21	22	23	21	22	23	24	25	26	27
26	27	28	29	30			24	25	26	27	28	29	31	28	29	30				
6:○	14:○	21:●	28:○				6:○	14:○	20:●	27:○				5:○	12:○	19:●	26:○			

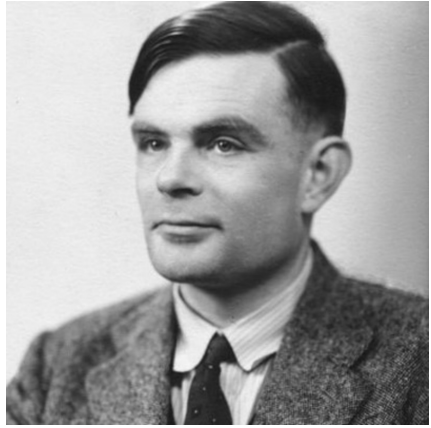


**Warning: Great Theoretical
Ideas**

When working with Turing computability, we assume we have some (usually finite) set Σ – our **alphabet**

We'll be computing with the set Σ^* of all finite strings/sequences/lists of elements of Σ (“strings over Σ ”). For instance,

$$\{a\}^* = \{\epsilon, a, aa, aaa, aaaa, aaaaa, aaaaaa, \dots\}$$



Computation is done by **machines**, which take some input (encoded as a string over some alphabet Σ) and either accepts it or rejects it.

$$M : \Sigma^* \rightarrow \{\mathbf{ACCEPT}, \mathbf{REJECT}\}$$

This is a partial function, because it could loop forever on some inputs.

For any Turing machine M , we write

$$\mathcal{L}(M) = \{s \in \Sigma^* \mid M(s) = \mathbf{ACCEPT}\}$$

for the **language** of M .

Question: Are all subsets $L \subseteq \Sigma^*$ **computable**: is there a Turing Machine M such that $L = \mathcal{L}(M)$?

Thm. (Turing, 1936) No!

- Σ^* is countably-infinite, so $\mathcal{P}(\Sigma^*)$ is uncountably infinite. But there are only countably-many possible Turing Machines
- There are explicit subsets of Σ^* which are not computable, e.g. **HALTS**

We can have a similar idea in functional programming:

$$M : t \rightarrow \text{bool}$$
$$M : \text{Sigma list} \rightarrow \text{bool}$$

where `Sigma` is the alphabet type (e.g. `char`).

Question: For which sets L of values of type `Sigma list` is there an SML function $M : \text{Sigma list} \rightarrow \text{bool}$ such that

$$L = \{v : \text{Sigma list} \mid M(v) \implies \text{true}\}?$$

Thm. For each finite set Σ (with corresponding SML type `Sigma`), and each subset L of Σ^* the following are equivalent:

- there exists a Turing Machine M such that $L = \mathcal{L}(M)$
- there exists an SML function `M : Sigma list -> bool` such that

$$L = \{v : \text{Sigma list} \mid M(v) \implies \text{true}\}$$

Decidable Equality

A type t is said to be an **equality type** if there is a total function $(op =) : t * t \rightarrow bool$ deciding whether elements of that type are equal or not.

Examples: `int`, `bool`, `char`, `string`, `int list`

Non-Examples: `real`, `int -> int`

We can specify that a polymorphic type variable *must* be instantiated to an equality type by writing it with double tick-marks:

```
(op =) : ''a * ''a -> bool
Fn.equal : ''a -> ''a -> bool
Fn.notEqual : ''a -> ''a -> bool
```

Module: Language

[https://github.com/smlhelp/aux-library/
blob/main/Language.sml](https://github.com/smlhelp/aux-library/blob/main/Language.sml)

aux-library/Language.sml

```
33 fun everything (x:'S list) = true
```

For each type `Sigma`, we identify each `L : Sigma language` with the set of values `s : Sigma list` such that `L(v) \implies true`.

aux-library/Language.sml

```
7  val everything : 'S language
8  val nothing   : 'S language
9  val singleton : ''S list -> ''S language
```

aux-library/Language.sml

11

```
val Or
```

12

```
  : 'S language * 'S language -> 'S language
```

13

```
val And
```

14

```
  : 'S language * 'S language -> 'S language
```

15

```
val Not
```

16

```
  : 'S language -> 'S language
```

17

```
val Xor
```

18



Goal: Write more complex
languages

5-minute break

Take `Sigma` to be `char`. We have a special way of dealing with `char lists...strings!`

```
String.explode  : string -> char list  
String.implode  : char list -> string
```

`aux-library/Language.sml`

1 Regular Expressions

Idea:

Come up with an SML datatype whose elements encode different languages

We'll define a datatype parametrized by a single equality type variable

```
datatype 's regexp = ...
```

such that each value $R : t \text{ regexp}$ defines a language

$$\mathcal{L}(R) = \{v : t \text{ list} \mid v \text{ "matches with" } R\}$$

Module: `Regex`

`https://github.com/smlhelp/aux-library/
blob/main/Regex.sml`

(requires `Language.sml` in the same directory)

```
26 end
27 structure Regexp : REGEXP =
28 struct
29   datatype 'S regexp =
30     Zero
31     | One
32     | Const of 'S
```

Note: you can't actually require that the parameter of a datatype be an equality type. SML will treat this the same as `datatype 'S regexp = ...`, but the `'S` reminds us to use this with equality types

A datatype like any other

aux-library/Regexp.sml

```
34     | Times of ''S regexp * ''S regexp
35     | Star  of ''S regexp
36
37 fun depth Zero = 0
38     | depth One = 0
39     | depth (Const _) = 0
40     | depth (Plus(R1,R2)) =
41         1 + Int.max(depth R1,depth R2)
42     | depth (Times(R1,R2)) =
```


Idea: For each $R : \text{Sigma regexp}$, define the language $\mathcal{L}(R)$ to be the set of all values of type Sigma list which “match” or “are accepted” by R . We’ll do this recursively based on R .

LL : ' 'S regexp -> ' 'S language

REQUIRES: true

ENSURES: LL (R) is a total function deciding $\mathcal{L}(R)$, i.e. LL R s \implies true for all $s \in \mathcal{L}(R)$, and LL R s \implies false for all $s \notin \mathcal{L}(R)$.

- `Const v` only matches with `[v]`

$$\mathcal{L}(\text{Const } v) = \{ [v] \}$$

- `One` only matches with `[]`

$$\mathcal{L}(\text{One}) = \{ [] \}$$

- `Zero` does not match with anything

$$\mathcal{L}(\text{Zero}) = \emptyset$$

- `Plus (R1 , R2)` matches with any list which matches either `R1` or `R2`

$$\mathcal{L}(\text{Plus}(R1, R2)) = \mathcal{L}(R1) \cup \mathcal{L}(R2)$$

- `Times (R1 , R2)` matches with any list consisting of an `R1` list appended to an `R2` list

$$\mathcal{L}(\text{Times}(R1, R2)) = \{v1@v2 \mid v1 \in \mathcal{L}(R1) \text{ and } v2 \in \mathcal{L}(R2)\}$$

- $\text{Star}(R)$ matches with any list consisting of finitely-many R -matching lists appended together

$$\mathcal{L}(\text{Star}(R)) = \{v_1 @ v_2 @ \dots @ v_n \mid n \in \mathbb{N} \text{ and } v_i \in \mathcal{L}(R) \text{ for each } 1 \leq i \leq n\}$$

Note: $[] \in \mathcal{L}(\text{Star}(R))$ for all R .

Next Time:

How Regex .LL is implemented

Thank you!