Regular expressions can be used to

- Specify a set of strings
- Search a text for patterns
- Generate responses for a simple discourse agent
- Specify an entire human language (like French)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 -

Write recursive algorithms

- An **alphabet** is a set of symbols, Σ .
- A string over an alphabet is a sequence of symbols from that alphabet. Σ* is the set of all strings over alphabet Σ.
- A language over an alphabet is a set of strings, that is, a subset of Σ*.

Regular expressions constitute a system for specifying languages. (J&M, "a language for specifying text search strings", pg 3.). An individual regular expression denotes a language, that is, a set of strings.

▲ロト ▲圖ト ▲画ト ▲画ト 三回 - のへで

base cases	\emptyset $arepsilon$	the empty set of strings the set containing the empty string, $\{""\}$ the set containing only the string with only <i>a</i> , for some $a \in \Sigma$, $\{"a"\}$
recursive cases	r s	the set of strings made from concatening strings from r and s , $\{x + y \mid x \in r \land y \in s\}$, for some regular expressions r and s the set of strings from r or s , $r \cup s$ for some regular expressions r and s the set of strings made from concatenating 0 or more strings from r for some regular expression r

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Abbreviation	Meaning	Equivalence
[abc]	One occurrence of any of these symbols	(a b c)
[a-c]	One occurrence of any symbol in this range	(a b c)
r?	Optionally an occurrence of a string defined by r	(r arepsilon)
r ⁵	5 occurrences of a string defined by r	rrrrr
r ^{3,5}	Between 3 and 5 occurrences of a string defined by r	(rrr rrrr rrrr)
r+	One or more occurrences of a string defined by <i>r</i>	rr*

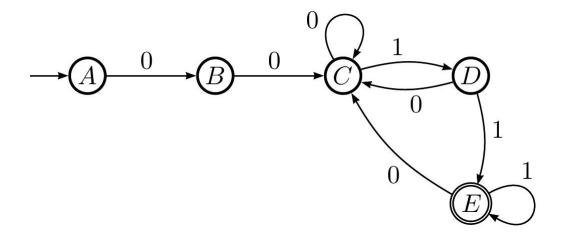
▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

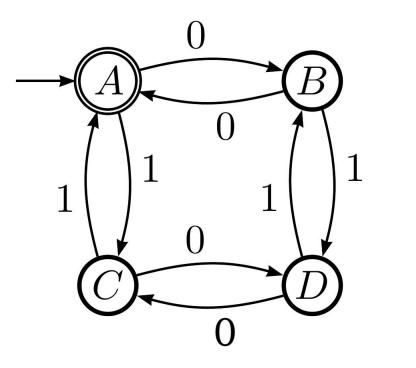
- DNA sequences: (A|C|G|T)*.
- Identifiers: (('|ε) [A-Za-z] [A-Za-z0-9_])|_.
- Phone numbers: $[2-9][0-9]^2 [2-9][0-9]^2 [0-9]^4$.
- ► Dates: ((1[0-2])|[1-9])/(30|31|([12][0-9])|[1-9])/[1-9][0-9]^{0,3}. |
- ► US Postal Addresses: [0-9] + [NSEW]^{0,2} [A-Z] [a-z] * (St|Ave|Rd|Ln|Dr| Blvd), ([A-Z] [a-z] *)*, [A-Z]² [0-9]⁵.

◆□▶ ◆□▶ ◆注▶ ◆注▶ 注 のへで

Lord, you have been our dwelling place in all generations.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?





▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで