

Lecture 4: Finite Automata

st select an open discussion section to be enrolled. Please
ately. You will not get off the waitlist until you do.
ould now have a Unix instructional account. If you don't,
ve teams formed by Friday, if possible, but next Monday
ork due next Wednesday (31 Jan).

Some Abbreviations

rm from the last slide is good for formal analysis, but
ng.
allow some abbreviations that are obviously expandable
c forms:

Abbreviation	Meaning
$\dots \mathcal{R}_n$	$A : \mathcal{R}_1$ \vdots $A : \mathcal{R}_n$
$\mathcal{R}) \dots$	$B : \mathcal{R}$ $A : \dots B \dots$
$" \mid \dots \mid "c_n"$	$[c_1 \dots c_n]$
other character classes)	

Public-Service Announcement

Experimental Social Science Laboratory (Xlab) invites
participate in social science studies! Experiments con-
ab (located in Hearst Gym, Suite 2) are computer-
n-making studies such as tasks, surveys, and games.
asionally offer remote online and mobile studies that
leted anywhere. Participants earn \$15/hour on av-
time they participate. For more information, visit
y.edu. To sign up, visit berkeley.sona-systems.com.

Formative Style for Describing Languages

giving a single pattern, we can give a set of rules of the

$$A : \alpha_1 \alpha_2 \dots \alpha_n, \quad n \geq 0,$$

mbol that is intended to stand for a language (set of
a *metavariable* or *nonterminal symbol*.

either a literal character (like "a") or a nonterminal

tation of this rule is

to form a string in $L(A)$ (the language denoted by A) is
enerate one string each from $L(\alpha_1), L(\alpha_2), \dots$.

) is just the language $\{ "c" \}$.

us-Naur Form (BNF). A set of rules is a *grammar*. One
rminals is designated as its *start symbol* denoting the
cribed by the grammar.

see that ' : ' written many different ways, such as ' : = ',
e'll just use the same notation our tools use.

A Big Restriction (for now)

being, we'll also add a restriction. In each rule:

$$A : \alpha_1\alpha_2 \cdots \alpha_n, \quad n \geq 0,$$

that if α_i is a nonterminal symbol, then either

es for that symbol have to occurred before all the rules

is the last item) and α_n is A .

n a restricted grammar a **Type 3** or *regular* grammar.
is definable by regular grammars are called *regular lan-*

regular languages are *exactly* the ones that can be de-
regular expressions.

Proof of Claim (II)

uct... .. with Q , where

$$Q : \\ Q : R \quad Q$$

$$Q : R \\ Q : R \quad Q$$

$$Q : \\ Q : R$$

Some Technicalities

inition, each nonterminal in a grammar defines a lan-
n, we are interested in just one of them (the *start*
the others are auxiliary definitions.

on of what a rule means ("One way to form a string in
eaves open the possibility that there are other ways to
n $L(A)$ than covered in the rule.

it freedom in order to allow multiple rules for A , but we
want to include strings that aren't covered by some rule.
athematical definitions throw in sentences like:

r defines the *minimal* languages that contain all strings
fy the rules.

Proof of Claim (I)

a regular expression, R , and make a (possibly not yet

(preceding) rule for each parenthesized expression.

re just the constructs ' X^* ', ' X^+ ', and ' $X^?$ '. What do we
n?

Side Note: The Empty Language

for the empty language is a bit non-intuitive:
strings, Q , satisfies this rule.
the implicit rule that we choose the *smallest* solution that
rules, Q represents the empty set.

Review: FA operation

ph whose nodes are *states (of memory)* and whose edges
transitions. There are a finite number of nodes.
the designated *start state*.
of the nodes are *final states*.
ion is labeled with a set of symbols (characters, etc.) or
lizes a string $c_1c_2 \dots c_n$ if there is a path (sequence of
the start state to a final state such that the labels
s in sequence, aside from ϵ edges, respectively contain
s leaving any node have disjoint sets of characters and
no ϵ nodes, FA is a DFA, else an NFA.

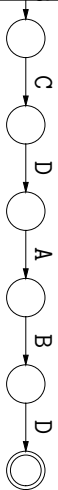
Example

regular expression $("+" | "-") ? ("0" | "1") +$
 $"+" | "-") ? ("0" | "1") +$ replace with ...
: "+" | "-"
: "0" | "1"
 $Q_1 ? Q_2 +$ replace with ...
 $Q_3 : \epsilon | Q_1$
 $Q_4 : Q_2 | Q_2 Q_4$
 $R : Q_3 Q_4$

cal Pattern-Matching Implementation

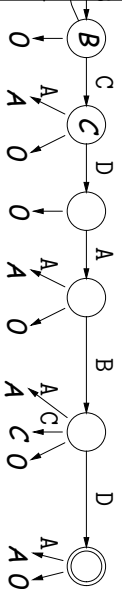
s, can generally make do with "classical" regular expres-
sion using *finite(-state) automata* or *FAs*. ("Finite state"
nary").
struction:
pression \Rightarrow nondeterministic FA (NFA)
nistic FA (DFA) \Rightarrow table-driven program.

Example: What does this NFA recognize?



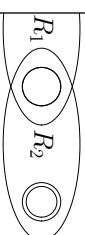
Strings of capitals ending in ABCDABD.

Simplest equivalent DFA you can think of?

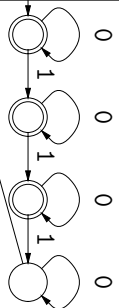


Labels mean "any character not covered by another edge."

Classical Regular Expressions to NFAs (I)



Example: What does this DFA recognize?

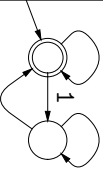


Bit strings with # of 1's divisible by 2 or 3.

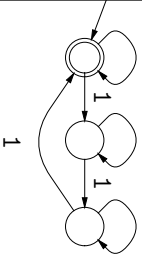
1

Simplest equivalent NFA you can think of?

0 0

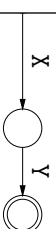


0 1 0 0



1

Example: What does this NFA recognize?



[XY]

Simplest equivalent DFA you can think of?

Extensions?

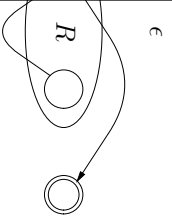
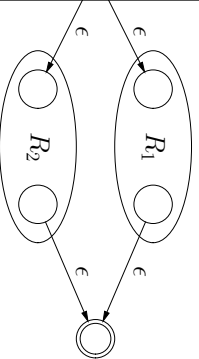
Can we translate ϕ (the empty language, containing no strings)

Can we translate 'R?' into an NFA?

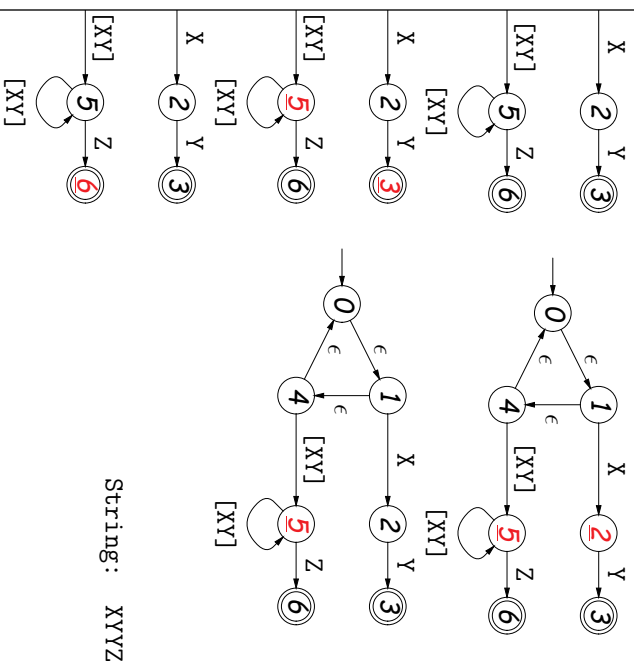
Can we translate 'R+' into an NFA?

Can we translate ' $R_1|R_2|\dots|R_n$ ' into an NFA?

Classical Regular Expressions to NFAs (II)

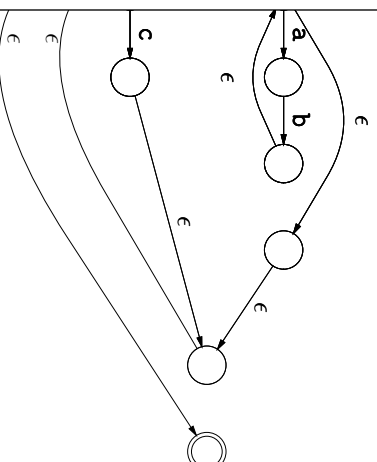


Abstract Implementation of NFAs



Example of Conversion

Translate $((ab)^*|c)^*$ into an NFA (using the construction)



DFA as Programs

DFA in program with control structure:

```
INITIAL;
input; *s != '\0'; s += 1) {
    (state):
    INITIAL;
    *s == 'a') state = A_STATE; break;
    STATE:
    *s == 'b') state = B_STATE; else state = INITIAL; break;
```

```
ce == FINAL1 || state == FINAL2;
```

structure (table driven):

```
INITIAL;
input; *s != '\0'; s += 1)
    transition[state][s];
    final[state];
```

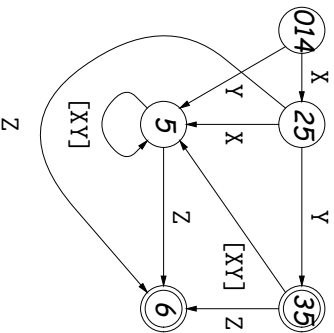
How Do They Do It?

use a DFA to recognize longest match?
 use DFA to act on first of equal-length matches?
 use a DFA to handle the R_1/R_2 pattern (matches just
 f followed by R_2 , like R_1 ($?=R_2$) in Python)?

Review: Converting to DFA

ON: The **set of states** that are marked (colored red)
 each character in a way that depends only on the set
 after.

ds, machine on previous slide acted like this DFA:



What Flex Does

n specification is giant regular expression of the form
 $R_1 | R_2 | \dots | R_n$, where none of the R_i match ϵ .
 ate labeled with some action.
 y previous methods, into a table-driven DFA.
 ticular DFA is used to recognize *prefixes* of the (re-
 t: initial portions that put machine in a final state.
 state(s) we end up in determine action. To deal with
 ons:
 ggest prefix ("maximum munch").
 are multiple matches, apply *first* rule in order.