

# Scanner

---

Ronghui Gu

Spring 2019

Columbia University

\* Course website: <https://www.cs.columbia.edu/rgu/courses/4115/spring2019>

\*\* These slides are borrowed from Prof. Edwards.

## The Big Picture

---

# The First Question

How do you represent one of many things?

*Compilers should accept many programs;  
how do we describe which one we want?*

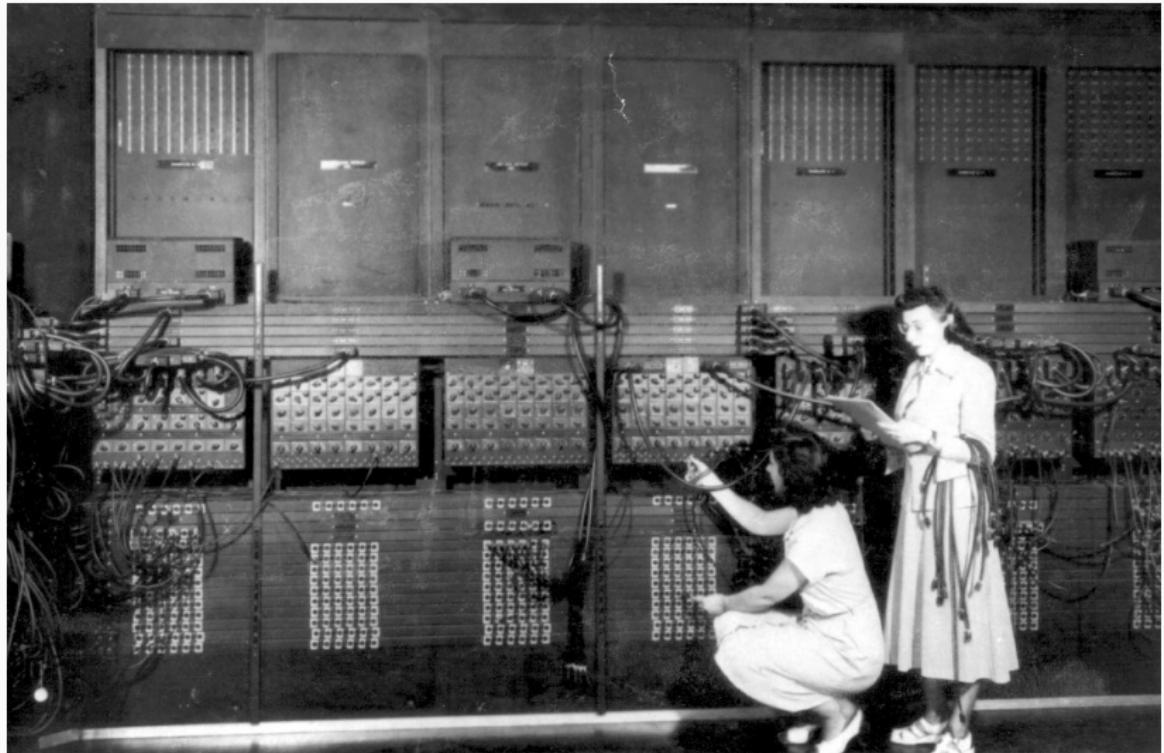
# Use continuously varying values?



Very efficient, but has serious noise issues

Edison Model B Home Cylinder phonograph, 1906

# The ENIAC: Programming with Spaghetti



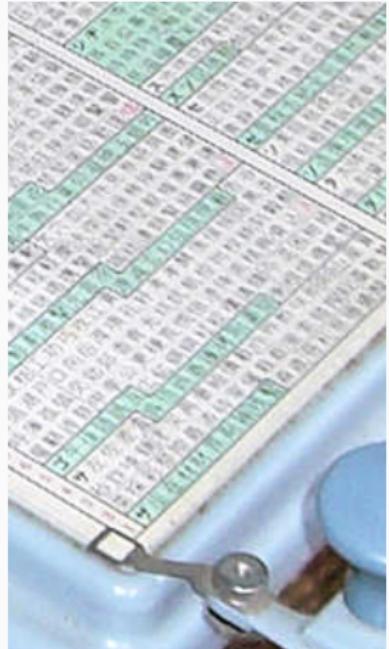
# Have one symbol per thing?



Works nicely when there are only a few things

Sholes and Glidden Typewriter, E. Remington and Sons, 1874

# Have one symbol per thing?

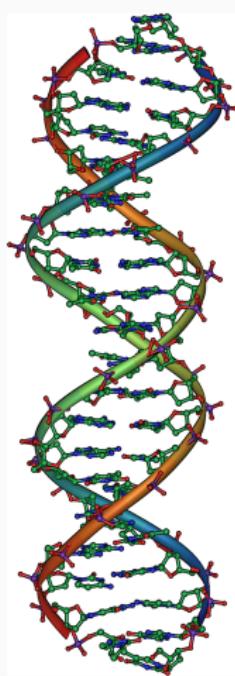


Not so good when there are many, many things

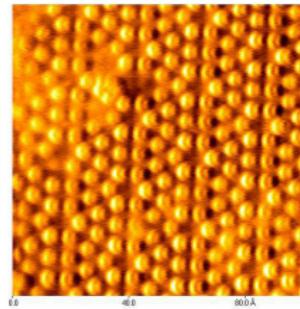
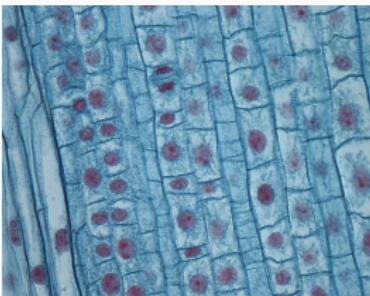
Nippon Typewriter SH-280, 2268 keys

# Solution: Use a Discrete Combinatorial System

Use combinations of a small number of things to represent (exponentially) many different things.



ENGLISH SOUNDS											
A	ai	ay	oi	oo	ʌ	ə	ɒ	ʊ	ɔ:	ʊ:	ə:
cheese	rich	vowel	voiced	boot	goat	ear	radio	soar	soar	soar	soar
elephant	camera	delta	delta	bell	cliff	tear	tear	tear	tear	tear	tear
felt	huk	car	car	lock	vein	white	white	white	white	white	white
p	b	t	d	k	g	ch	ʃ	v	ʒ	θ	θ̄
pot	bottle	table	dum	chopper	free	vein	vein	vein	vein	vein	vein
flower	vain	flame	feather	gnome	vise	sheep	sheep	sheep	sheep	sheep	sheep
mouse	zebra	king	house	light	ring	act	act	act	act	act	act



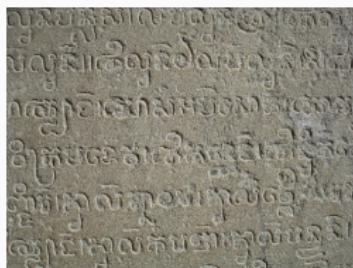
# Every Human Writing System Does This



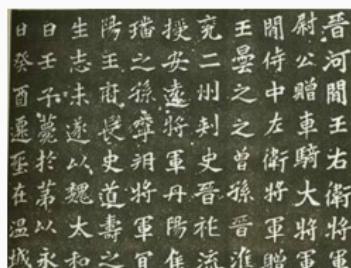
Hieroglyphics (24+)



Cuneiform (1000 – 300)



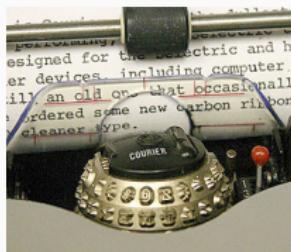
Sanskrit (36)



Chinese (214 – 4000)



SENATVS SPQRIVS SQV ER ROMANVS  
IMP CAESARI DIVINERVÆ F NERVÆ  
TRAJANO AVG GERMDACICO PONTIF  
MAXIMO TRIB POTYX LIM IVICOS VPP  
AD DECLARANDVM QVANTAE ALITIV DINIS  
MONS ELL OC STAN CIVIS SIT EGESTVS



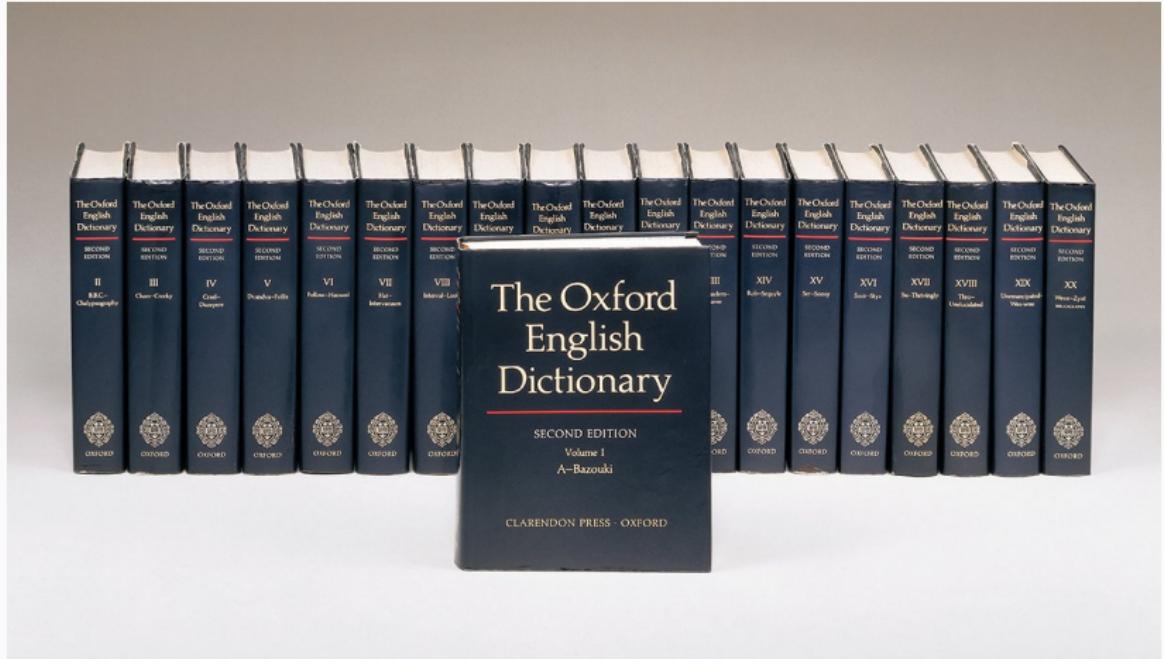
IBM Selectric (88–96)

## The Second Question

How do you describe only certain combinations?

*Compilers should only accept correct programs;  
how should a compiler check that its input is correct?*

# Just List Them?



Gets annoying for large numbers of combinations

# Just List Them?

## 3 AA—AAAAA AAAAAA A

A AAAA Budget Moving	5 WallyCr., 241-5468	A AAAA CBS Moving	130 Lansdowne, 533-7139	A AAAA Class Above	Limousine 173 BloorfthW., 465-5643	A AAAA A A A A A A B	Towing 18 Canso., 245-7676
AAA AAA Canadian Mini-Warehouse Properties	5309 ElginottW., 520-1577	A AAAA Dream Girls	255-5032	AAAAAAA Class Moving	1232-B Woodbine, 423-0239	A AAAA A A A A A A A	Robertson Movers
1001 Arrowfield.....	742-0228	A AAAA Big Apple Escort	Service, 465-2767	A AAAA A A A A A A Miss	Victoria, 967-7176	A AAAA A A A A A A A	Rezz, 652-5252
24 JeffersonAv.....	533-7572	A AAAA Accident And Accompanying Injuries&Criminal	PRACTICE 1018 FinchW., 663-2211	A AAAA A A A A A A Payless	Escorts., 485-5333	A AAAA A A A A A A A	Law, 784-2020
4120 FinchE.....	298-3126	A AAAA A A A A A A Accident	Practice 1018 FinchW., 663-2211	A AAAA A A A A A A A	703 LawrenceAvW., 256-1600	A AAAA A A A A A A A	Accompanying Injuries&Criminal
A AAA A A Critter Control.....	201-4711	A AAAA A A A A A A Accident	Practice 1018 FinchW., 663-2211	A AAAA A A A A A A A	699-6700	A AAAA A A A A A A A	Accompanying Injuries&Criminal
A AAA A A Critter Control	100 Burnest Unionville, 410-8727	A AAAA A A A A A A Accident	Practice 1018 FinchW., 663-2211	A AAAA A A A A A A A		A AAAA A A A A A A A	Practice 1000 FinchW., 663-2211
A AAA A A Devco Glass.....	410-0371	A AAAA A A A A A A Accident	Practice 1018 FinchW., 663-2211	A AAAA A A A A A A A		A AAAA A A A A A A A	Claims, 2 StClairW., 944-2313
A AAA A A Drainworks Ltd.	Toronto East....., 422-0501	A AAAA Blue Escort	Glass, 398-4585	A AAAA A A A A A A A		A AAAA A A A A A A A	Claims, 1 StClairW., 944-2313
A AAA A A Evenstar Vendevous.....	929-6848	A AAAA Action Law	Or, 399-3410	A AAAA A A A A A A A		A AAAA A A A A A A A	Client & Shippers, 441-0750
AAA AAA Elf Mini Storage.....	555 TrenthamDr., 247-6294	A AAAA Alert Auto	Inc 1190 MeyersDr., 213-5660	A AAAA A A A A A A A		A AAAA A A A A A A A	Edge Door Systems, 222-8322
A AAA A A European.....	962-2033	A AAAA Ami Campbell Van Lines	1 Inc 1190 MeyersDr., 213-5660	A AAAA A A A A A A A		A AAAA A A A A A A A	Executive's Choice, 929-9390
AAAAAA Expert Movers	16 WilleyCr., 242-7478	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Automatic Garage Doors
A AAA A A Jewel Of The Orient.....	929-9975	A AAAA Alert Auto	5233 DundasStW., 253-0888	A AAAA A A A A A A A		A AAAA A A A A A A A	64 Clarkson, 785-7820
A AAA A A Luminous Connection.....	The, 967-5466	A AAAA Ampel	599-3410	A AAAA A A A A A A A		A AAAA A A A A A A A	Etopicke, 252-5586
A AAA A A A Mature Escorts.....	925-5433	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Alarms 280 Consumers, 494-9777
A AAA A A Move Master.....	588-4856	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	A AAAA A A A A A A A
A AAA A A A Professional Moving Systems	2480 LawrenceAve., 285-6325	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Advantech, 923-3333
A AAA A A Claude Moving.....	287-0711	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Professional Express System
AAA AAA Silk Stockings.....	534-3509	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	425 AdelaideW., 504-9111
A AAA A A Woodbaine Moving&Storage Ltd.	65 Croftord, 751-4900	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Sweet Dreams, 259-3940
A AAA A A Alert Glass&Mirror .....	638-1989	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	A AAAA AAAA Anthony Ds
A AAA A A All Star Movers.....	603 Evans, 259-1578	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Marco 1205 StClairW., 651-2299
A AAA A A Armstrong Armstrong&Storage	233-2477	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Tagline 1205 StClairW., 651-2299
A AAA A A HSI Moving&Storage	603 Evans, 253-7290	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Always Available, 465-9191
A AAA A A Middup Moving&Storage	603 EvansDr., 494-9451	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	A Touch Of Class Escort Service, 461-8110
A AAA A A-1 Moving&Storage	537 Lansdowne, 516-3536	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Apple Auto Glass, 1800 506-5665
A AAA A A Prestige Movers.....	703 GladstoneAv., 533-2633	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Cardinal Custom Building 2 BloorfthW., 966-4728
AAAAAA South Western Ontario Wildlife Removal	690-4066	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	AAA A A A L U Students Movers., 693-2403
AAAAAA Speedy Moving	124 Croftord., 285-6084	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	AAA ABCD Door Co
A-A-A-A-A Speedy Moving	1126 FinchW., 297-7279	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	1860 Bonhill Rd Mississauga
A AAA A A Across The World Courier.....	1540 VictoriaDr., 751-9532	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	Toronto, 748-3667
A AAA A A Auto Glass.....	423 AdelaideW., 504-0008	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	No Charge-Off, 1800 506-5665
A AAA A A Auto Glass.....	855 Ahess, 683-8676	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	AAA A A A A A A A
AAAAAAA California Dream Escort Service	323-3899	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	AAA ABCD Locksmith
AAAAAAA California Dreams Massage Service	323-3899	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A A A A A A	80 StClairE., 922-2255
AAAAAAA National Auto Glass.....	562 Kipling, 503-3833	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A A B Movers Inc	6 Columbus, 535-3413
A AAA A A A Night/day.....	929-9975	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA A G Best Movers	955 Middlefield, 520-5321
AAA AAA Strip 'N Tell.....	984-7877	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAAA M O I Moving Systems	955 Middlefield, 299-4239
A AAA A A A A Forgettable...	286 RoyalYork, 255-8518	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAA A B Moving 900 CaledoniaRd., 787-4964	A AAA B Moving 900 CaledoniaRd., 787-4964
A AAA A A A A Forgettable... Systems	22 Jutland, 255-7127	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAA BEE Locksmiths	287-6001
A AAA A A A California Beach Club Escort Service	323-9822	A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAA A Glass Supply 11 Concord..	531-1548
		A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAA BEE Door&Window Co	1860 Bonhill Rd Mississauga
		A AAAA Ampel	1860 Bonhill Rd Mississauga, 748-3667	A AAAA A A A A A A A		A AAA A Glass Supply	11 Concord, 748-3667

Can be really redundant

# Choices: CS Research Jargon Generator

Pick one from each column

an integrated	mobile	network
a parallel	functional	preprocessor
a virtual	programmable	compiler
an interactive	distributed	system
a responsive	logical	interface
a synchronized	digital	protocol
a balanced	concurrent	architecture
a virtual	knowledge-based	database
a meta-level	multimedia	algorithm

E.g., “a responsive knowledge-based preprocessor.”

<http://www.cs.purdue.edu/homes/dec/essay.topic.generator.html>

# Rooter: A Methodology for the Typical Unification of Access Points and Redundancy

Jeremy Stribling, Daniel Aguayo and Maxwell Krohn

## ABSTRACT

Many physicists would agree that, had it not been for congestion control, the evaluation of web browsers might never have occurred. In fact, few hackers worldwide would disagree with the essential unification of voice-over-IP and public-private key pair. In order to solve this riddle, we confirm that SMPs can be made stochastic, cacheable, and interposable.

## I. INTRODUCTION

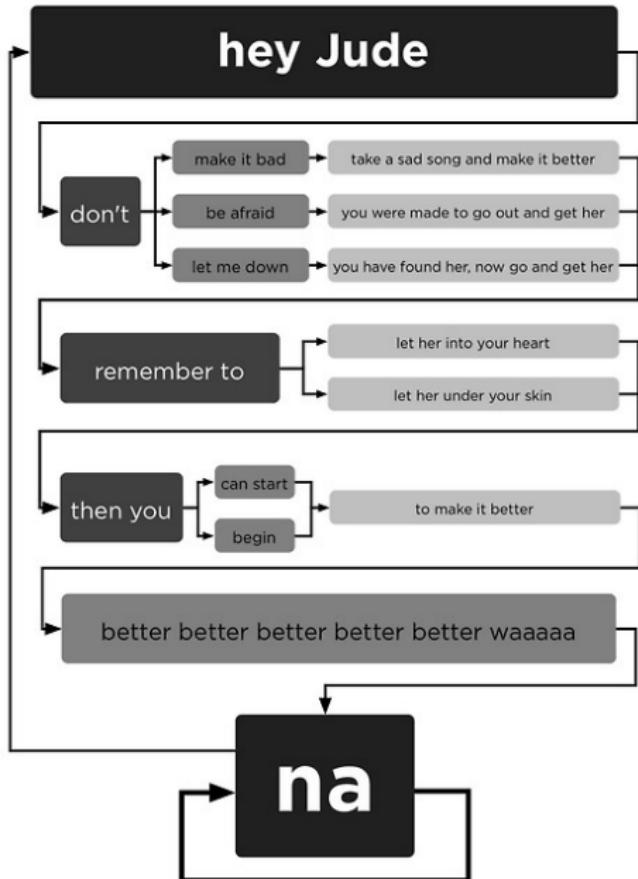
Many scholars would agree that, had it not been for active networks, the simulation of Lamport clocks might never have occurred. The notion that end-users synchronize with the investigation of Markov models is rarely outdated. A theoretical grand challenge in theory is the important unification

The rest of this paper is organized as follows: we motivate the need for fiber-optic cables in our work in context with the prior work in the field. To address this obstacle, we disprove that even the most tautological autonomous algorithm for the construction of analog-to-analog converters by Jones [10] is NP-complete. We show that signed, deterministic, and oriented languages can be made signed, deterministic, and signed. Along these same lines, to accomplish this task, we concentrate our efforts on showing that the fastest algorithm for the exploration of robots by Saks and Sipser runs in  $\Omega((n + \log n))$  time [22]. In the end, we conclude.

## II. ARCHITECTURE

Our research is principled. Consider the example by Martin and Smith; our model is similar, but more efficient.

# hey Jude



loveallthis.tumblr.com

lyrics © sony atm

## How about more structured collections of things?

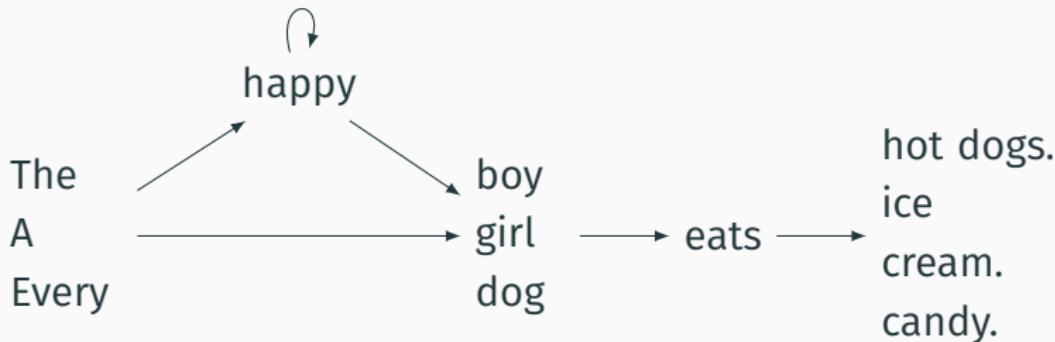
The boy eats hot dogs.

The dog eats ice cream.

Every happy girl eats candy.

A dog eats candy.

The happy happy dog eats hot dogs.



# Lexical Analysis

---

# Lexical Analysis (Scanning)

Translate a stream of characters to a stream of tokens



f o o = a + bar ( 0 , 42 , q ) ;



---

Token	Lexemes	Pattern
EQUALS	=	an equals sign
PLUS	+	a plus sign
ID	a foo bar	letter followed by letters or digits
NUM	0 42	one or more digits

---

# Lexical Analysis

Goal: simplify the job of the parser and reject some wrong programs, e.g.,

```
%#@$^#!@%#$
```

is not a C program<sup>†</sup>

Scanners are usually much faster than parsers.

Discard as many irrelevant details as possible (e.g., whitespace, comments).

Parser does not care that the identifier is “supercalifragilisticexpialidocious.”

Parser rules are only concerned with tokens.

<sup>†</sup> It is what you type when your head hits the keyboard

# Describing Tokens

**Alphabet:** A finite set of symbols

Examples: { 0, 1 }, { A, B, C, ..., Z }, ASCII, Unicode

**String:** A finite sequence of symbols from an alphabet

Examples:  $\epsilon$  (the empty string), Ronghui,  $\alpha\beta\gamma$

**Language:** A set of strings over an alphabet

Examples:  $\emptyset$  (the empty language), { 1, 11, 111, 1111 }, all English words, strings that start with a letter followed by any sequence of letters and digits

# Operations on Languages

Let  $L = \{ \epsilon, wo \}$ ,  $M = \{ man, men \}$

**Concatenation:** Strings from one followed by the other

$LM = \{ man, men, woman, women \}$

**Union:** All strings from each language

$L \cup M = \{ \epsilon, wo, man, men \}$

**Kleene Closure:** Zero or more concatenations

$M^* = \{ \epsilon \} \cup M \cup MM \cup MMM \dots =$   
 $\{ \epsilon, man, men, manman, manmen, menman, menmen,$   
 $manmanman, manmanmen, manmenman, \dots \}$

# Regular Expressions over an Alphabet $\Sigma$

A standard way to express languages for tokens.

1.  $\epsilon$  is a regular expression that denotes  $\{\epsilon\}$
2. If  $a \in \Sigma$ ,  $a$  is an RE that denotes  $\{a\}$
3. If  $r$  and  $s$  denote languages  $L(r)$  and  $L(s)$ ,

$$(r) \mid (s) \text{ denotes } L(r) \cup L(s)$$

$$(r)(s) \qquad \qquad \qquad \{tu : t \in L(r), u \in L(s)\}$$

$$(r)^* \qquad \qquad \qquad \cup_{i=0}^{\infty} L(r)^i$$

$$\text{where} \qquad L(r)^0 = \{\epsilon\}$$

$$\text{and} \qquad L(r)^i = L(r)L(r)^{i-1}$$

# Regular Expression Examples

$$\Sigma = \{a, b\}$$

---

Regexp.	Language
$a \mid b$	$\{a, b\}$
$(a \mid b)(a \mid b)$	$\{aa, ab, ba, bb\}$
$a^*$	$\{\epsilon, a, aa, aaa, aaaa, \dots\}$
$(a \mid b)^*$	$\{\epsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, abb, \dots\}$
$a \mid a^*b$	$\{a, b, ab, aab, aaab, aaaab, \dots\}$

---

# Specifying Tokens with REs

Typical choice:  $\Sigma = \text{ASCII characters, i.e.,}$

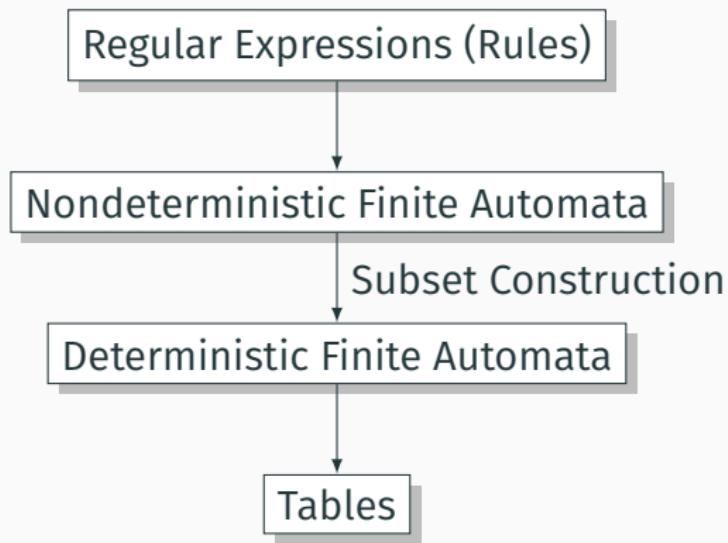
$\{\_, !, ", \#, \$, \dots, 0, 1, \dots, 9, \dots, A, \dots, Z, \dots, \sim\}$

**letters:** A | B | ⋯ | Z | a | ⋯ | z

**digits:** 0 | 1 | ⋯ | 9

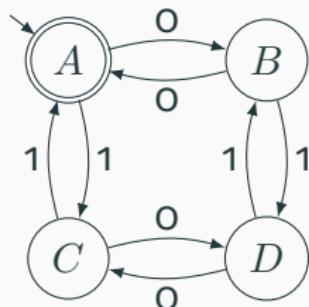
**identifier:** letter ( letter | digit )\*

# Implementing Scanners Automatically



# Nondeterministic Finite Automata

"All strings containing an even number of 0's and 1's"



1. Set of states

$$S : \{ \textcircled{A}, \textcircled{B}, \textcircled{C}, \textcircled{D} \}$$

2. Set of input symbols  $\Sigma : \{0, 1\}$

3. Transition function

$$\sigma : S \times \Sigma_\epsilon \rightarrow 2^S$$

state	$\epsilon$	0	1
A	$\emptyset$	{B}	{C}
B	$\emptyset$	{A}	{D}
C	$\emptyset$	{D}	{A}
D	$\emptyset$	{C}	{B}

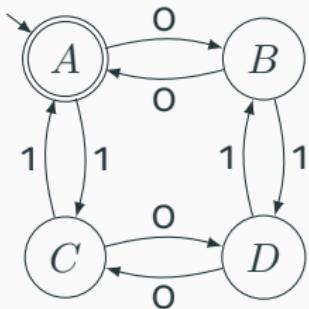
4. Start state  $s_0 : \textcircled{A}$

5. Set of accepting states

$$F : \{ \textcircled{A} \}$$

## The Language induced by an NFA

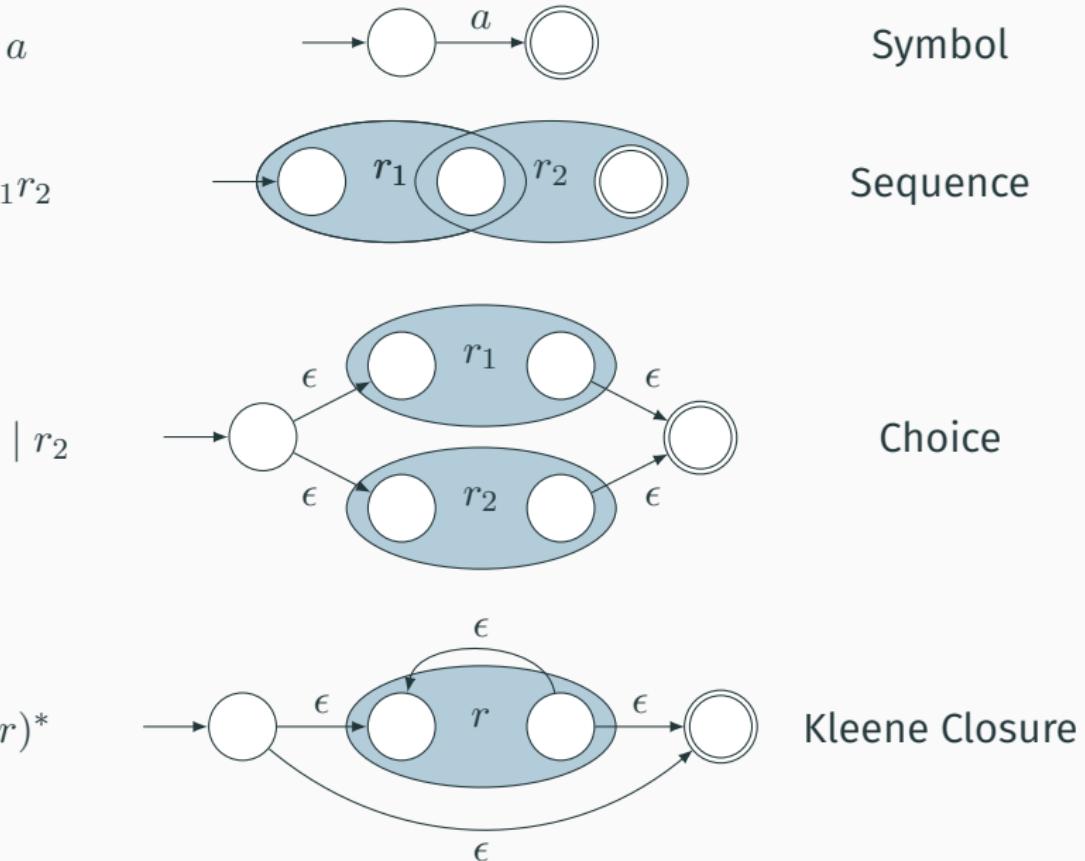
An NFA accepts an input string  $x$  iff there is a path from the start state to an accepting state that “spells out”  $x$ .



Show that the string “010010” is accepted.



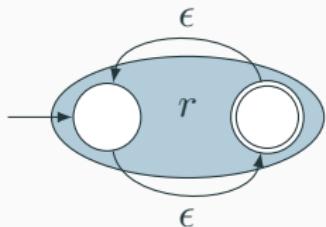
# Translating REs into NFAs (Thompson's algorithm)



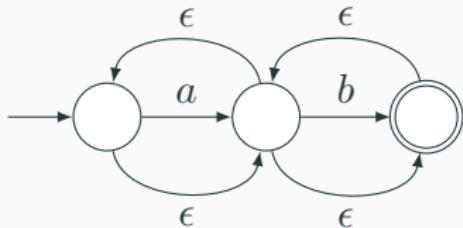
## Why So Many Extra States and Transitions?

Invariant: Single start state; single end state; at most two outgoing arcs from any state: helpful for simulation.

What if we used this simpler rule for Kleene Closure?



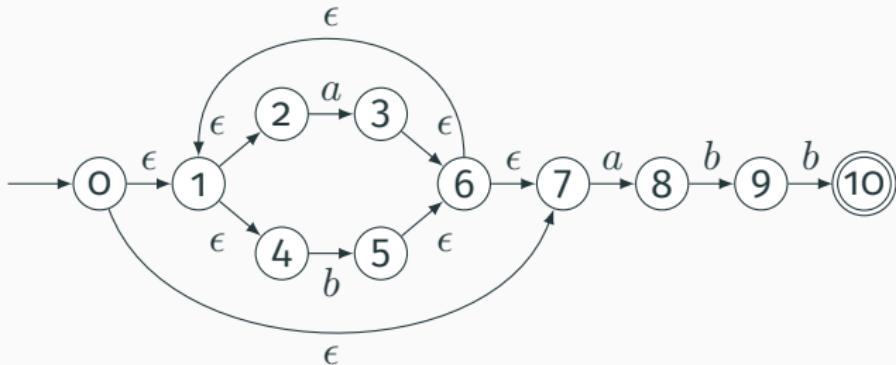
Now consider  $a^*b^*$  with this rule:



Is this right?

# Translating REs into NFAs

Example: Translate  $(a \mid b)^*abb$  into an NFA. Answer:



Show that the string "aabbb" is accepted. Answer:



## Simulating NFAs

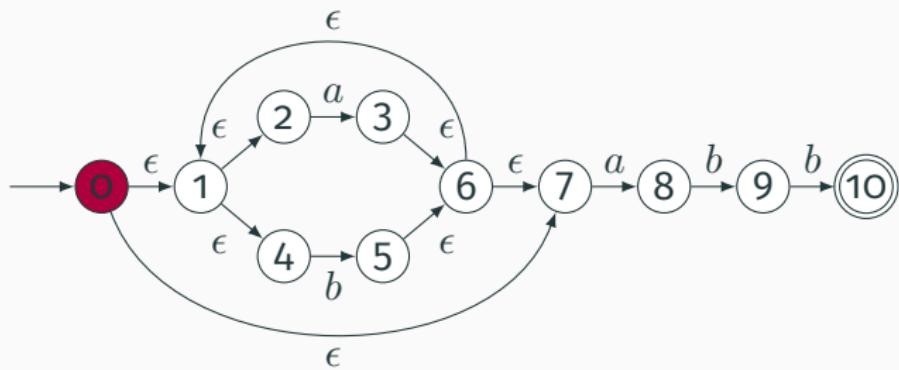
Problem: you must follow the “right” arcs to show that a string is accepted. How do you know which arc is right?

Solution: follow them all and sort it out later.

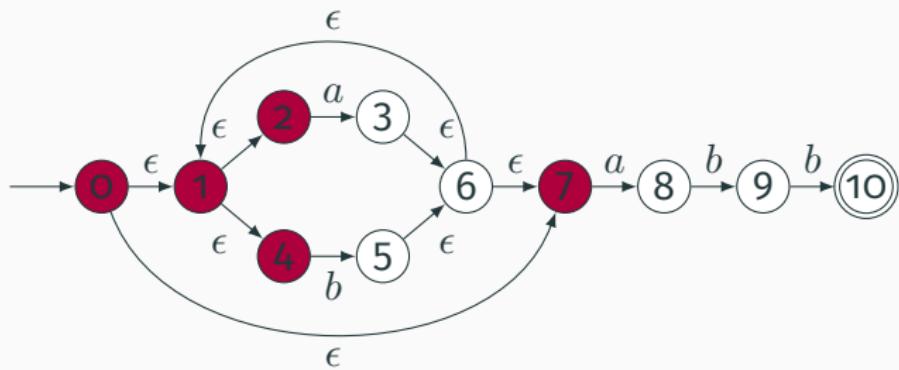
“Two-stack” NFA simulation algorithm:

1. Initial states: the  $\epsilon$ -closure of the start state
2. For each character  $c$ ,
  - New states: follow all transitions labeled  $c$
  - Form the  $\epsilon$ -closure of the current states
3. Accept if any final state is accepting

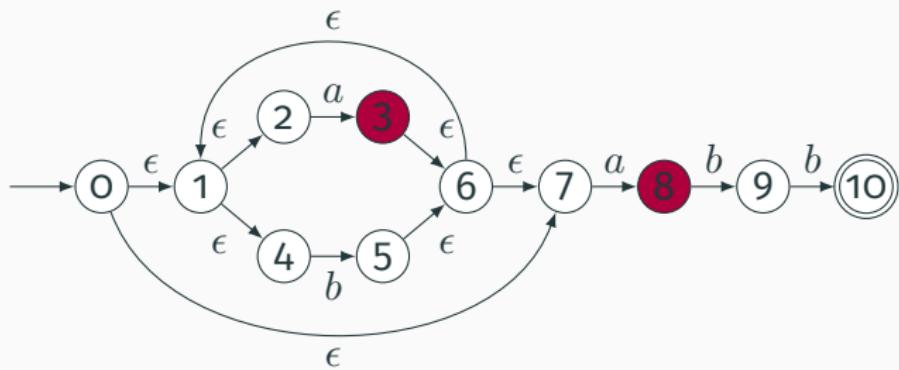
## Simulating an NFA: ·*aabb*, Start



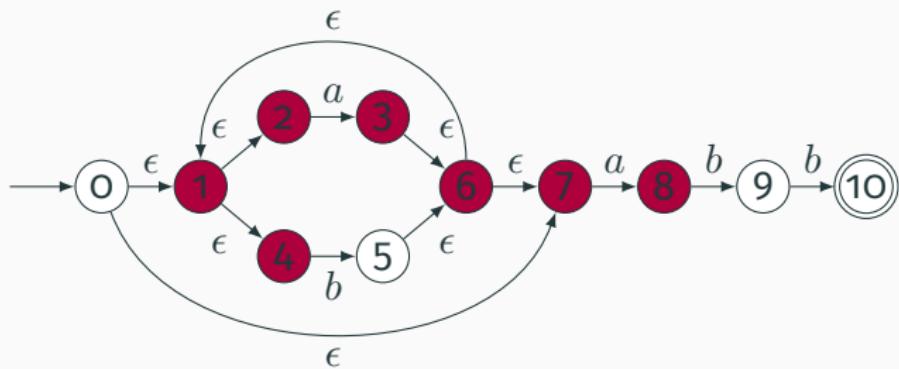
## Simulating an NFA: $\cdot aabb$ , $\epsilon$ -closure



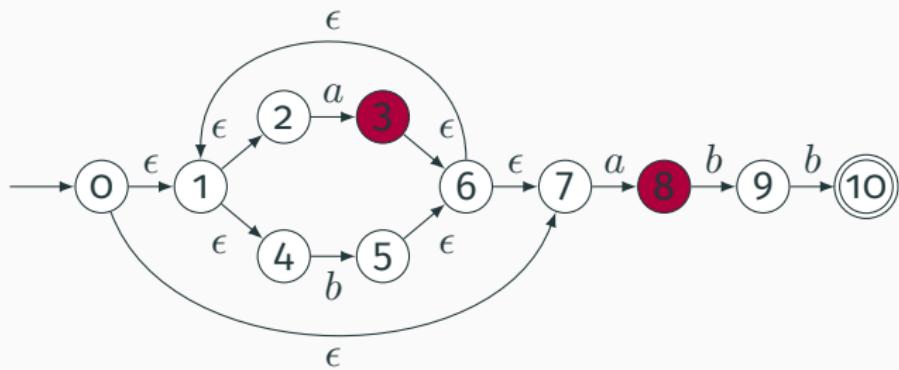
## Simulating an NFA: $a \cdot abb$



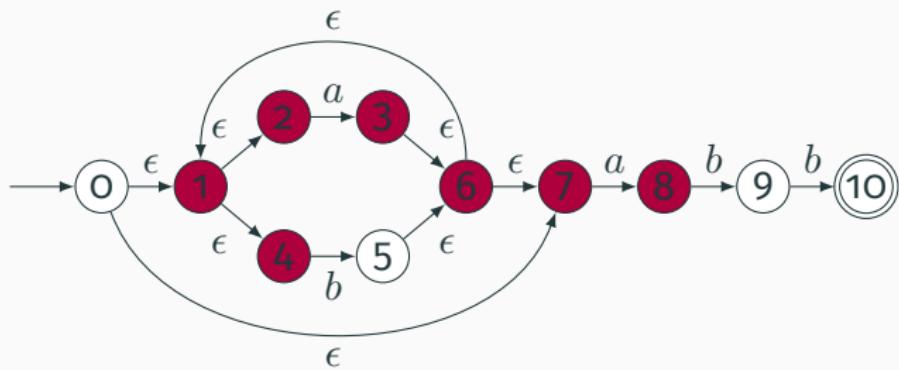
## Simulating an NFA: $a \cdot abb$ , $\epsilon$ -closure



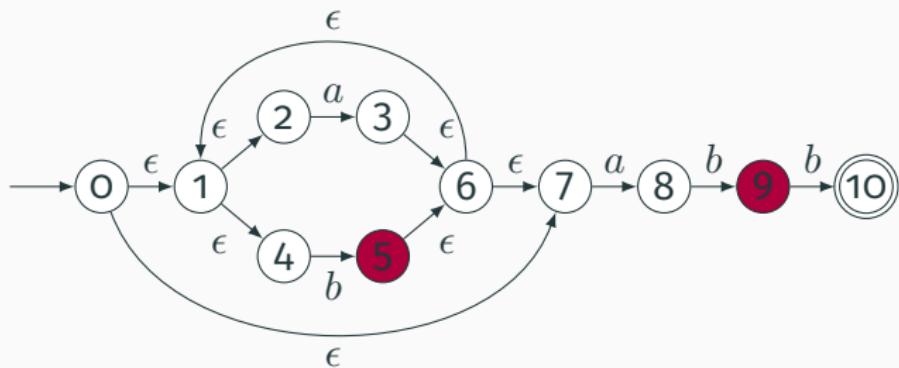
## Simulating an NFA: $aa \cdot bb$



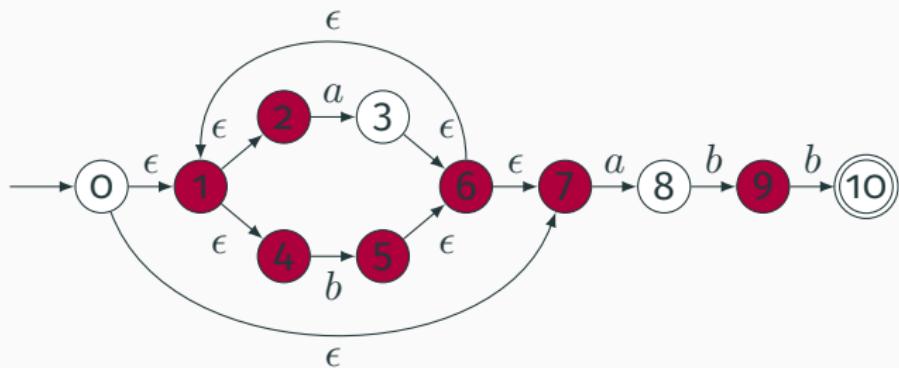
## Simulating an NFA: $aa \cdot bb$ , $\epsilon$ -closure



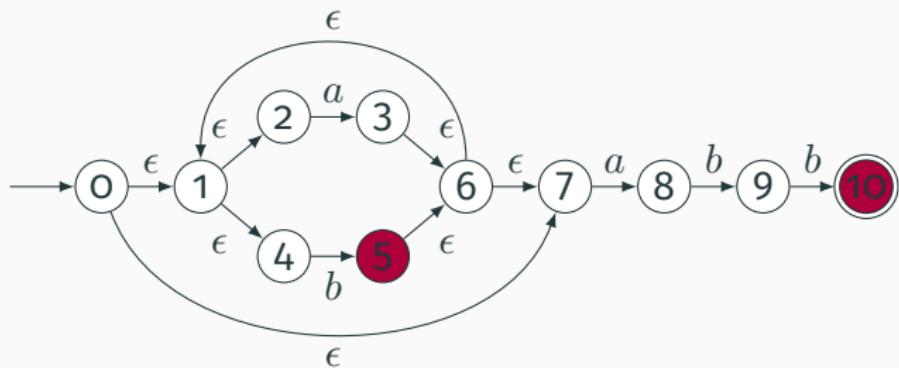
## Simulating an NFA: $aab \cdot b$



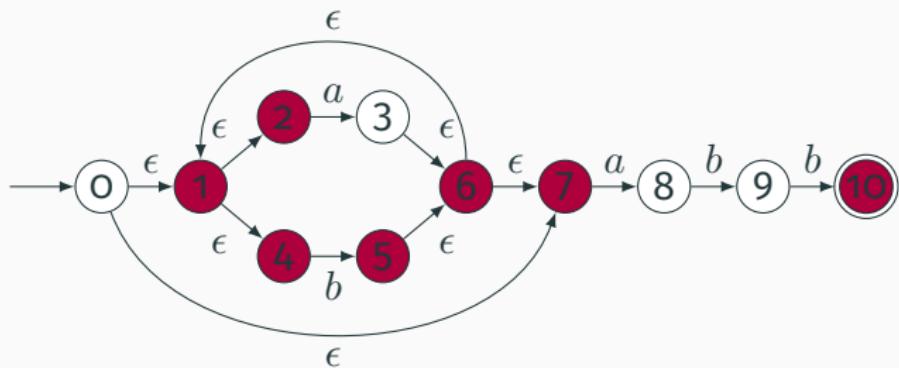
## Simulating an NFA: $aab \cdot b$ , $\epsilon$ -closure



## Simulating an NFA: $aabb$ .



## Simulating an NFA: $aabb\cdot$ , Done



# Deterministic Finite Automata

Restricted form of NFAs:

- No state has a transition on  $\epsilon$
- For each state  $s$  and symbol  $a$ , there is at most one edge labeled  $a$  leaving  $s$ .

Differs subtly from the definition used in COMS W3261 (Sipser,  
*Introduction to the Theory of Computation*)

Very easy to check acceptance: simulate by maintaining current state. Accept if you end up on an accepting state. Reject if you end on a non-accepting state or if there is no transition from the current state for the next symbol.

# Deterministic Finite Automata

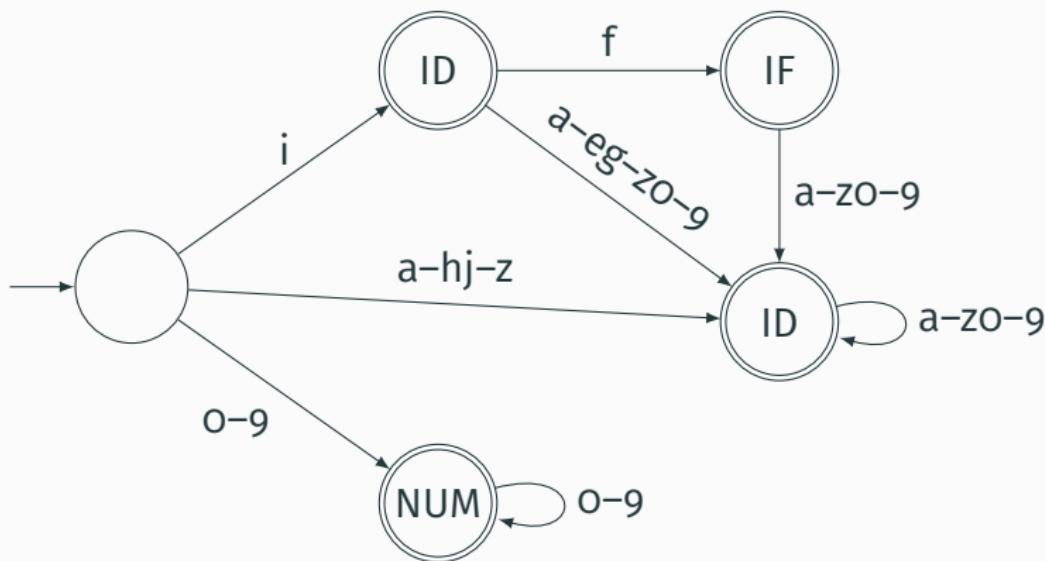
```
{  
    type token = ELSE | ELSEIF  
}  
  
rule token =  
    parse "else" { ELSE }  
    | "elseif" { ELSEIF }
```



# Deterministic Finite Automata

```
{ type token = IF | ID of string | NUM of string }
```

```
rule token =
  parse "if"
  | [ 'a' - 'z' ] [ 'a' - 'z' ] '0' - '9' ]* as lit { IF }
  | [ '0' - '9' ]+ as num { ID(lit) }
  | [ '0' - '9' ]+ as num { NUM(num) }
```



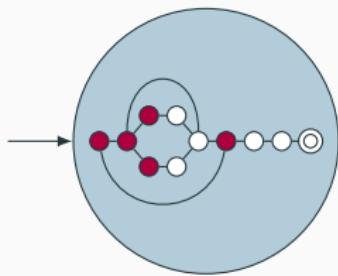
## Building a DFA from an NFA

Subset construction algorithm

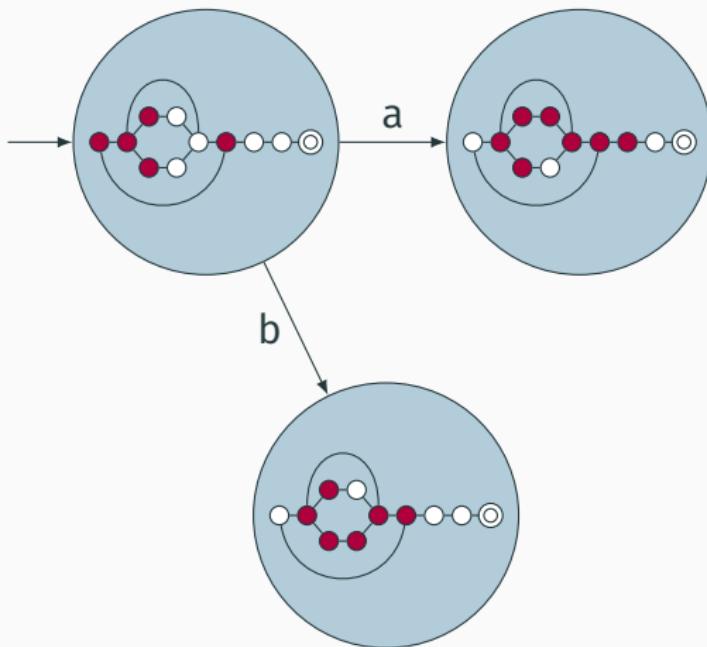
Simulate the NFA for all possible inputs and track the states that appear.

Each unique state during simulation becomes a state in the DFA.

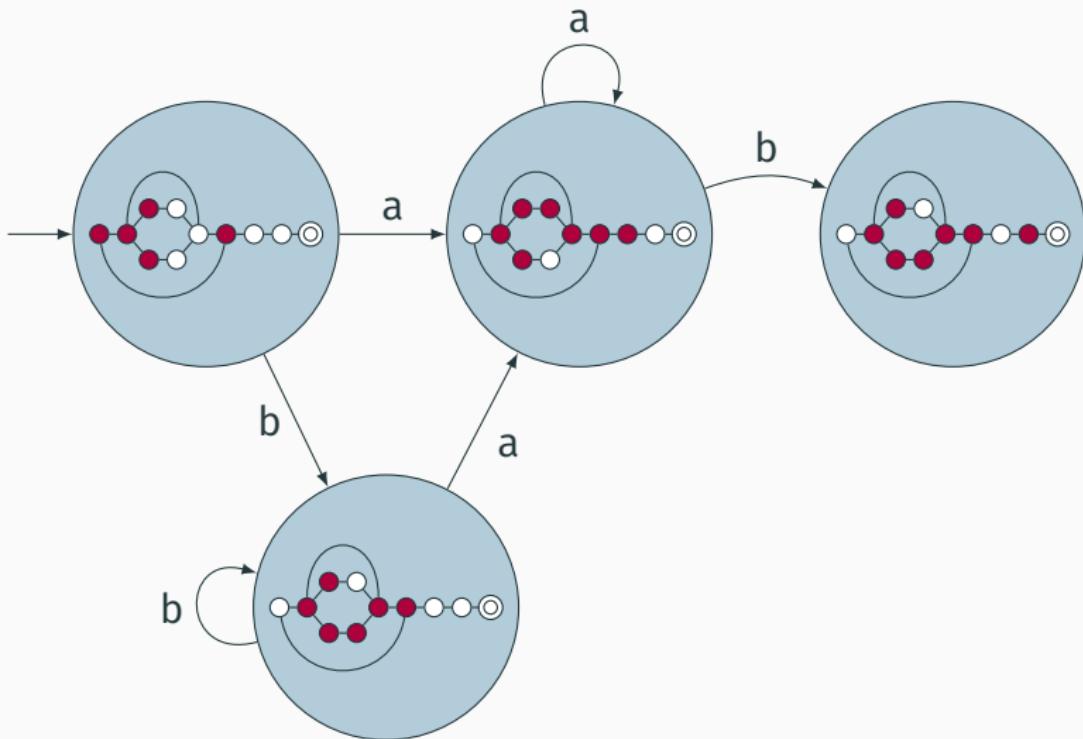
## Subset construction for $(a \mid b)^*abb$



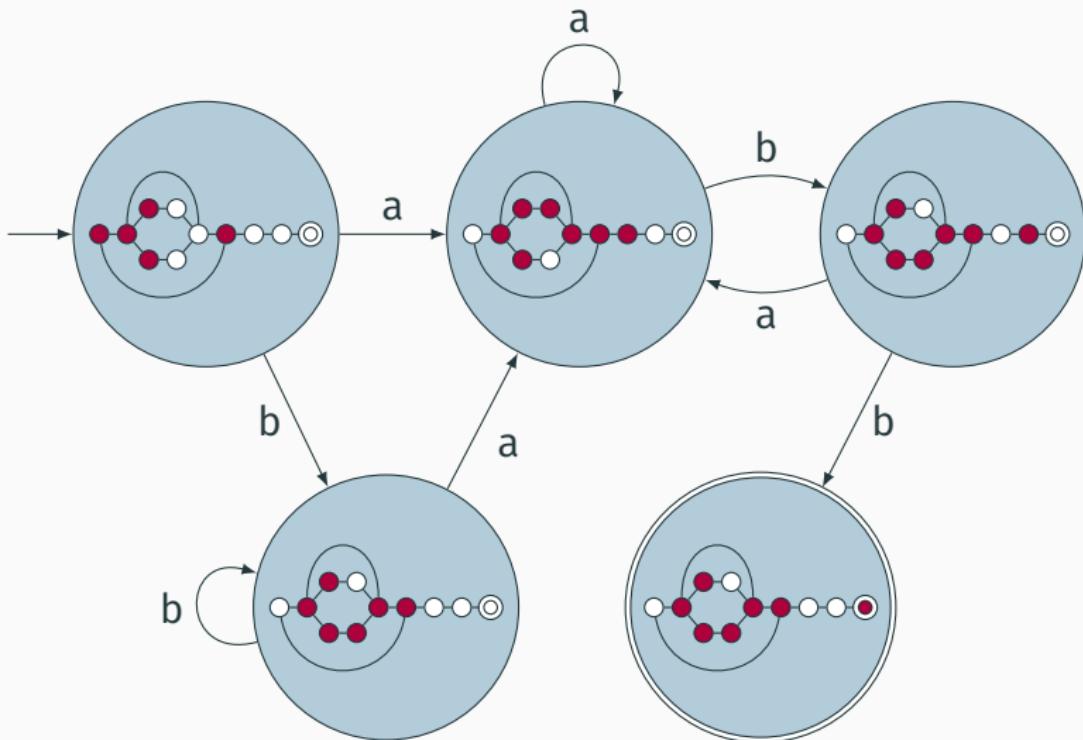
## Subset construction for $(a \mid b)^*abb$



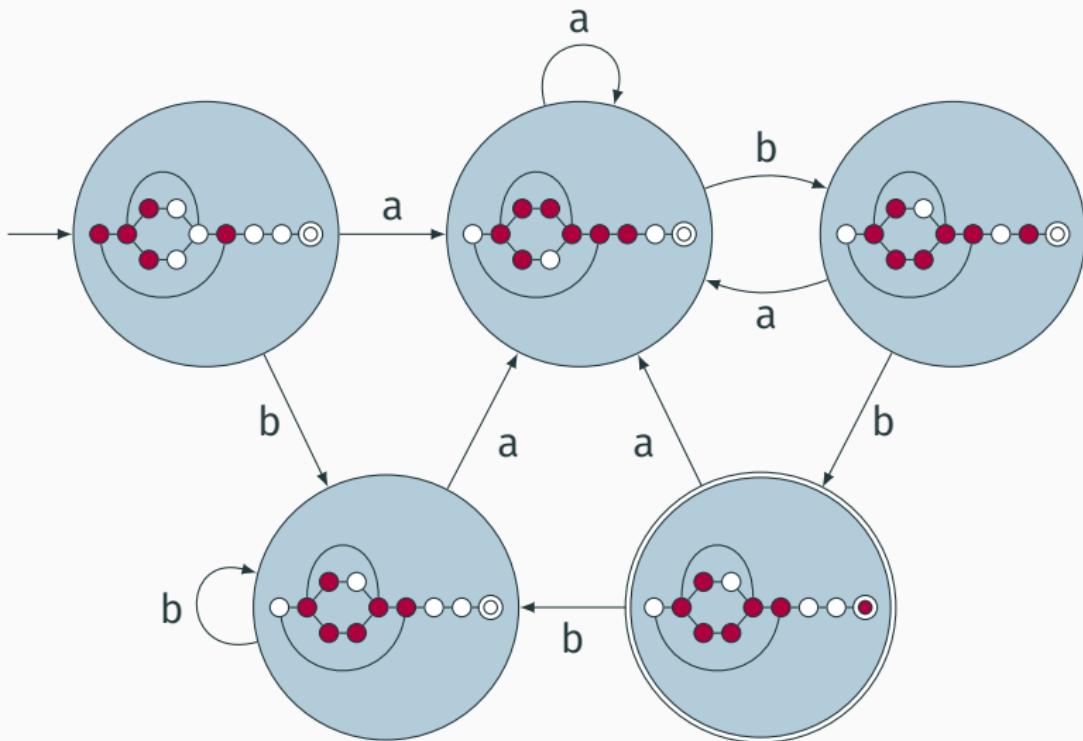
## Subset construction for $(a \mid b)^*abb$



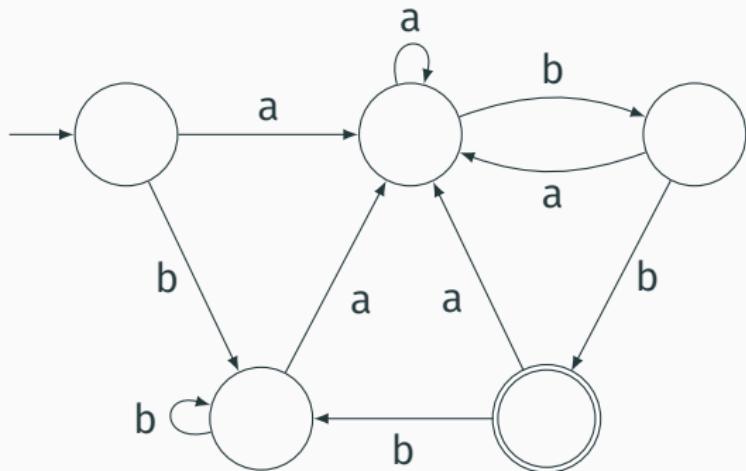
## Subset construction for $(a \mid b)^*abb$



## Subset construction for $(a \mid b)^*abb$

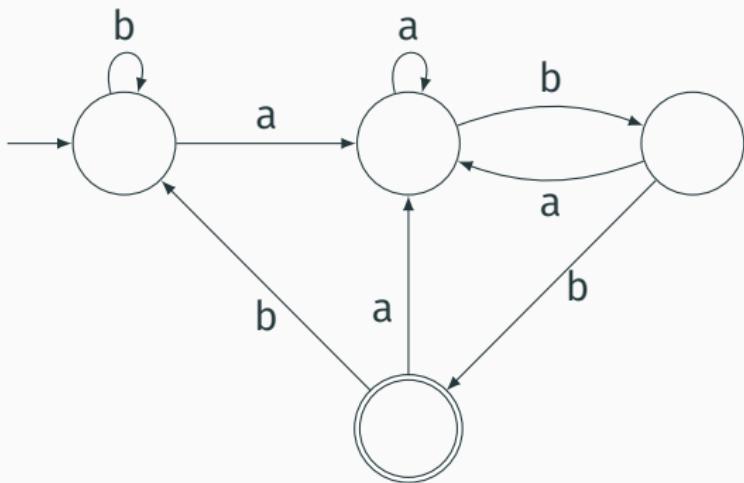


## Result of subset construction for $(a \mid b)^*abb$



*Is this minimal?*

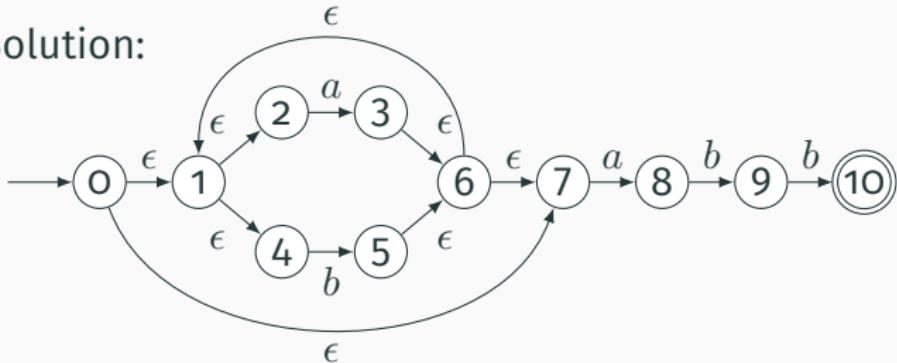
## Minimized result for $(a \mid b)^*abb$



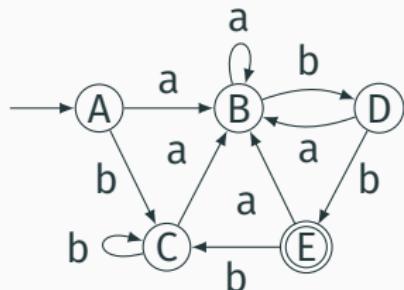
# Transition Table Used In the Dragon Book

Problem: Translate  $(a \mid b)^*abb$  into an NFA and perform subset construction to produce a DFA.

Solution:



NFA State	DFA State	a	b
{0,1,2,4,7}	A	B	C
{1,2,3,4,6,7,8}	B	B	D
{1,2,4,5,6,7}	C	B	C
{1,2,4,5,6,7,9}	D	B	E
{1,2,4,5,6,7,10}	E	B	C



## Subset Construction

An DFA can be exponentially larger than the corresponding NFA.

$n$  states versus  $2^n$

Tools often try to strike a balance between the two representations.

# Lexical Analysis with Ocamllex

---

# Constructing Scanners with Ocamllex



An example:

scanner.mll

```
{ open Parser }

rule token =
  parse [ ',' '\t' '\r' '\n' ] { token lexbuf }
  | '+,' { PLUS }
  | '-,' { MINUS }
  | '*,' { TIMES }
  | '/,' { DIVIDE }
  | ['0'-'9']+ as lit { LITERAL(int_of_string lit) }
  | eof { EOF }
```

# Ocamlex Specifications

```
{ (* Header: verbatim OCaml code; mandatory *)
}

(* Definitions: optional *)
let ident = regexp
let ...

(* Rules: mandatory *)
rule entrypoint1 [arg1 ... argn] =
  parse pattern1 { action (* OCaml code *) }
  | ...
  | patternn { action }
and entrypoint2 [arg1 ... argn] = ...
and ...

{
  (* Trailer: verbatim OCaml code; optional *)
}
```

# Patterns (In Order of Decreasing Precedence)

Pattern	Meaning
'c'	A single character
_	Any character (underline)
eof	The end-of-file
"foo"	A literal string
[ '1' '5' 'a'-'z' ]	"1," "5," or any lowercase letter
[^ '0'-'9' ]	Any character except a digit
( pattern )	Grouping
<i>identifier</i>	A pattern defined in the let section
<i>pattern</i> *	Zero or more <i>patterns</i>
<i>pattern</i> +	One or more <i>patterns</i>
<i>pattern</i> ?	Zero or one <i>patterns</i>
<i>pattern</i> <sub>1</sub> <i>pattern</i> <sub>2</sub>	<i>pattern</i> <sub>1</sub> followed by <i>pattern</i> <sub>2</sub>
<i>pattern</i> <sub>1</sub>   <i>pattern</i> <sub>2</sub>	Either <i>pattern</i> <sub>1</sub> or <i>pattern</i> <sub>2</sub>
<i>pattern</i> as <i>id</i>	Bind the matched pattern to variable <i>id</i>

# An Example

```
{ type token = PLUS | IF | ID of string | NUM of int }

let letter = [ 'a' - 'z' 'A' - 'Z' ]
let digit = [ '0' - '9' ]

rule token =
  parse [ ' ' '\n' '\t' ] { token lexbuf } (* Ignore whitespace *)
  | '+' { PLUS } (* A symbol *)
  | "if" { IF } (* A keyword *)
  | letter (letter | digit | '_')* as id { ID(id) } (* Identifiers *)
  | digit+ as lit { NUM(int_of_string lit) } (* Numeric literals *)
  | /*/* { comment lexbuf } (* C-style comments *)
  | _ { comment lexbuf } (* Ignore other characters *)
```

and comment =  
parse /\*/\* { token lexbuf } (\* Return to normal scanning \*)  
| \_ { comment lexbuf } (\* Ignore other characters \*)

# Free-Format Languages

Typical style arising from scanner/parser division

Program text is a series of tokens possibly separated by whitespace and comments, which are both ignored.

- keywords (if while)
- punctuation (, ( +)
- identifiers (foo bar)
- numbers (10 -3.14159e+32)
- strings ("A String")

# Free-Format Languages

Java      C      C++      C#      Algol      Pascal

Some deviate a little (e.g., C and C++ have a separate preprocessor)

But not all languages are free-format.

# Python

The Python scripting language groups with indentation

```
i = 0
while i < 10:
    i = i + 1
    print i      # Prints 1, 2, ..., 10
```

```
i = 0
while i < 10:
    i = i + 1
print i          # Just prints 10
```

This is succinct, but can be error-prone.

How do you wrap a conditional around instructions?

# Syntax and Language Design

Does syntax matter? Yes and no

More important is a language's *semantics*—its meaning.

The syntax is aesthetic, but can be a religious issue.

But aesthetics matter to people, and can be critical.

Verbosity does matter: smaller is usually better.

Too small can be problematic: APL is a succinct language with its own character set.

There are no APL programs, only puzzles.

# Syntax and Language Design

Some syntax is error-prone. Classic FORTRAN example:

```
DO 5 I = 1,25 ! Loop header (for i = 1 to 25)
DO 5 I = 1.25 ! Assignment to variable DO5I
```

Trying too hard to reuse existing syntax in C++:

```
vector< vector<int> > foo;
vector<vector<int>> foo; // Syntax error
```

C distinguishes `>` and `>>` as different operators.

Bjarne Stroustrup tells me they have finally fixed this.