

Generating Lexers

Lexer generators

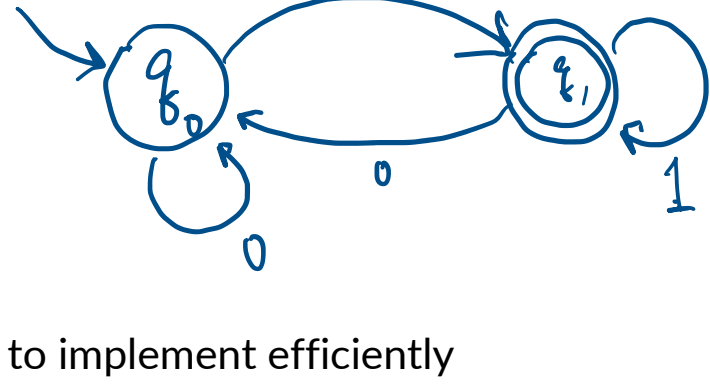
Spec = RE1 { Action1 }
 | RE2 { Action2 }
 | RE3 { Action3 }
 ... // in decreasing priority

Idea: Spec --> NFA --> DFA --> code = *loop + table*

Review: DFA

machine with finite # of states
 four-tuple (Q, Σ, F, δ)

$$\delta(q, x) = q'$$



DFA for recognizes odd bit numbers

Easy to implement efficiently

```
q = q0
while (i ≤ n) {
  q = δ(q, input[i])
  ++i
}
if (q ∈ F) return accept
else return fail
```

How to generate δ from regexes?

R1 | R2 | ... R1 IF
 R2 identifiers

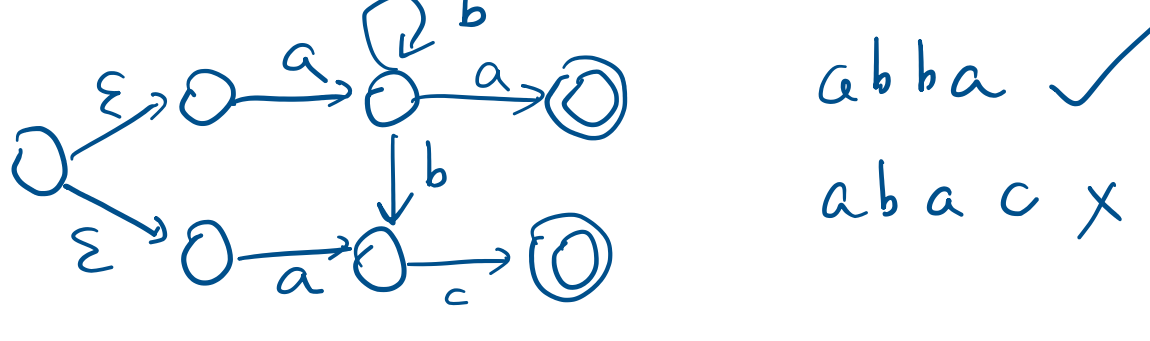
"i"

Idea: NFA, and then NFA --> DFA

Review: NFA

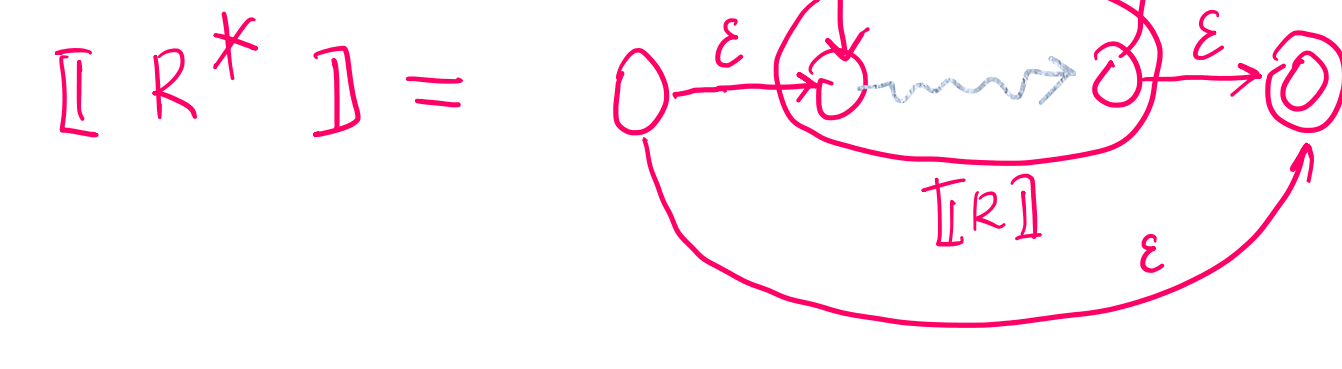
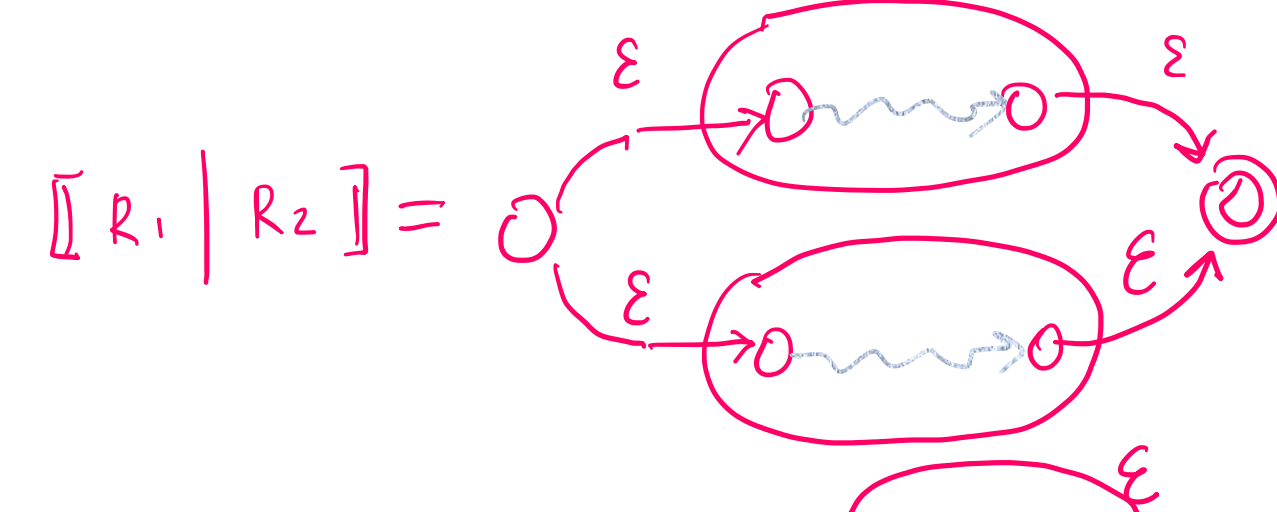
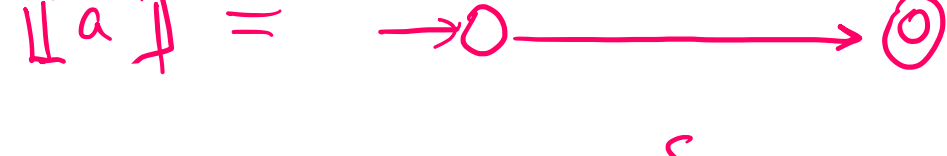
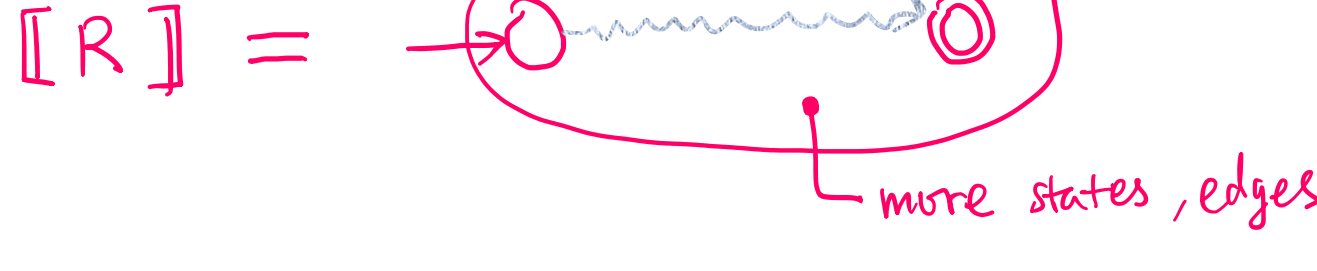
DFA + ...

- arrows can be ϵ
- exit arrows from a state can share same input symbol
- angelic nondeterminism

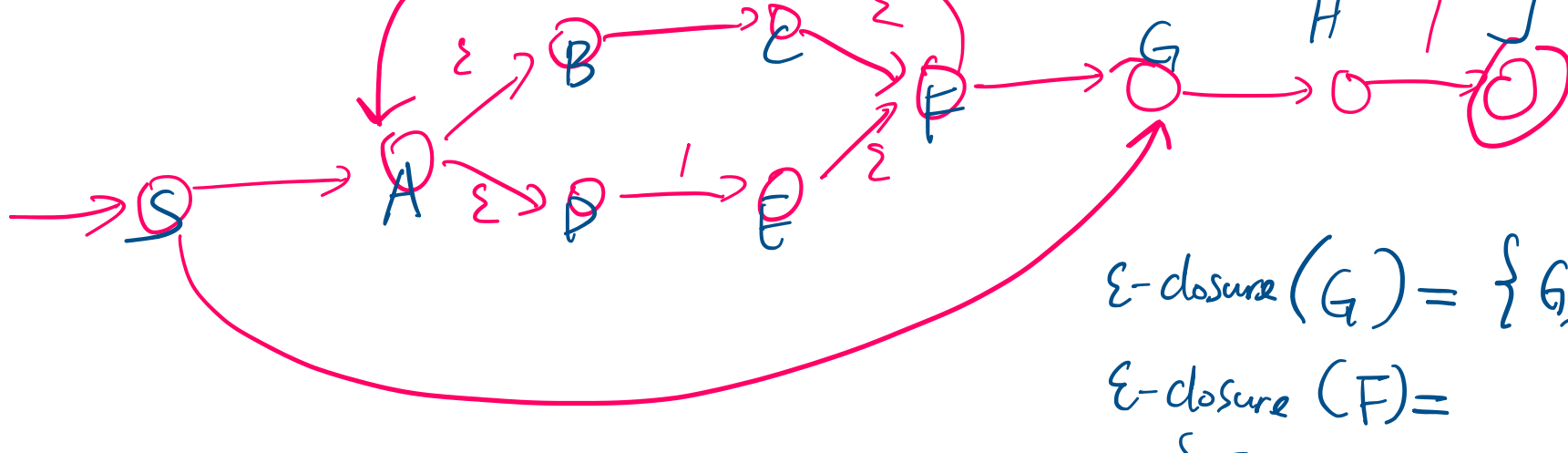


RE --> NFA

Defined recursively.



$[(0|1)^* 1]$



$\epsilon\text{-closure}(G) = \{G, H\}$
 $\epsilon\text{-closure}(F) = \{F, G, H, A, B, D\}$

NFA --> DFA

Idea: a DFA state = set of all NFA states that can possibly be reached by reading the same sequence of input symbols.

$\epsilon\text{-closure}(q)$ = Set of states reachable from q using zero or more ϵ edges
 Inductively defined.

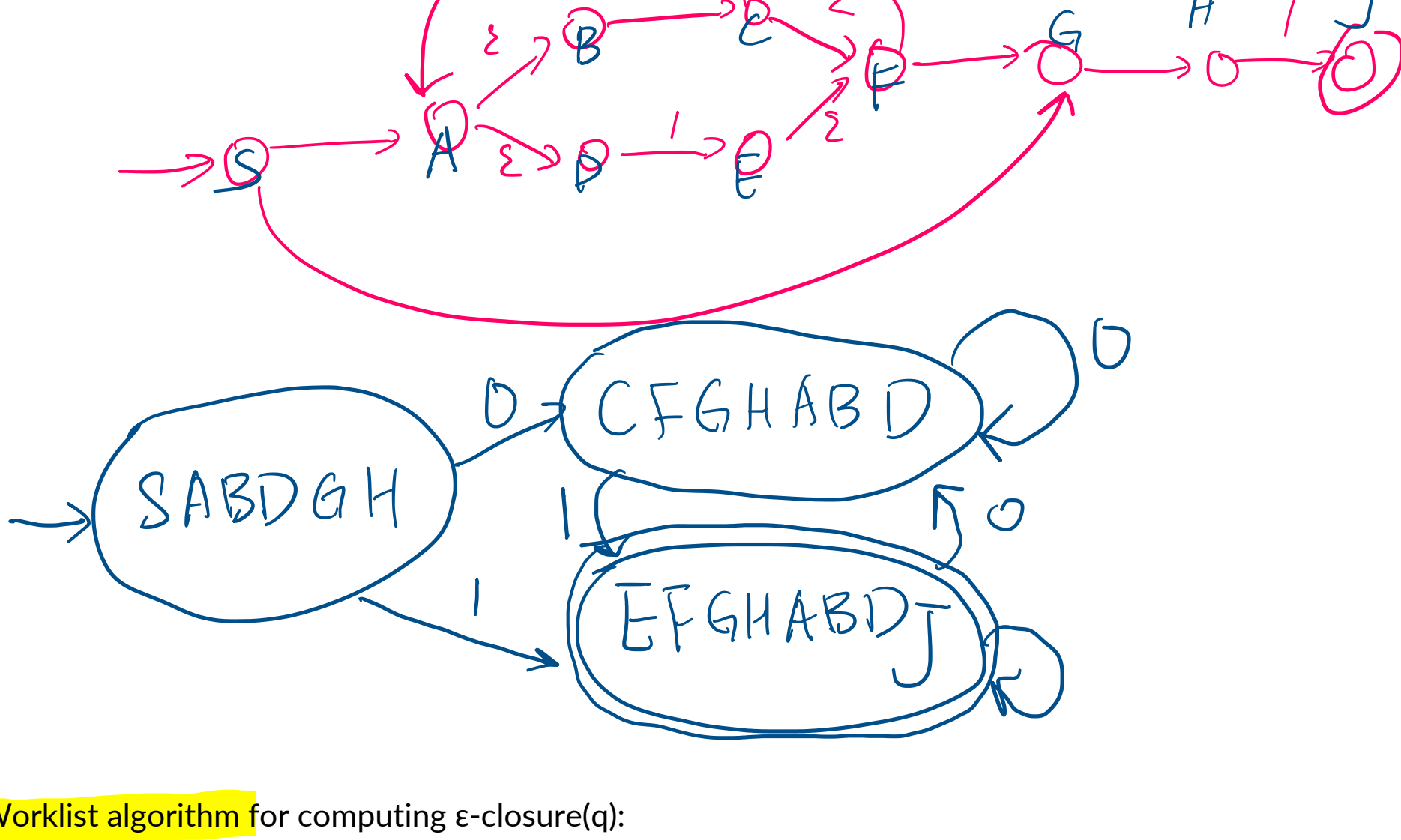
$$\frac{q \in \epsilon\text{-closure}(q)}{q \in \epsilon\text{-closure}(q)} \quad \frac{q'' \in \epsilon\text{-closure}(q) \quad q'' \xrightarrow{\epsilon} q'}{q' \in \epsilon\text{-closure}(q)}$$

Construct DFA from NFA:

Idea: a DFA state = set of all NFA states that can possibly be reached by reading the same sequence of input symbols

- Initial states: $Q_0 = \epsilon\text{-closure}(q_0)$
- States and transitions: $\delta(Q, x) = \epsilon\text{-closure}(\bigcup_{q \in Q} \{q' \mid q \xrightarrow{x} q'\})$
- Final state: Q is a final state of the DFA if any $q \in Q$ is a final NFA state

$[(0|1)^* 1]$



Worklist algorithm for computing $\epsilon\text{-closure}(q)$:

$$\epsilon\text{-closure}(q) [q'] \in B$$

$$\epsilon\text{-closure}(q) [q'] = \text{true} \quad \text{if } q = q'$$

$$\epsilon\text{-closure}(q) [q'] = \bigvee_{q'' \xrightarrow{\epsilon} q'} \epsilon\text{-closure}(q) [q'']$$

$$\epsilon\text{-closure}(q) = \{q, \mapsto \text{false}, \dots\}$$

worklist = { q } // empty on termination
 // Invariant: equations about q' not satisfied
 $\Rightarrow q' \in \text{worklist}$

```
while (worklist not empty) {
  remove q' from worklist
  update ε-closure(q) [q'] per equations
  if ε-closure(q) [q'] changed :
    for each q'' s.t. q' → q''
      add q'' to worklist
}
```

$O(N^2)$

Correctness?

N

Partial correctness + Termination

DFA minimization

Some DFA states can be equivalent after conversion.

$$\frac{Q_1 \in \text{Final} \quad Q_2 \notin \text{Final}}{Q_1 \neq Q_2} \quad \frac{Q_1 \notin \text{Final} \quad Q_2 \in \text{Final}}{Q_1 \neq Q_2} \quad \frac{\delta(Q_1, x) \neq \delta(Q_2, x)}{Q_1 \neq Q_2}$$

Algorithm: start with every pair of states being equivalent, and unequate them.

When to unequate two states?

Myhill-Nerode theorem: for any DFA, there is a unique minimal DFA that accepts the same input.

Putting it all together

Spec = R1 { Action1 }
 | R2 { Action2 }
 | R3 { Action3 }
 ... // in decreasing priority



How to implement the "Longest-matching token" rule?

NFA

R1 = abc

R2 = (abc)* d

Input = abc | abc | ... abc x

$O(n^2)$

$$n + n-1 + n-2 + \dots = O(n^2)$$