

# Formal Language Theory & Finite State Automata

COMP90042

Natural Language Processing

Lecture 13

Semester 1 2022 Week 7  
Jey Han Lau



THE UNIVERSITY OF  

---

MELBOURNE

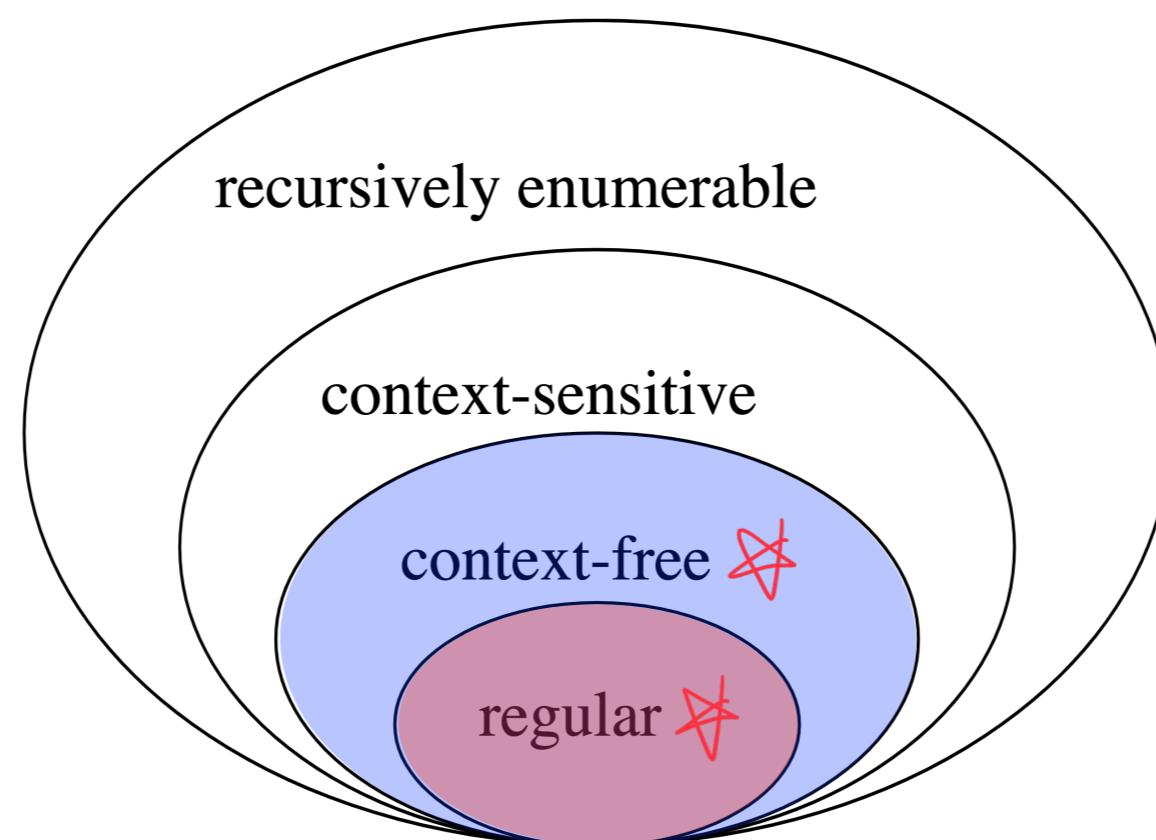
# What Have We Learnt?

- Methods to process sequence of words:
  - N-gram language Model
  - Hidden Markov Model
  - Recurrent Neural Networks
- Nothing is fundamentally linguistic about these models

# Formal Language Theory

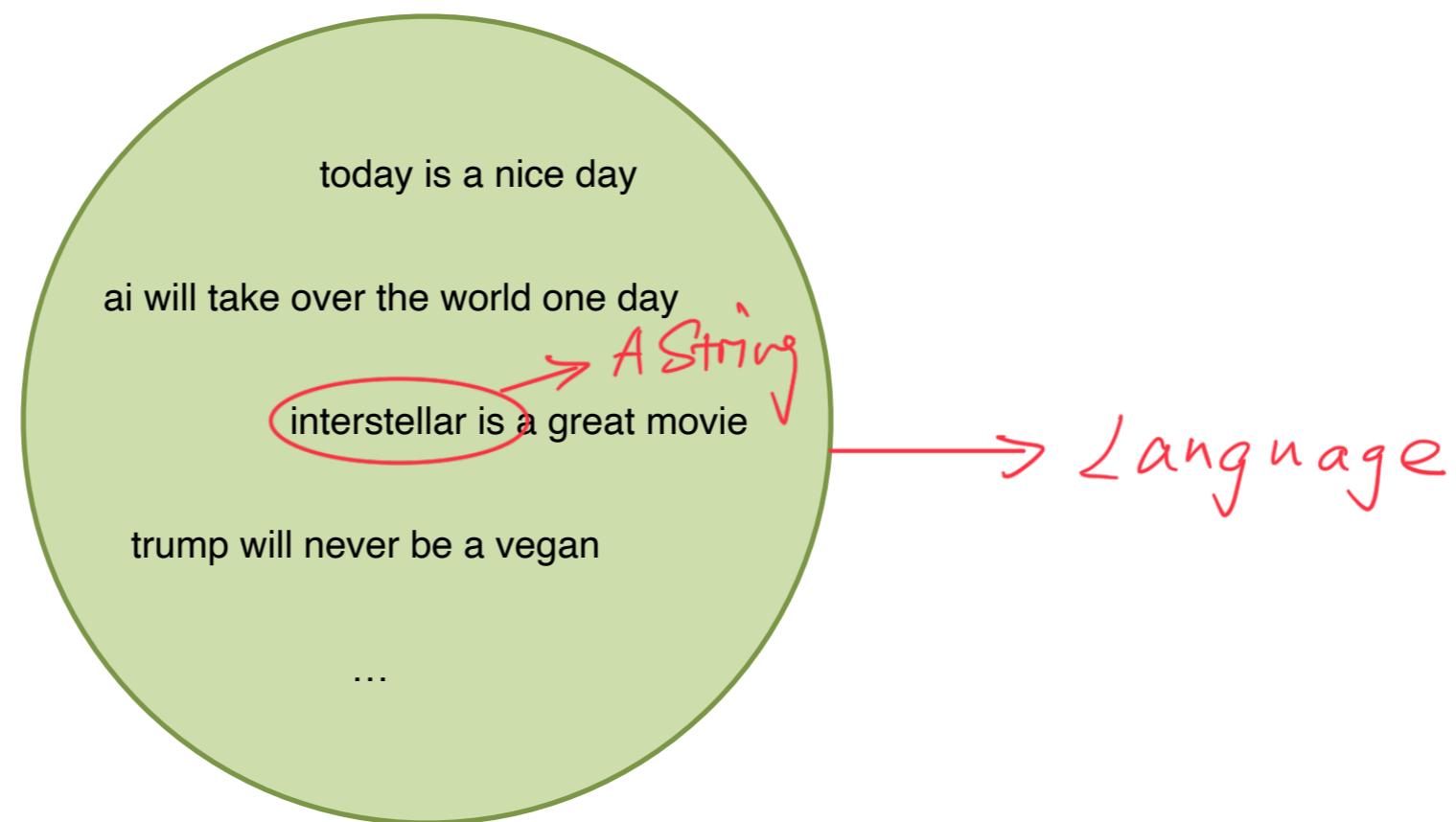
语言的种类

- Studies **classes of languages** and their computational properties
  - Regular language (this lecture)
  - Context free language (next lecture)



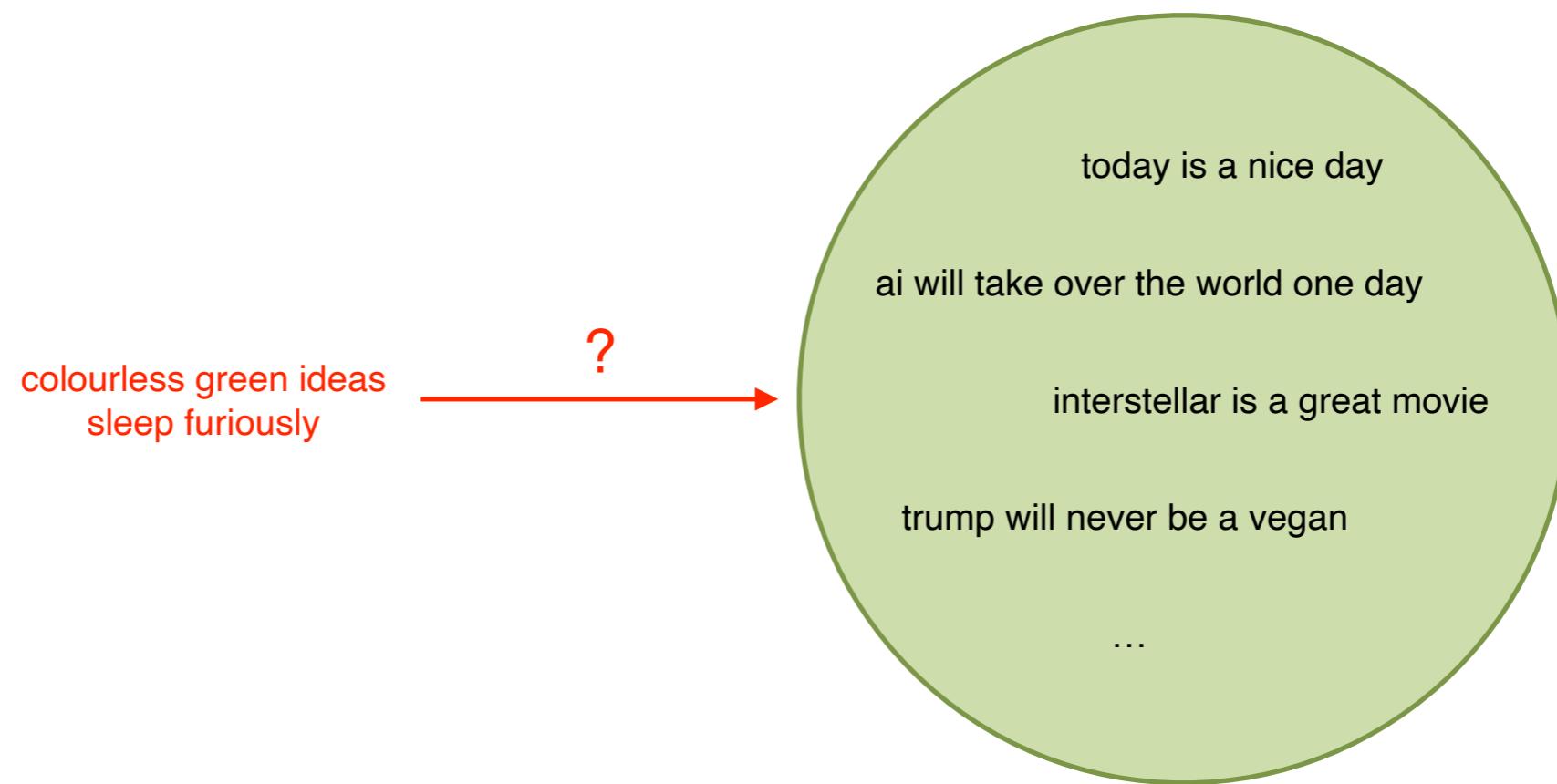
# Formal Language Theory

- A language = set of **strings** (*words*)
- A string = sequence of elements from a finite alphabet (AKA vocabulary)



# Why?

- Main goal is to solve the **membership problem**
  - Whether a string is in a language
- How? By defining its **grammar**



# Examples of Language

- Binary strings that start with 0 and end with 1 语言定义规则
  - ▶ { 01, 001, 011, 0001, ... } ✓
  - ▶ ( 1, 0, 00, 11, 100, ... ) ✗
- Even-length sequences from alphabet {a, b}
  - ▶ { aa, ab, ba, bb, aaaa, ... } ✓
  - ▶ { aaa, aba, bbb, ... } ✗
- English sentences that start with wh-word and end in ?
  - ▶ { *what* ?, *where my pants* ?, ... } ✓
  - ▶ { *hello how are you?*, *why is the dog so cute!* }

# Beyond Membership Problem...

- Membership 离散问题 (Binary)
  - Is the string part of the language? Y/N
- Scoring 连续问题 (What's the prob to be a language?)
  - Graded membership
  - “How acceptable is a string?” (language models!)
- Transduction
  - “Translate” one string into another (stemming!)

# Outline

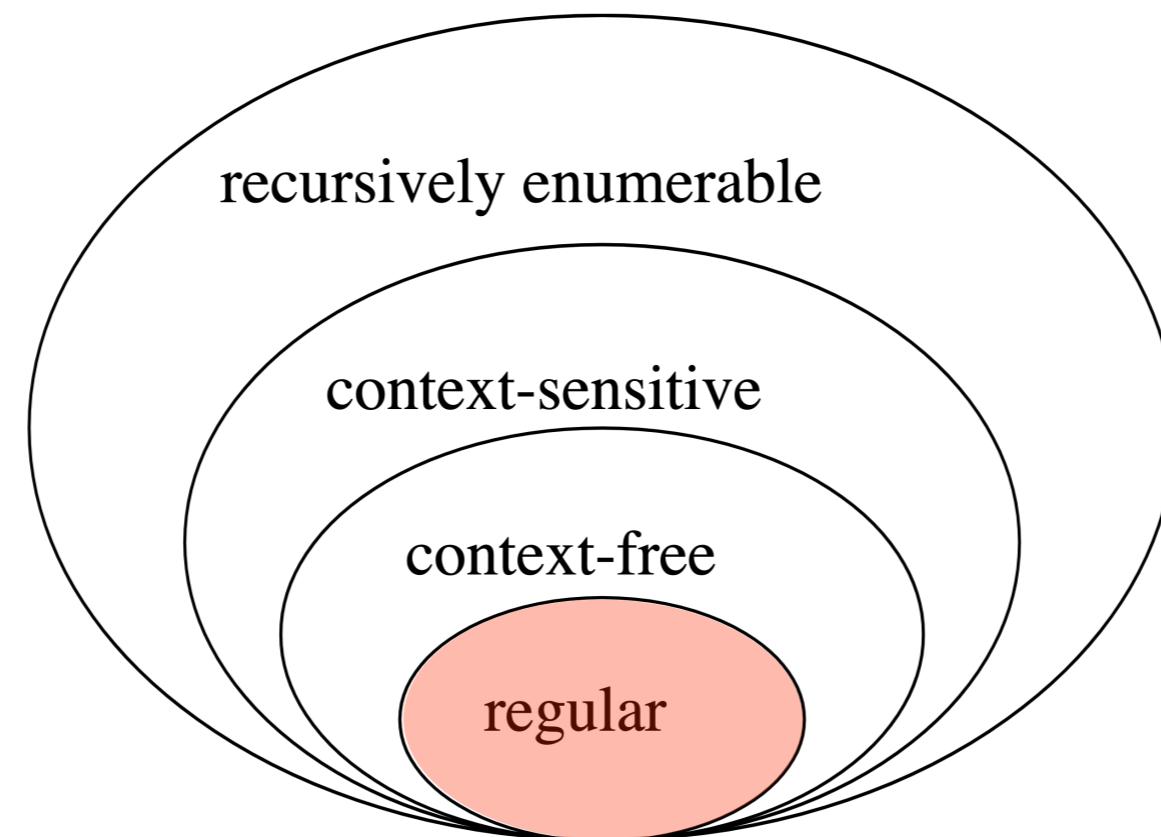
- Regular language
- Finite state acceptor
- Finite state transducer

# Regular Language

# Regular Language

只能用正则表达式表达的语言都是正则语言.

- The simplest class of languages
- Any regular expression is a regular language
  - ▶ Describes what strings are part of the language  
(e.g. ‘ $0(0|1)^*1$ ’)



# Regular Languages

- Formally, a regular expression includes the following operations/definitions:
    - Symbol drawn from **alphabet**,  $\Sigma$
    - **Empty string**,  $\epsilon$
    - **Concatenation** of two regular expressions,  $RS$
    - **Alternation** of two regular expressions,  $R|S$
    - Kleene star for **0 or more repeats**,  $R^*$
    - Parenthesis  $()$  to define **scope of operations**
-

# Examples of Regular Languages

- Binary strings that start with 0 and end with 1
  - $0(0|1)^*1$  0, 1 之間 0 or 1 不斷 repeat
- Even-length sequences from alphabet {a, b} 偶數長度
  - $((aa)|(ab)|(ba)|(bb))^*$  aa, ab, ba, bb  
每一段都可以是上列中的某一个，重複 0 次 or N 次
- English sentences that start with wh-word and end in ?
  - $((what)|(where)|(why)|(which)|(whose)|(whom)) \Sigma^* ?$

# Properties of Regular Languages

2021 F Exam → Answer in Workshop 8.

- **Closure:** if we take regular languages L<sub>1</sub> and L<sub>2</sub> and merge them, is the resulting language regular?
- RLs are closed under the following:
  - **concatenation and union**
  - **intersection:** strings that are valid in both L<sub>1</sub> and L<sub>2</sub>
  - **negation:** strings that are not in L
- Extremely versatile! Can have RLs for different properties of language, and use them together



# Finite State Acceptor

# Finite State Acceptor

- Regular expression **defines a regular language**
- But it doesn't give an algorithm to check whether a string belongs to the language
- **Finite state acceptors** (FSA) describes the computation involved for **membership checking**

State transfer ( condition ? )

# Finite State Acceptors

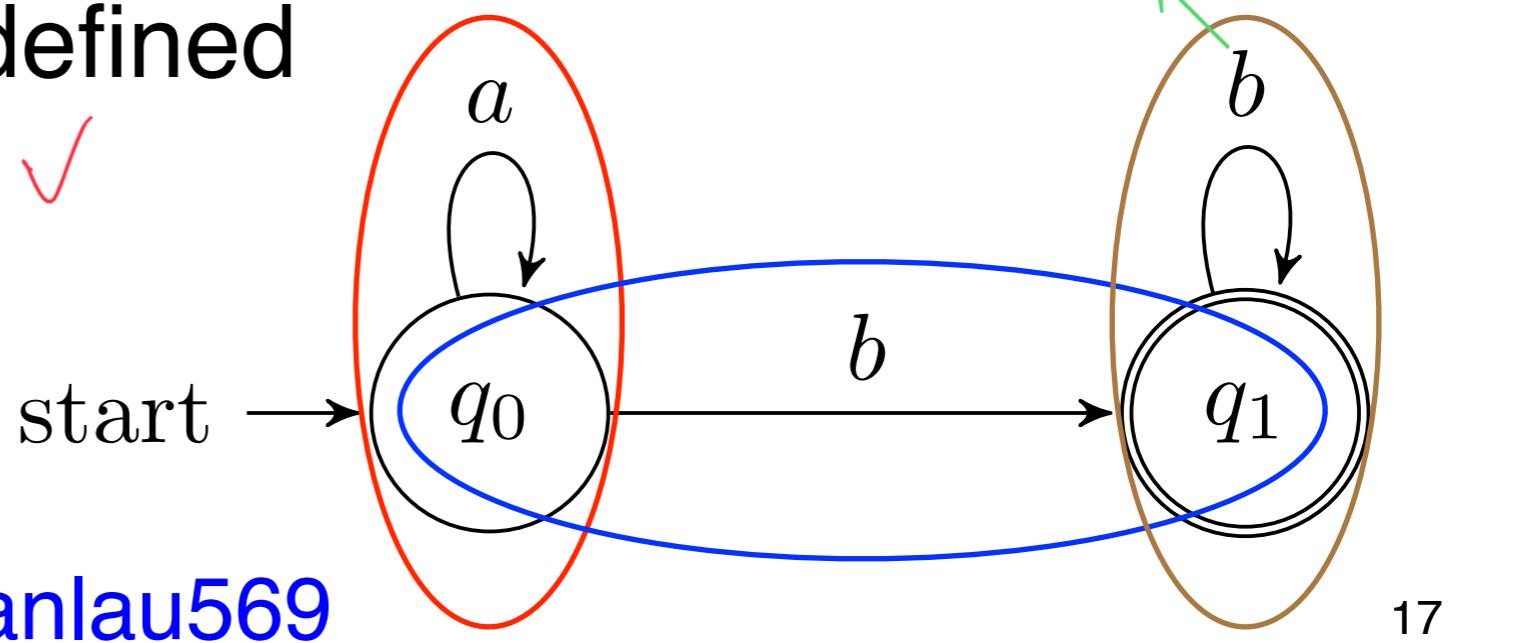
- FSA consists:
  - alphabet of input symbols,  $\Sigma$
  - set of states,  $Q$
  - start state,  $q_0 \in Q$
  - final states,  $F \subseteq Q$
  - transition function: symbol and state  $\rightarrow$  next state
- Accepts strings if there is path from  $q_0$  to a final state with transitions matching each symbol
  - Dijkstra's shortest-path algorithm,  $O(V \log V + E)$

# Example FSA

- Input alphabet  $\{a, b\}$
  - States  $\{q_0, q_1\}$
  - Start, final states  $q_0, \{q_1\}$
  - Transition function  $\{(q_0, a) \rightarrow q_0, (q_0, b) \rightarrow q_1, (q_1, b) \rightarrow q_1\}$
  - Regular expression defined by this FSA?
- $a^* b b^*$  ✓



[PollEv.com/jeyhanlau569](https://PollEv.com/jeyhanlau569)



# Derivational Morphology

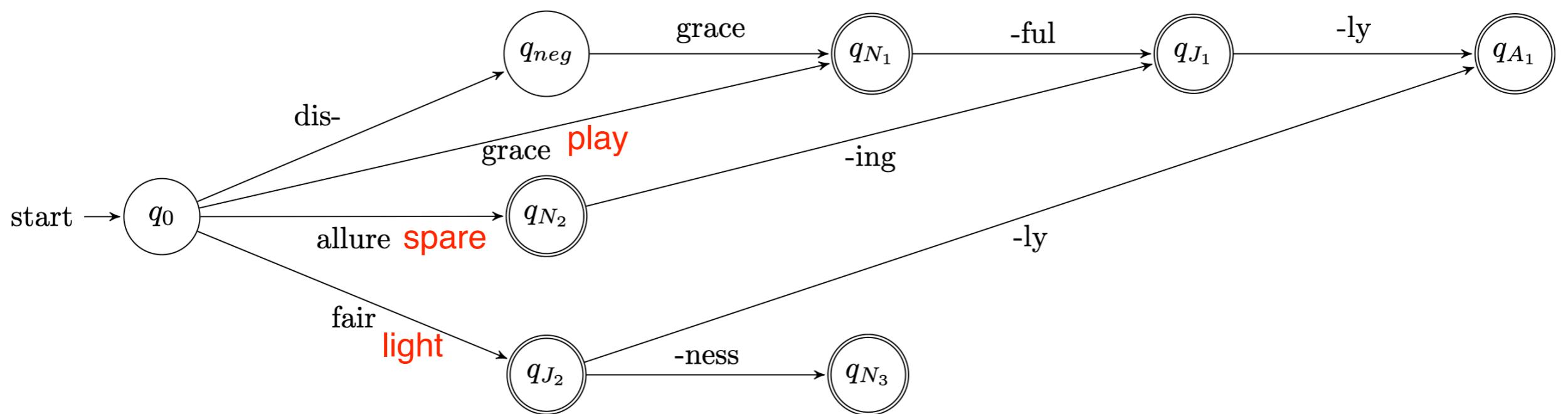
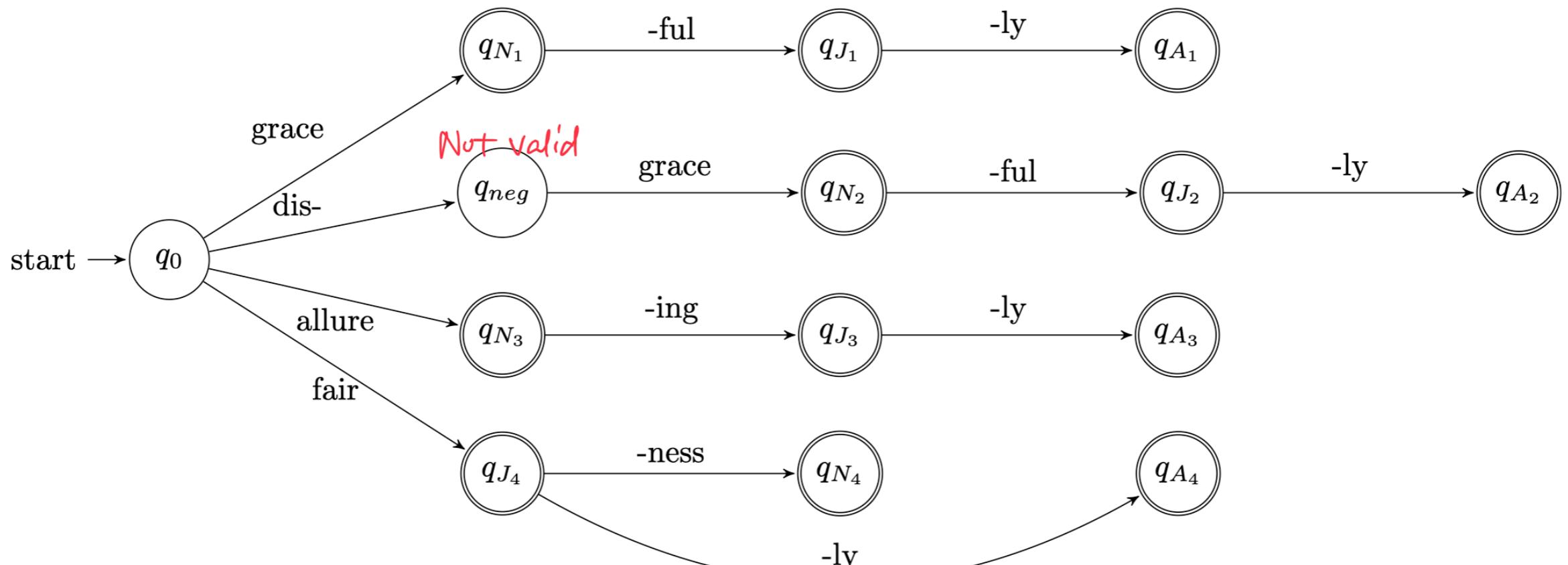
- Use of affixes to change word to another grammatical category
- *grace* → *graceful* → *gracefully*
- *grace* → *disgrace* → *disgracefully*
- *allure* → *alluring* → *alluringly*
- *allure* → *\*allureful*
- *allure* → *\*disallure*

} × ( 没有这些词 )      ✓

# FSA for Morphology

- Fairly consistent process
  - want to accept valid forms (*grace* → *graceful*)
  - reject invalid ones (*allure* → \**allureful*) 
  - generalise to other words, e.g., nouns that behave like *grace* or *allure*

# FSA for Word Morphology



# Weighted FSA

- Some words are **more plausible** than others
    - *fishful* vs. *disgracelyful*
    - *musicky* vs. *writey*
  - Graded measure of acceptability — weighted FSA changes the following:
    - start state weight function,  $\lambda: Q \rightarrow \mathbb{R}$
    - final state weight function,  $\rho: Q \rightarrow \mathbb{R}$
    - transition function,  $\delta: (Q, \Sigma, Q) \rightarrow \mathbb{R}$
- 

# WFSA Shortest-Path

- Total score of a path  $\pi = t_1, \dots, t_N$

$$\lambda(t_0) + \sum_{i=1}^N \delta(t_i) + \rho(t_N)$$

*Start*                            *transition*                            *Final*

- $t$  is an edge
- Use **shortest-path algorithm** to find  $\pi$  with minimum cost
  - $O(V \log V + E)$ , as before

# Finite State Transducer

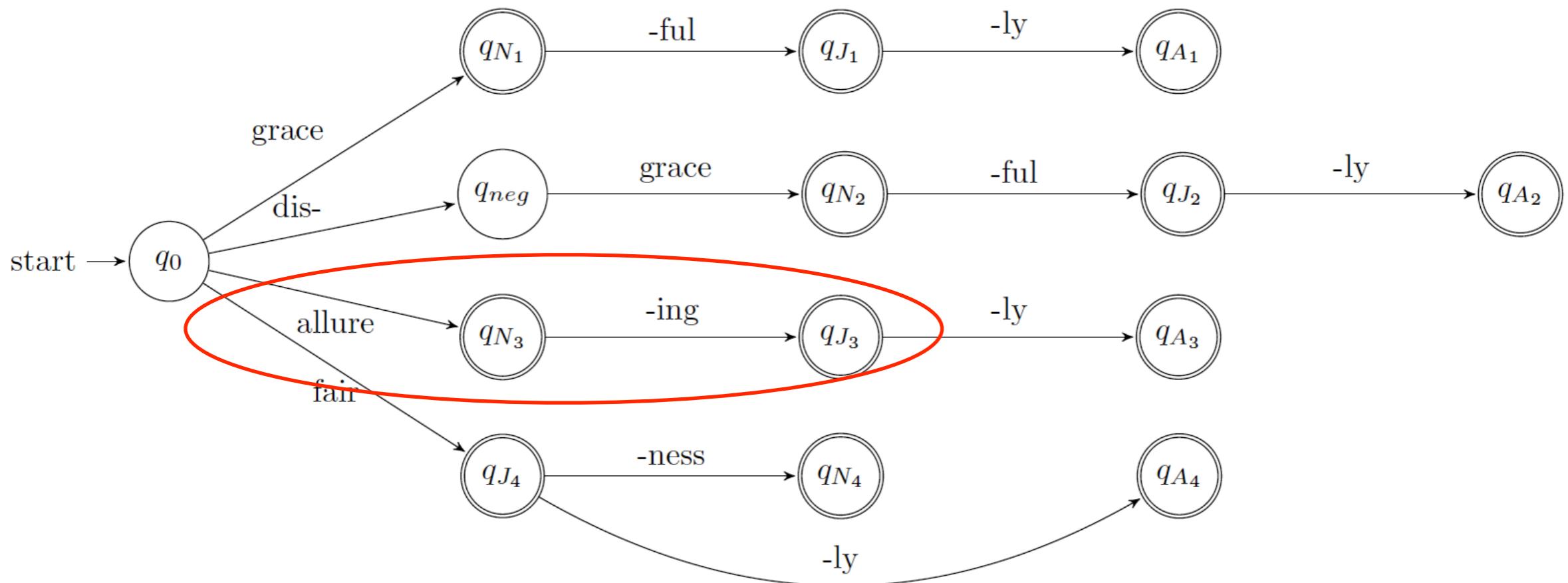
# Finite State Transducers (FST)

- Often don't want to just accept or score strings
  - want to **translate** them into another string

# Finite State Transducer

- FST add **string output** capability to FSA
  - includes an **output alphabet**
  - transitions now take input symbol and **emit output symbol** ( $Q, \Sigma, \Sigma, Q$ )
- Can be **weighted** = WFST
  - Graded scores for transition
- E.g., edit distance as WFST
  - distance to transform one string to another

# FSA for Word Morphology



Finite state **acceptor**: allure + ing = allureing

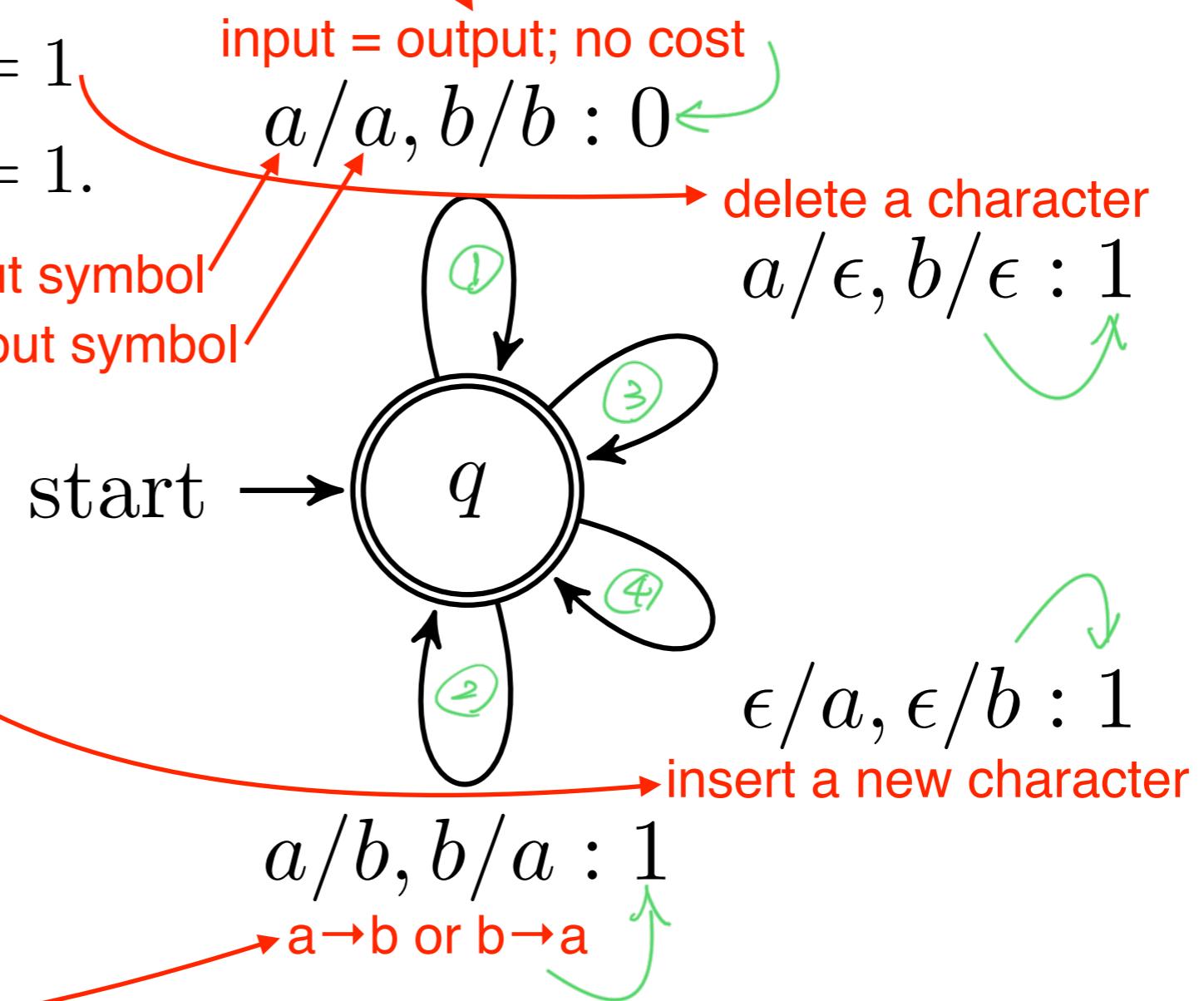
Finite state **transducer**: allure + ing = alluring

# Edit Distance Automata

transition weight

- ①  $\delta(q, a, a, q) = \delta(q, b, b, q) = 0$
- ②  $\delta(q, a, b, q) = \delta(q, b, a, q) = 1$
- ③  $\delta(q, a, \epsilon, q) = \delta(q, b, \epsilon, q) = 1$
- ④  $\delta(q, \epsilon, a, q) \stackrel{\text{empty string}}{=} \delta(q, \epsilon, b, q) = 1.$

1  
 $ab \rightarrow bb : 1$   
 $ab \rightarrow \underline{aaab} : 2$   
 insert 2 symbols

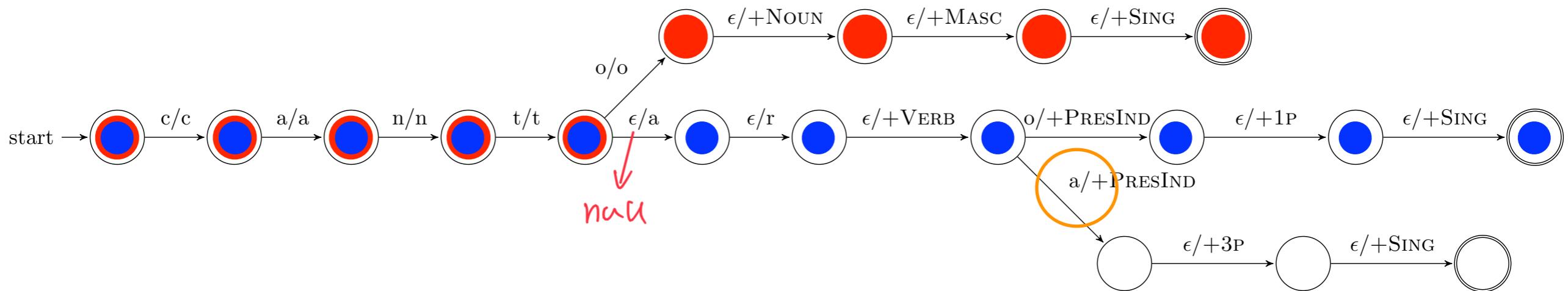


# FST for Inflectional Morphology

- Verb inflections in Spanish must match the subject in person & number
- Goal of morphological analysis:
  - canto* → *cantar+VERB+present+1P+singular*

	<b>cantar</b>	<b>to sing</b>
1P singular	yo canto	I sing
2P singular	tu cantas	you sing
3P singular	ella canta	she sings
1P plural	nostotros cantamos	we sing
2P plural	vosotros cantáis	you sing
3P plural	ellas cantan	they sing

# FST for Spanish Inflection



*canto* →

- *canto+Noun+Masc+Sing*
- *cantar+Verb+PresInd+1P+Sing*

*canta* → *cantar+VERB+PresInd+3P+Sing*

# Is natural language regular?

- Yes
- ~~No~~
- Maybe
- Don't know

[PollEv.com/jeyhanlau569](https://PollEv.com/jeyhanlau569)



# Sometimes...

- Example:

*the mouse that ran.*

*the cat that killed the mouse that ran.*

...

*the lion that bullied the hyena that bit the dog that chased  
the cat that killed the mouse that ran*

...

- Length is unbounded, but structure is local

- ▶ (Det Noun Prep Verb)\*

- ▶ Can describe with FSA

# Non-Regular Languages

- Arithmetic expressions with balanced parentheses
  - $(a + (b \times (c / d)))$  无法被 RL & FSA/FST 描述. : 括号太多增长无法规定
  - Can have arbitrarily many opening parentheses
  - Need to remember how many open parentheses, to produce the same number of closed parentheses
  - Can't be done with finite number of states
- $a^n b^n$

# Center Embedding

- Center embedding of relative clauses
  - The cat loves Mozart
  - The cat the dog chased loves Mozart
  - The cat the dog the rat bit chased loves Mozart
  - The cat the dog the rat the elephant admired bit chased loves Mozart
- Need to remember the *n* subject nouns, to ensure *n* verbs follow (and that they agree etc)
- Requires (at least) **context-free grammar** (next lecture!)

# Summary

- Concept of a language
- Regular languages
- Finite state automata: acceptors, transducers
- Weighted variants
- Application to edit distance, morphology

# Reading

- E18, Chapter 9.1