#### Regular Languages

Natural Language Processing CS 6120—Spring 2014 Northeastern University

David Smith with material from Jason Eisner, Andrew McCallum, and Lari Karttunen

# Brief History: 1950s

- Early NLP on machines less powerful than pocket calculators
  - E.g., how to compress a word list into memory
- Foundational work on automata, formal languages, information theory
- First speech systems (Bell Labs)
- Machine translation heavily funded by military, basically just word substitution
- Little formalization of syntax, semantics, pragmatics

## Brief History: 1960s

- ALPAC report (Alvey, 1966) ends funding for MT in U.S.
  - Lack of practical results, recommends basic research
- ELIZA and other early AI dialogue systems
  - Risibly easy Turing tests
- Early corpora: Brown Corpus (Kučera & Francis)

# Brief History: 1970s

- Winograd's SHRDLU (1971): existence proof of NLP (in tangled Lisp code)
  - Interpreted language about "blocks world"
    - Which cube is sitting on the table?
    - The large green one which supports the red pyramid.
    - Is there a large block behind the pyramid?
    - Yes, three of them. A large red one, a large green cube, and the blue one.
    - Put a small one onto the green cube which supports a pyramid.
    - OK.
- Hidden Markov models for speech recognition

### Brief History: 1980s

- Procedural  $\rightarrow$  declarative
  - Grammars, logic programming
  - Separation of processing (parser) from description of linguistic knowledge
- Representations of meaning: procedural semantics (SHRDLU), semantic nets (Schank), logic (starting in 1970s, Montague, Partee)
- Knowledge representation (Lenat: Cyc, still going!)
- MT in limited domains (METEO)
- HMMs for part-of-speech tagging (independently, Church & DeRose)

## Brief History: 1990s

- Probabilistic paradigm shift
  - Speech recognition methods take over the world
- IR-style evaluations take over the world
- Finite-state methods in speech and beyond
- Large amounts of monolingual and multilingual text become available, esp. on WWW
- Classification problems and *ambiguity resolution* (in syntax, lexical semantics, translation, etc.)

#### Brief History: Now

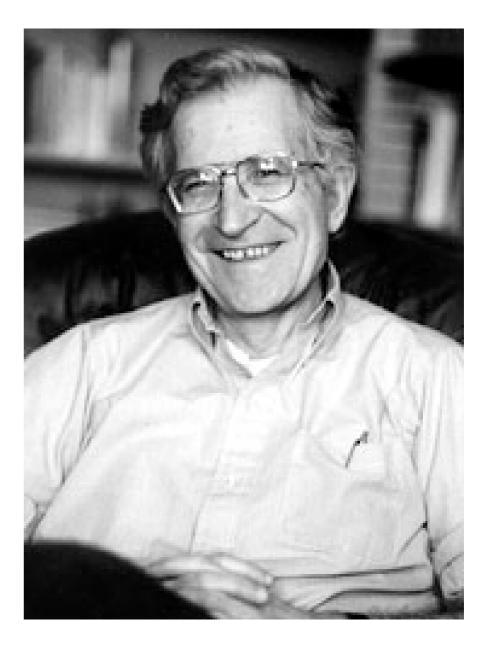
- Even more machine learning
  - Successful unsupervised systems
- Even more data, and tasks
- Widely usable—and used—speech recognition and machine translation
- Widely usable syntactic parsing
- Some usable dialog systems
- Convergence with IR, question answering, probabilistic knowledge representation

#### Noam Chomsky 1928–

Formal languages (Chomsky hierarchy)

Generative grammar

Anarcho-Socialist



### A Language

- Some sentences in the language
  - The man took the book.
  - Colorless green ideas sleep furiously.
  - This sentence is false.
- Some sentences not in the language
  - \*The girl, the sidewalk, the chalk, drew.
  - \*Backwards is sentence this.
  - \* \*Je parle anglais.

#### Languages as Rewriting Systems

- Start with some "non-terminal" symbol **S**
- Expand that symbol, using a rewrite rule.
- Keep applying rules until all non-terminals are expanded to terminals.
- The string of terminals is a sentence of the language.

## Chomsky Hierarchy

- Let Caps = nonterminals; lower = terminals; Greek = strings of terms/nonterms
- Recursively enumerable (Turing equivalent)
  - \* Rules:  $\alpha \rightarrow \beta$
- Context-sensitive
  - \* Rules:  $\alpha A \beta \rightarrow \alpha \gamma \beta$
- Context-free
  - ✤ Rules: A→α
- Regular (finite-state)
  - \* Rules:  $A \rightarrow aB$ ;  $A \rightarrow a$

# Regular Language Example

- Nonterminals: S, X
- Terminals: m, o
- Rules:
  - S→mX
  - X→oX
  - X→o

One expansion

S mX moX mooX mooo

• Start symbol: S

# Another Regular Language

- Strings in and not in this language
  - In the language:
    - "ba!", "baa!", "baaaaaaaa!"
  - Not in the language:
    - "ba", "b!", "ab!", "bbaaa!", "alibaba!"

а

s3

s4

- Regular expression: baa\*!
- Finite state automaton

s2

b

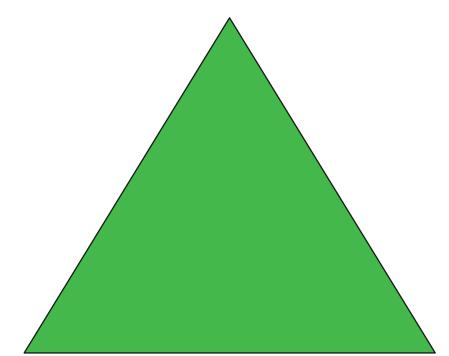
**s1** 

double circle indicates "accept state"

#### Regular Languages

#### **Regular Languages**

the accepted strings



#### **Finite-state Automata**

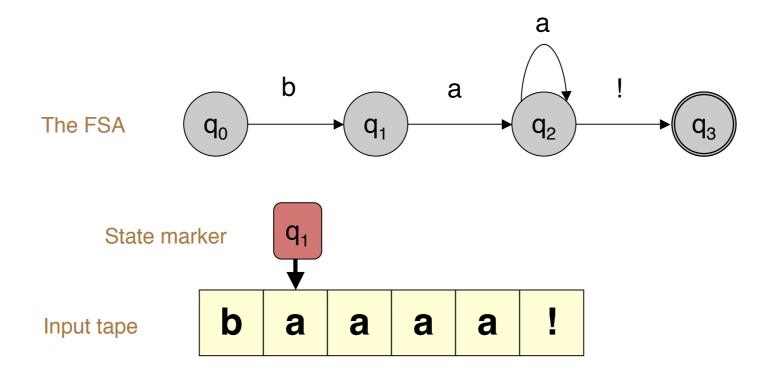
machinery for accepting

**Regular Expressions** 

a way to type the automata

#### Finite-State Automata

- A (deterministic) finite-state automaton is a 5-tuple (Q,  $\Sigma$ , q<sub>0</sub>, F,  $\delta$ (q,i))
  - Q: finite set of states  $q_0, q_1, q_2, ..., q_N$
  - \*  $\Sigma$ : finite set of terminals
  - \*  $\delta(q,i)$ : transition function (relation if non-deterministic)
  - ✤ q<sub>0</sub>: start state
  - ✤ F: set of final states



#### **Transition Table**

	Input			
State	b	а	!	
0 1 2 3	1 Ø Ø	Ø 2 2 Ø	Ø Ø 3 Ø	

### Regular Expressions

- Two types of characters
- Literal
  - Every "normal" alphanumeric character is an RE, and matches itself
- Meta-characters
  - Special characters that allow you to combine REs in various ways
- Example:
  - \* a matches a
  - a\* matches ε or a or aa or aaa or ...

### Regular Expressions

	Pattern	Matches
Concatenation	abc	abc
Disjunction	a b (a bb)d	a b ad bbd
Kleene star	a* c(a bb)*	ε a aa aaa ca cbba

Regular expressions / FSAs are closed under these operations

## Practical Applications

- Word processing find & replace
- Validate fields in database (dates, email, ...)
- Searching for linguistic patterns
- Finite-state machines
  - Language modeling in speech recognition (where things need to be real-time or better)
  - Information extraction
  - Morphology

	Pattern	Matches
Character Concat	went	went
Alternatives	(go went) [aeiou]	go went a o u
disjunc. negation wildcard char	[^aeiou]	bcdfg az&
Loops & skips	a*	ε a aa aaa
one or more	a+	a aa aaa
zero or one	colou?r	color colour

- Special characters
  - $\ t$ tab $\ v$ vertical tab $\ n$ newline $\ r$ carriage return
- Aliases (shorthand)

— \d	digits	[0-9]
— \D	non-digits	[^0-9]
— \w	alphabetic	[a-zA-Z]
— \W	non-alphabetic	[^a-zA-Z]
— \s	whitespace	[\t\n\r\f\v]
— \w	alphabetic	[a-zA-Z]

- Examples
  - \d+ dollars 3 dollars, 50 dollars, 982 dollars
  - $\mathbf{w*oo}\mathbf{w*}$  food, boo, oodles
- Escape character
  - \ is the general escape character; e.g. \. is not a wildcard, but matches a period .
  - if you want to use  $\ in a string it has to be escaped <math>\$

- Anchors. AKA, "zero width characters".
- They match positions in the text.
  - ^ beginning of line
  - \$ end of line
  - $\mathbf{b}$  word boundary, i.e. location with  $\mathbf{w}$  on one side but not on the other.
  - **B** negated word boundary, i.e. any location that would not match \b
- Examples:
  - \bthe\b <u>the</u> together
- Counters {1}, {1,2}, {3,}

- Grouping
  - a (good|bad) movie
  - He said it (again and )\*again.
- Parens also indicate Registers (saved contents)

 $-b(\w+)h\1$ 

matches *boohoo* and *baha*, but not *boohaa* The digit after the \ indicates which of multiple paren groups, as ordered by when then were opened.

Grouping without the cost of register saving
 - He went (?:this|that) way.

- Non-consuming tests
  - (?=...) Positive lookAHEAD
  - (?!...) Negative lookAHEAD
  - (?<=...) Positive lookBEHIND</pre>
  - (?<!...) Negative lookBEHIND</pre>
- Example
  - $-(? \leq [Hh]e) \setminus w + ed (? \leq w + ly)$

# ELIZA (Weizenbaum 1966)

User: Men are all alike.

ELIZA: IN WHAT WAY?

- User: They're always bugging us about something or other.
- ELIZA: CAN YOU THINK OF A SPECIFIC EXAMPLE?
- User: Well, my boyfriend made me come here.
- ELIZA: YOUR BOYFRIEND MADE YOU COME HERE.
- User: He say's I'm depressed much of the time.
- ELIZA: I'M SORRY TO HEAR THAT YOU ARE DEPRESSED.

#### Implemented with regular expression substitution!

s/.\* I'm (depressed|sad) .\*/I AM SORRY TO HEAR THAT YOU ARE \1/ s/.\* always .\*/CAN YOU THINK OF A SPECIFIC EXAMPLE?/

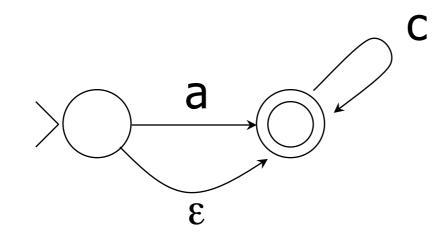
### Reading

 Karttunen, Chanod, Grefenstette, Schiller.
 Regular expressions for language engineering. JNLE, 1997.

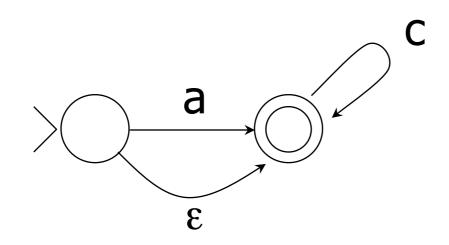
<u>http://www.stanford.edu/~laurik/</u> <u>publications/jnle-97/rele.pdf</u>

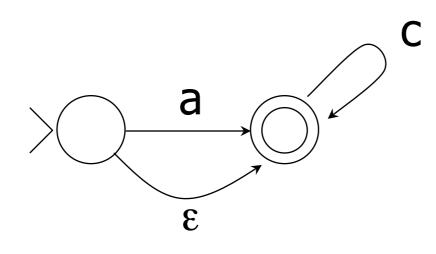
• RE/FSA background: Jurafsky & Martin, c.2

#### Finite-State Machines: Acceptors and Transducers

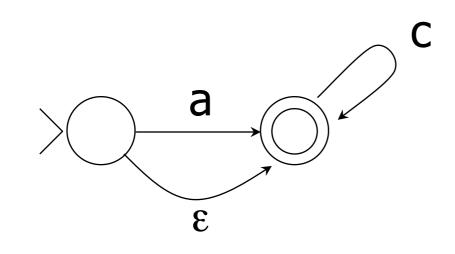




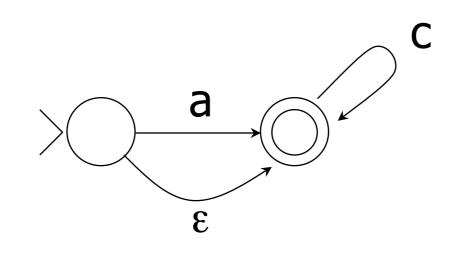




 Regexps
 Union, Kleene \*, concat, intersect, complement, reversal



- Regexps
- Union, Kleene \*, concat, intersect, complement, reversal
- Determinization, minimization



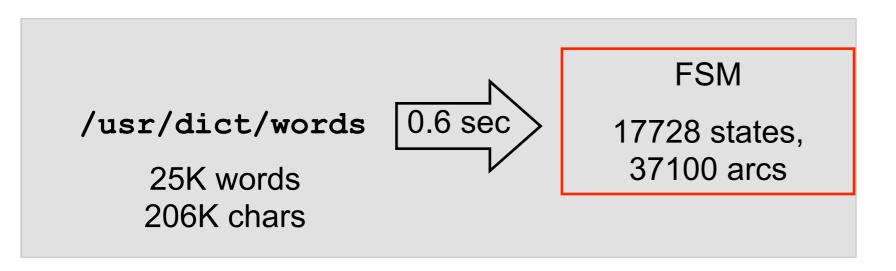
- Regexps
- Union, Kleene \*, concat, intersect, complement, reversal
- Determinization, minimization
- Pumping, Myhill-Nerode

slide courtesy of L. Karttunen (modified)

#### A useful FSA ...

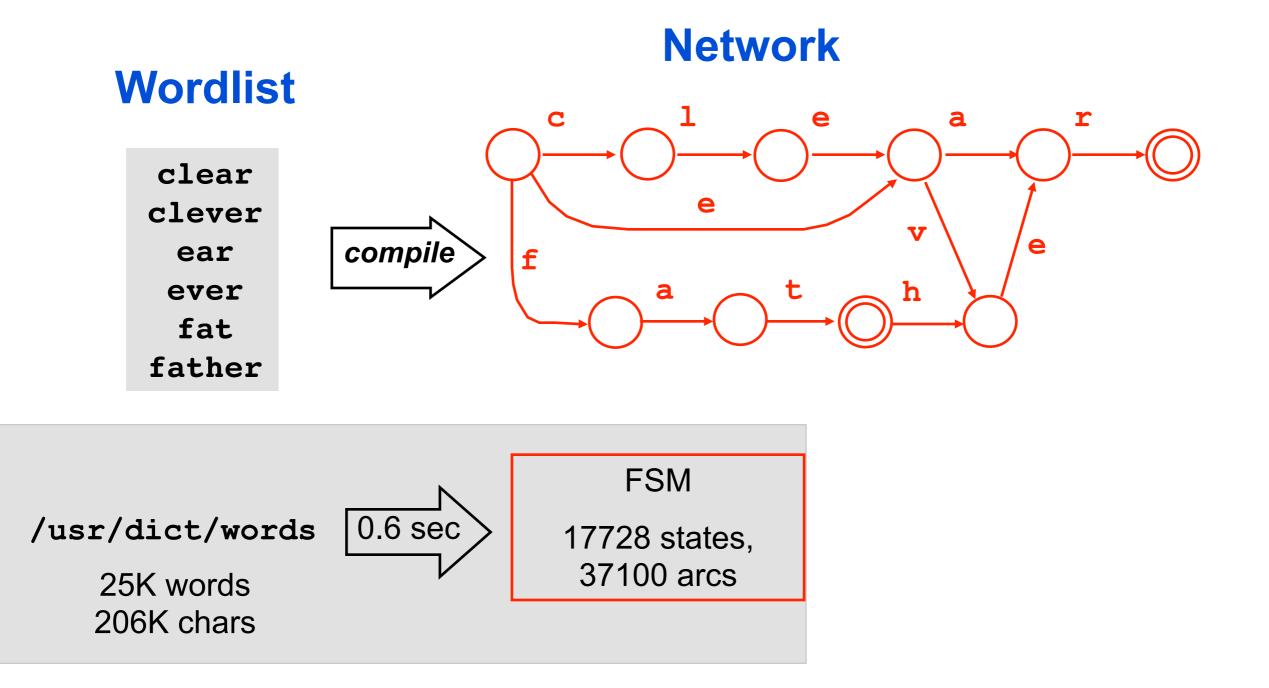
#### Wordlist

clear clever ear ever fat father

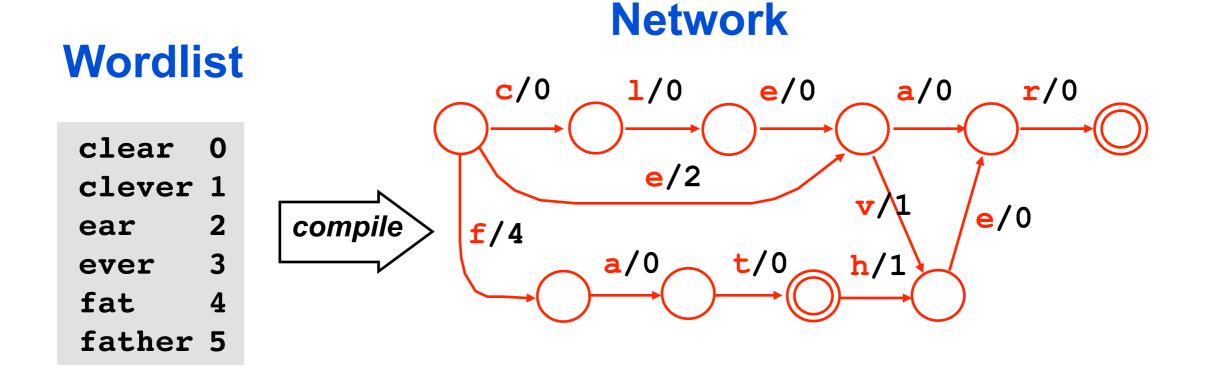


slide courtesy of L. Karttunen (modified)

#### A useful FSA ...

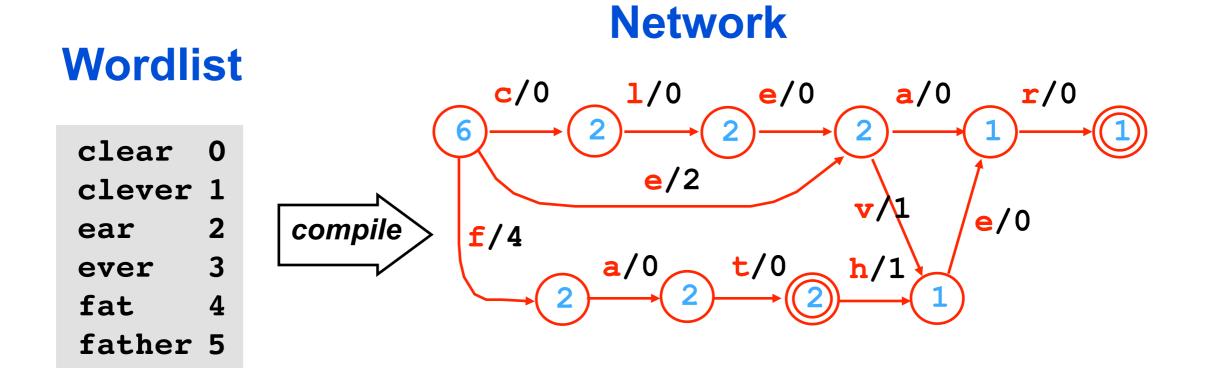


#### Weights are useful here too!

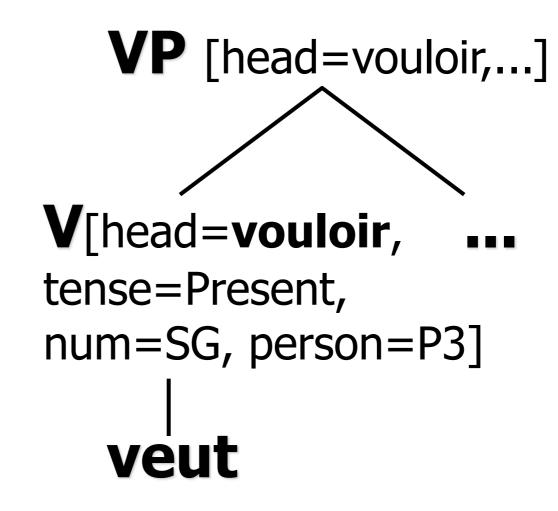


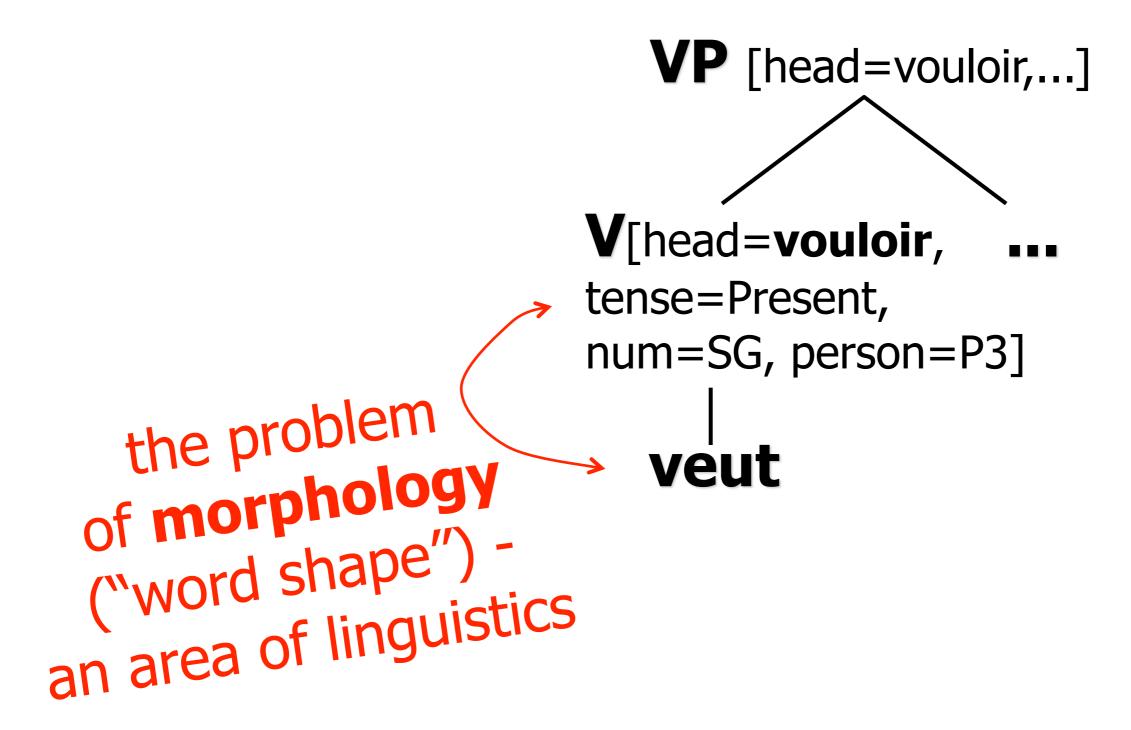
#### Computes a perfect hash!

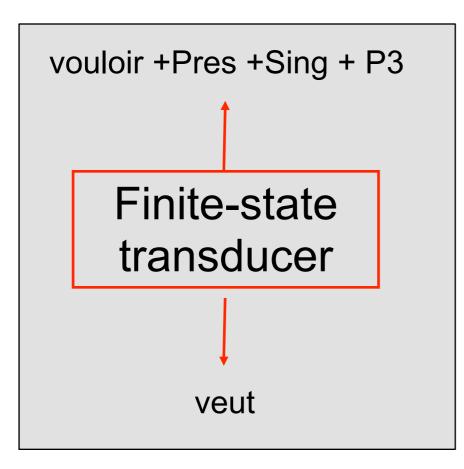
#### **Example: Weighted acceptor**

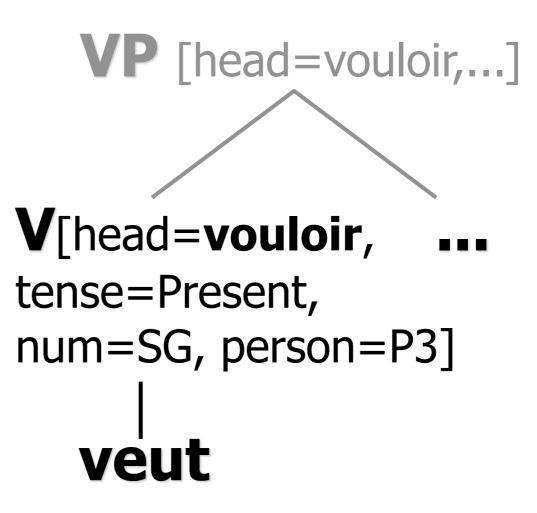


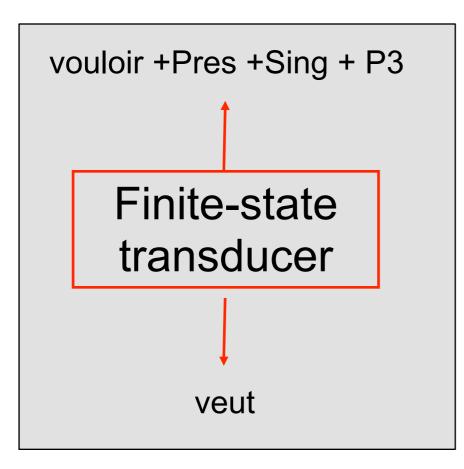
- Compute number of paths from each state (Q: how?) A: recursively, like DFS
- Successor states partition the path set
- Use offsets of successor states as arc weights
- **Q**: Would this work for an arbitrary numbering of the words?

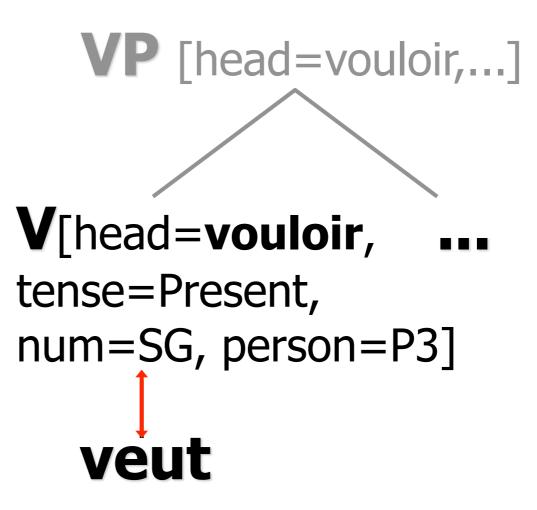


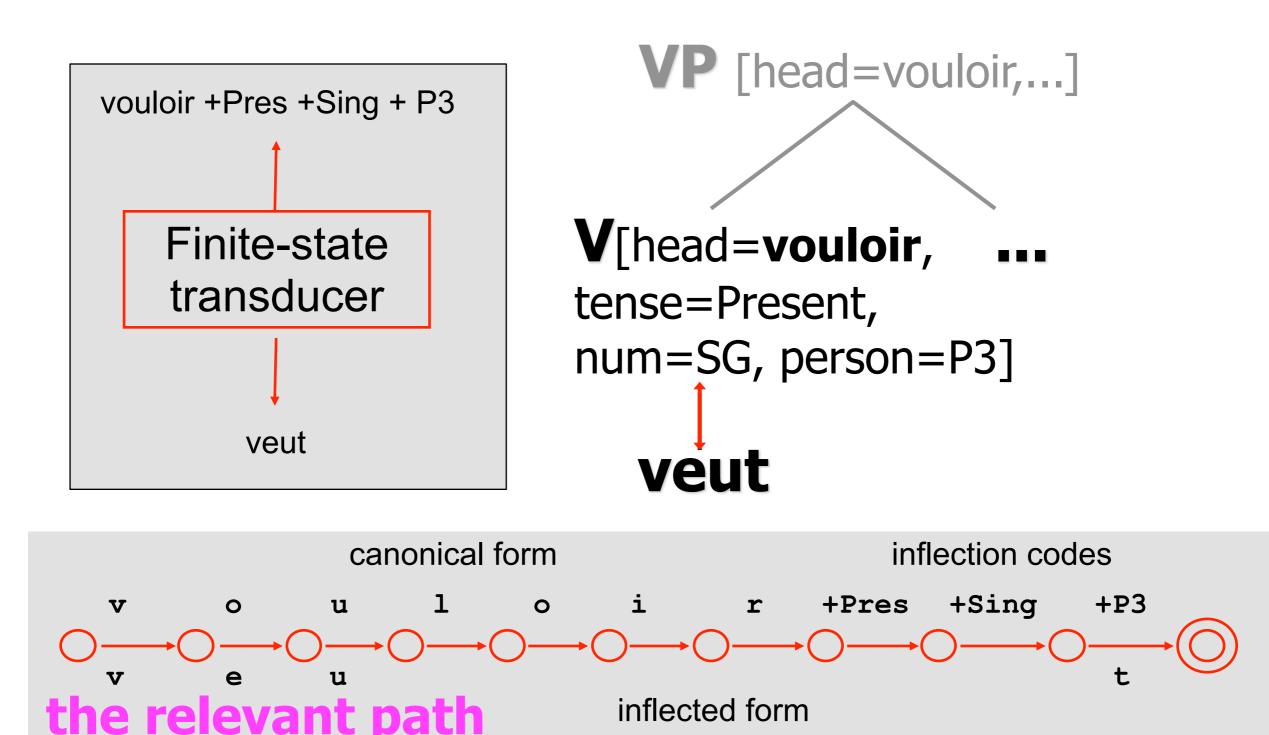


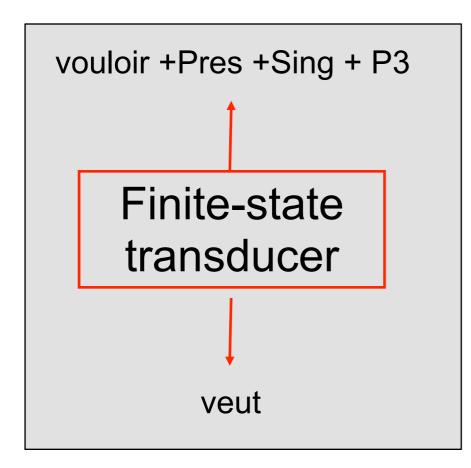




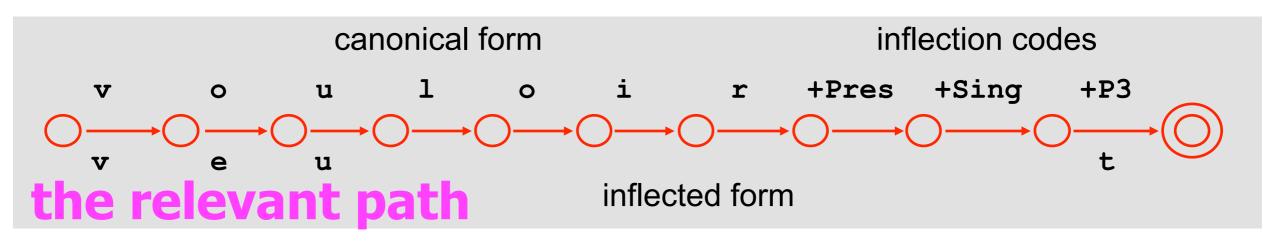


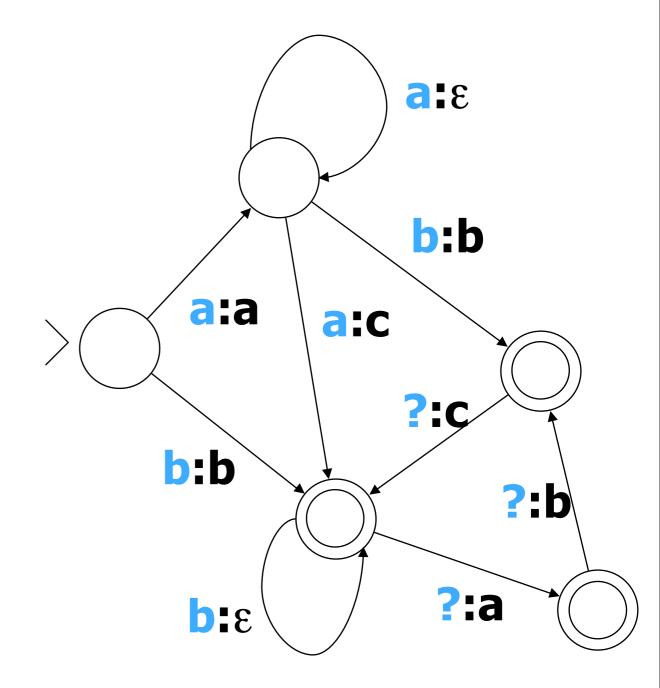




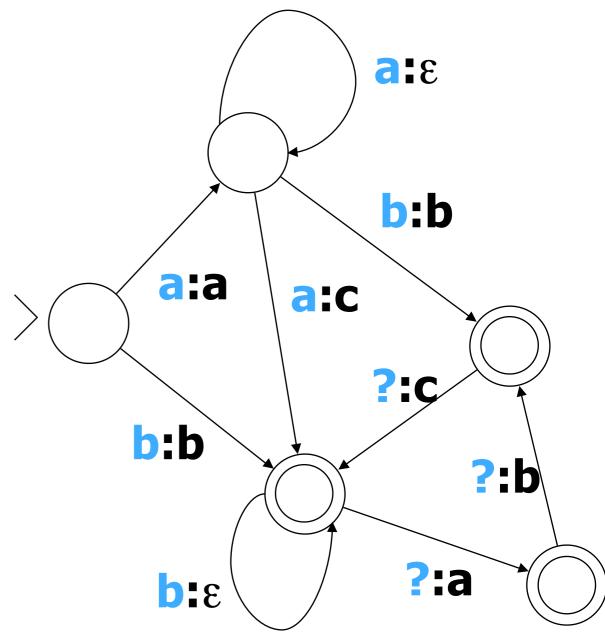


- Bidirectional: generation or analysis
- Compact and fast
- Xerox sells for about 20 languges including English, German, Dutch, French, Italian, Spanish, Portuguese, Finnish, Russian, Turkish, Japanese, ...
- Research systems for many other languages, including Arabic, Malay

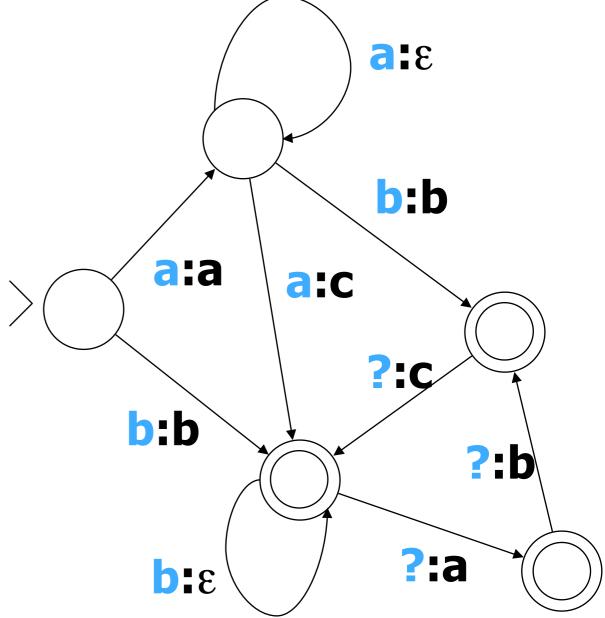




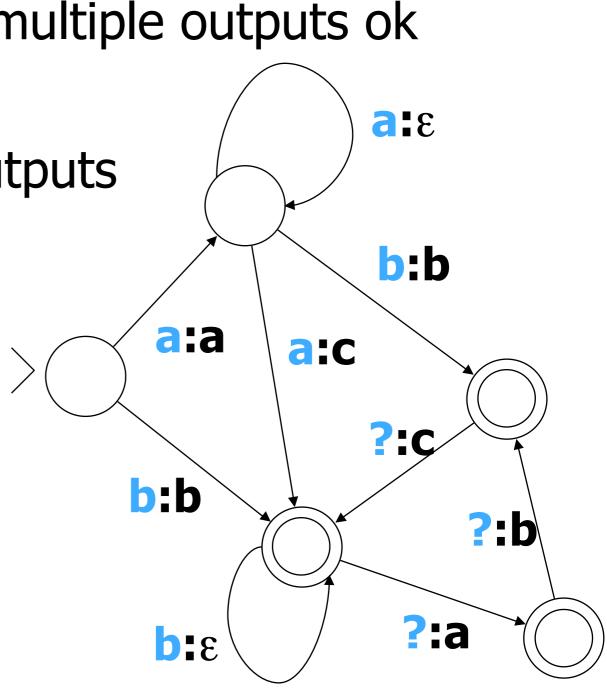
Relation: like a function, but multiple outputs ok



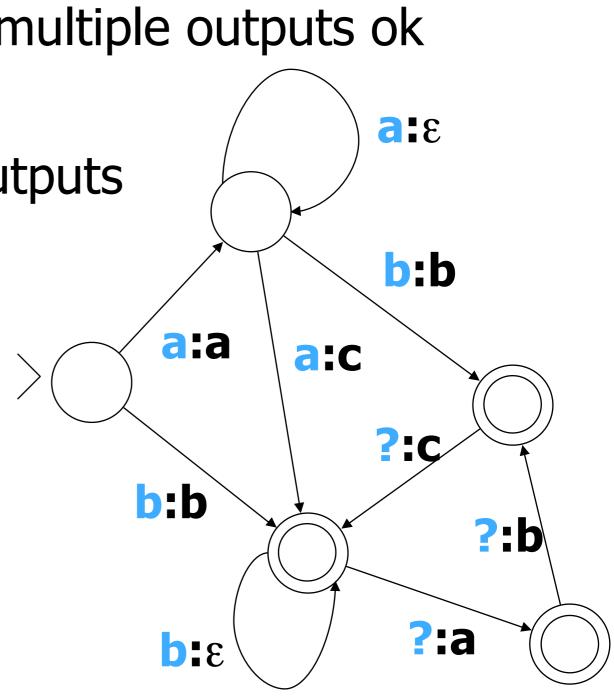
Relation: like a function, but multiple outputs ok
Regular: finite-state



- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs

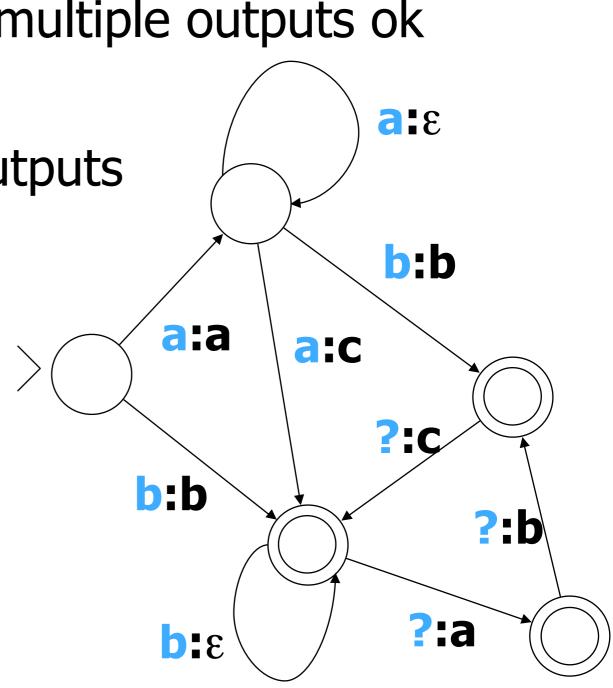


- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs
- $\bullet \rightarrow ? \qquad a \rightarrow ?$



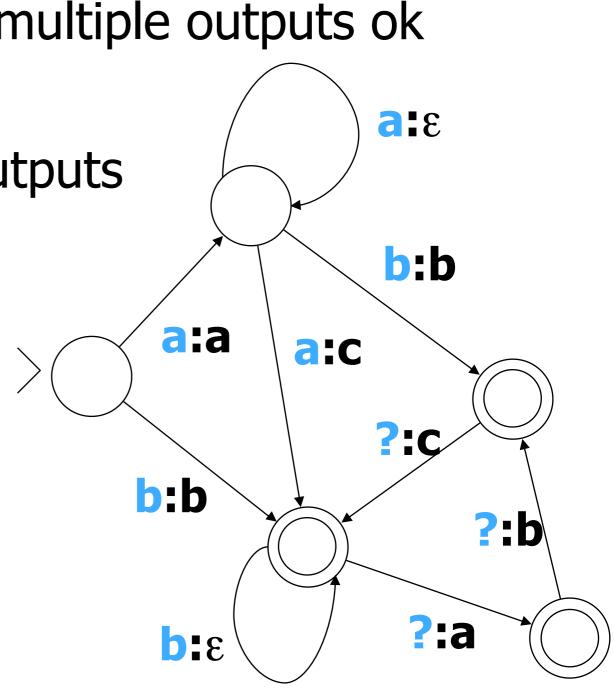
- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs

• b 
$$\rightarrow$$
 {b} a  $\rightarrow$  ?

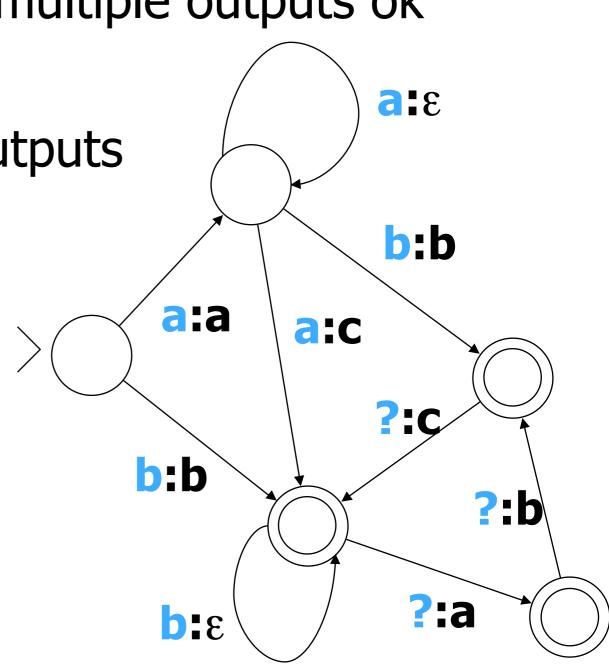


- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs

• b 
$$\rightarrow$$
 {b} a  $\rightarrow$  {}

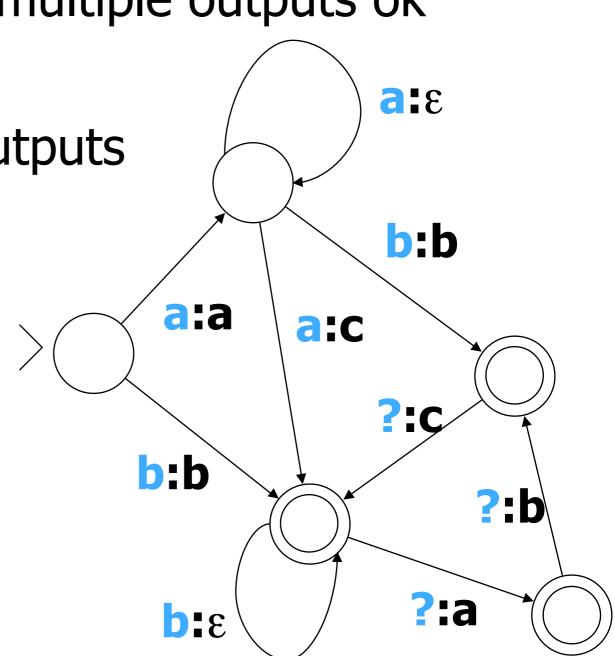


- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs



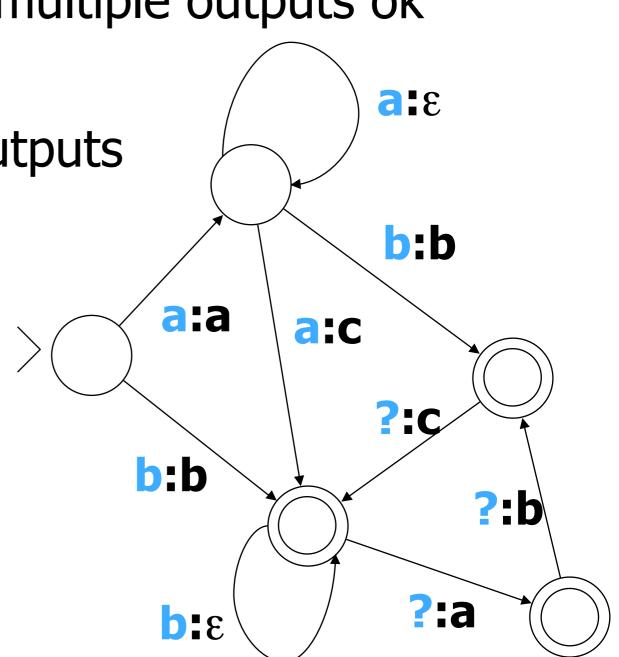
- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs

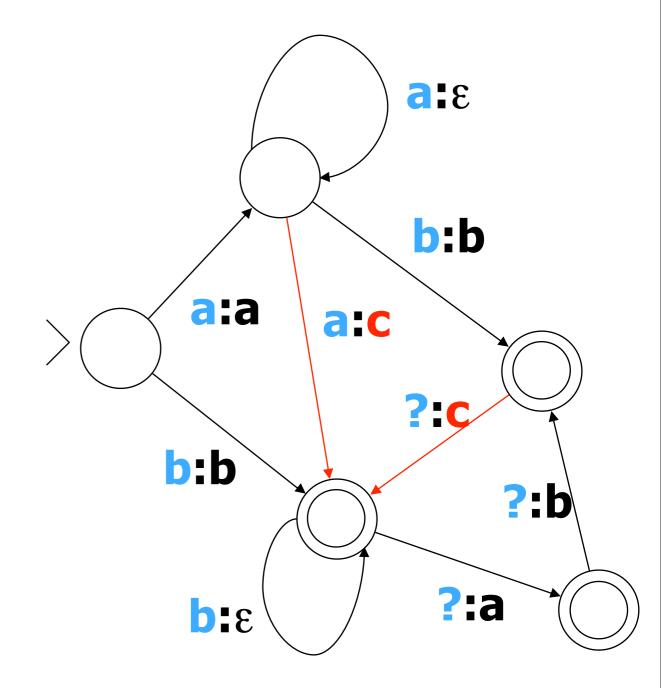
Invertible?



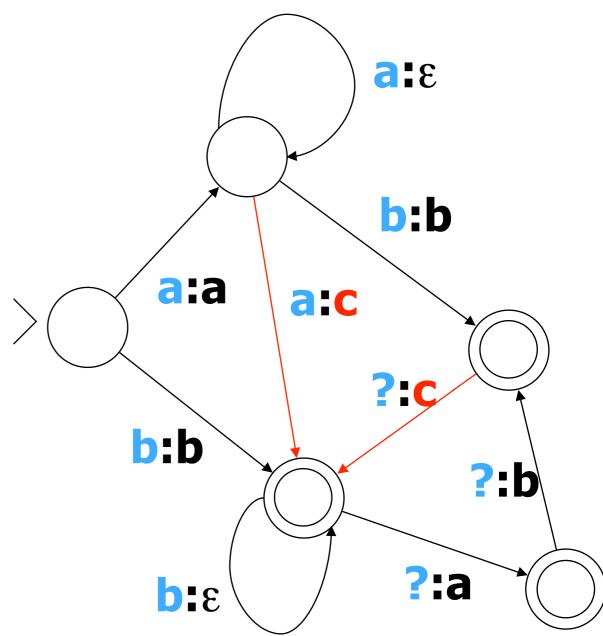
- Relation: like a function, but multiple outputs ok
- Regular: finite-state
- Transducer: automaton w/ outputs

- Invertible?
- Closed under composition?

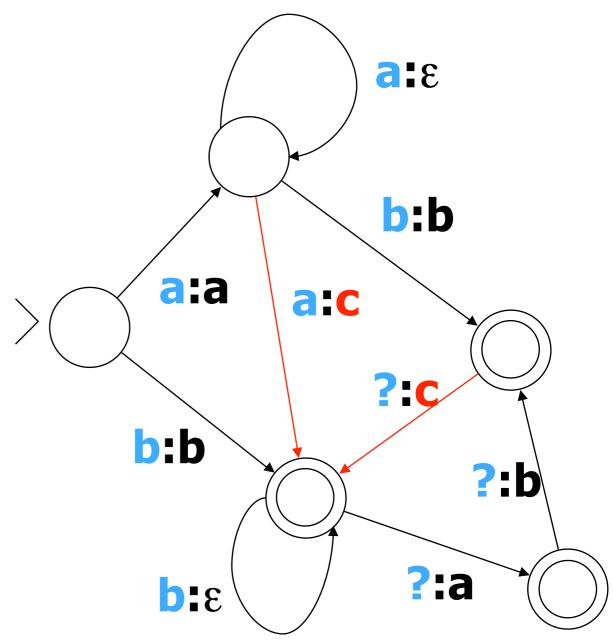




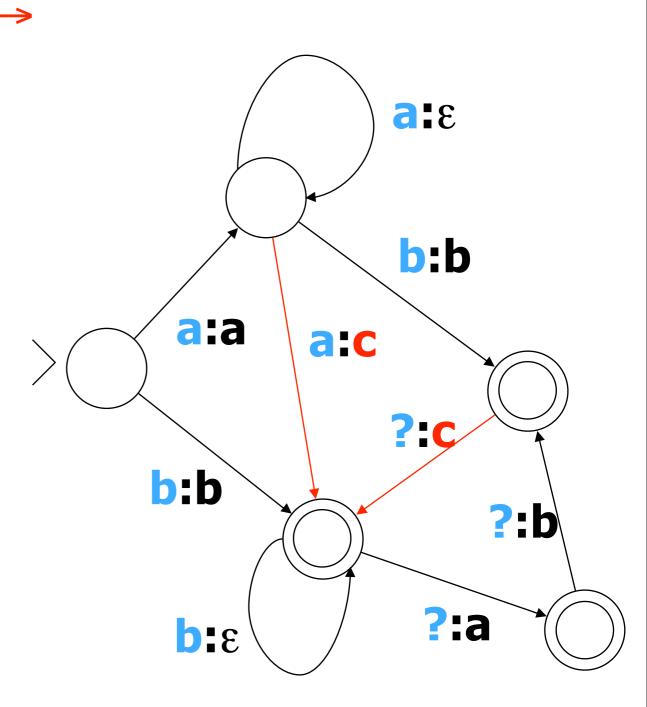
• Can weight the arcs:  $\rightarrow$  vs.  $\rightarrow$ 



Can weight the arcs:  $\rightarrow$  vs.  $\rightarrow$ b  $\rightarrow$  {b} a  $\rightarrow$  {}

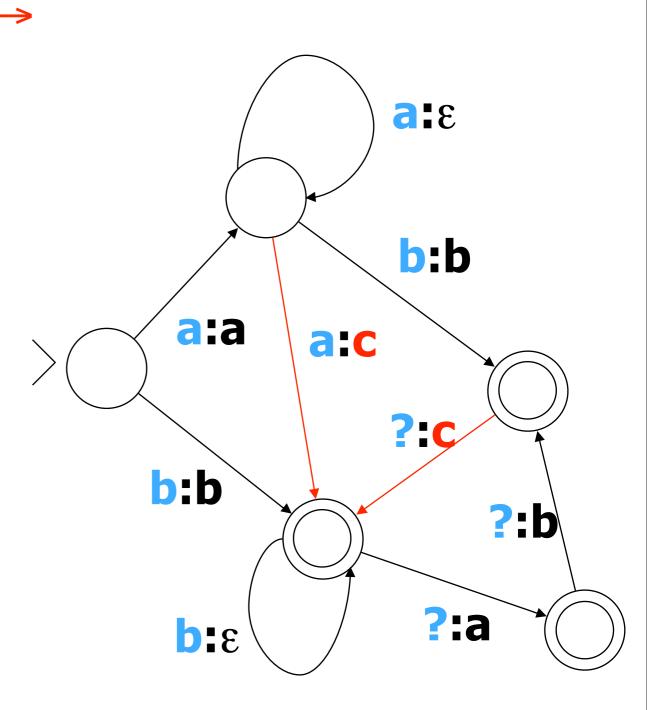


Can weight the arcs: → vs. →
b → {b} a → {}
aaaaa → {ac, aca, acab, acabc}



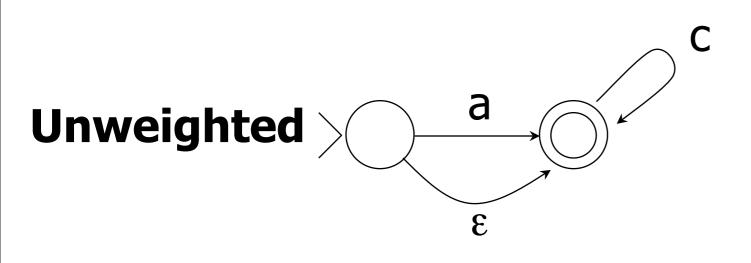
Can weight the arcs: → vs. →
b → {b} a → {}
aaaaa → {ac, aca, acab, acabc}

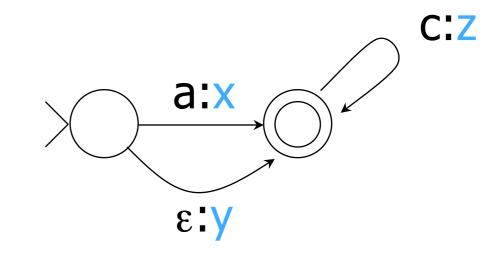
How to find <u>best</u> outputs?

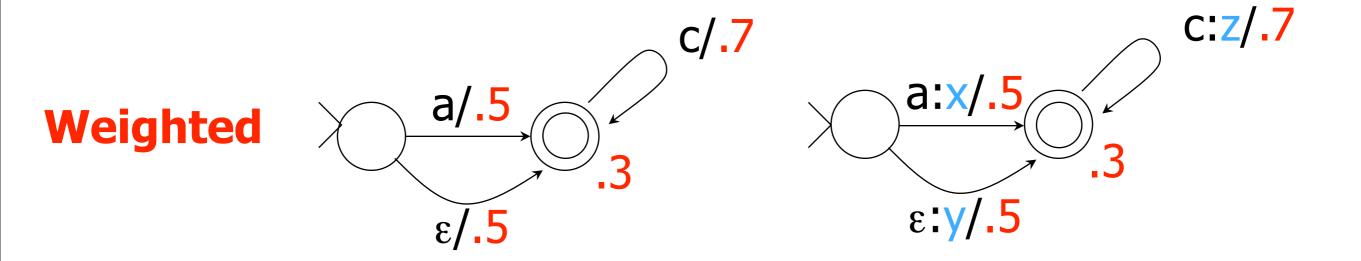


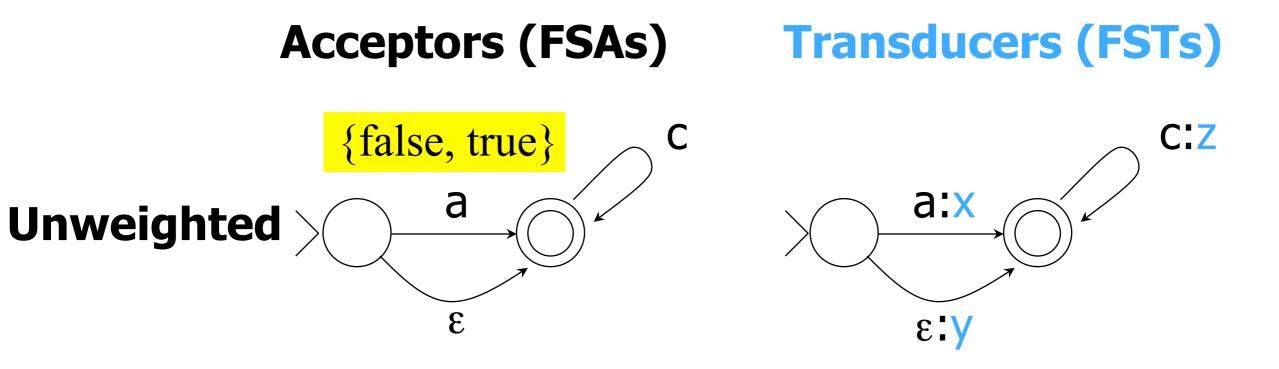
Acceptors (FSAs)

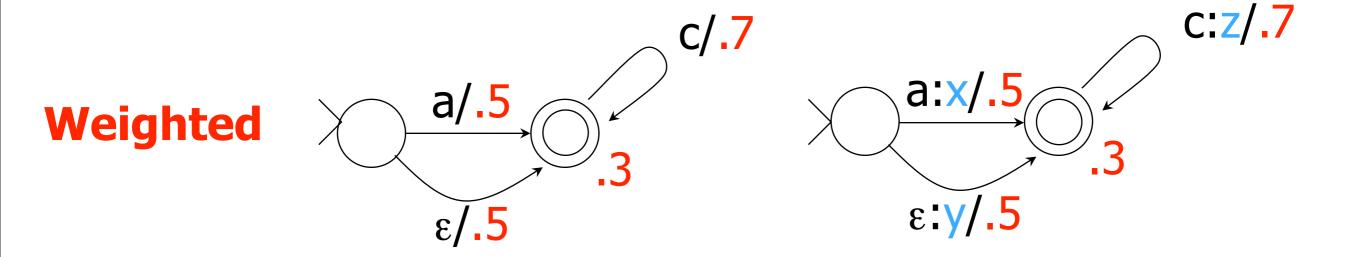
**Transducers (FSTs)** 

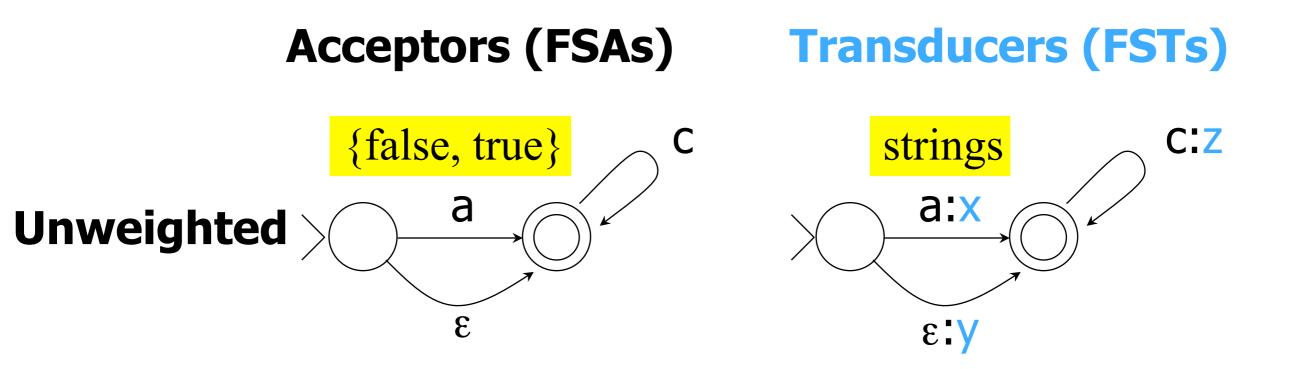


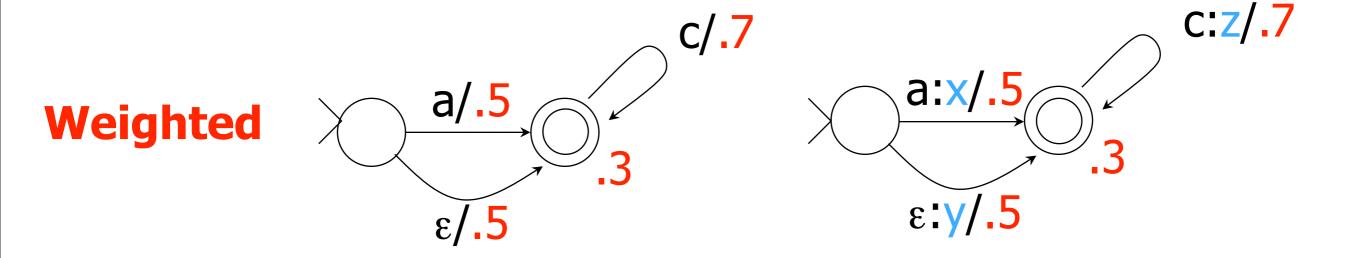


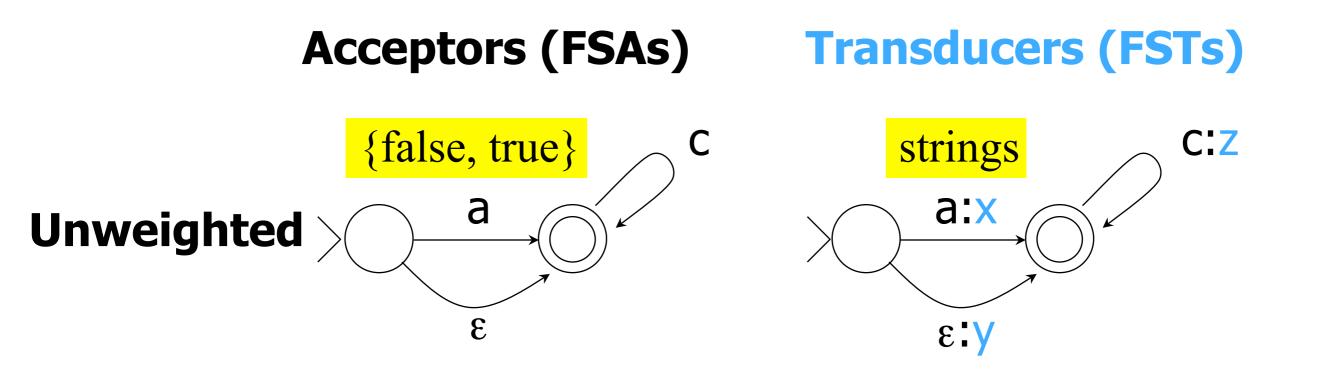


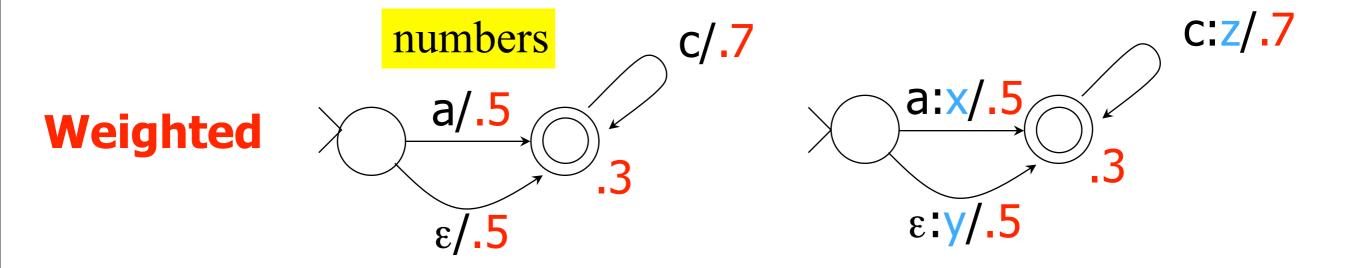


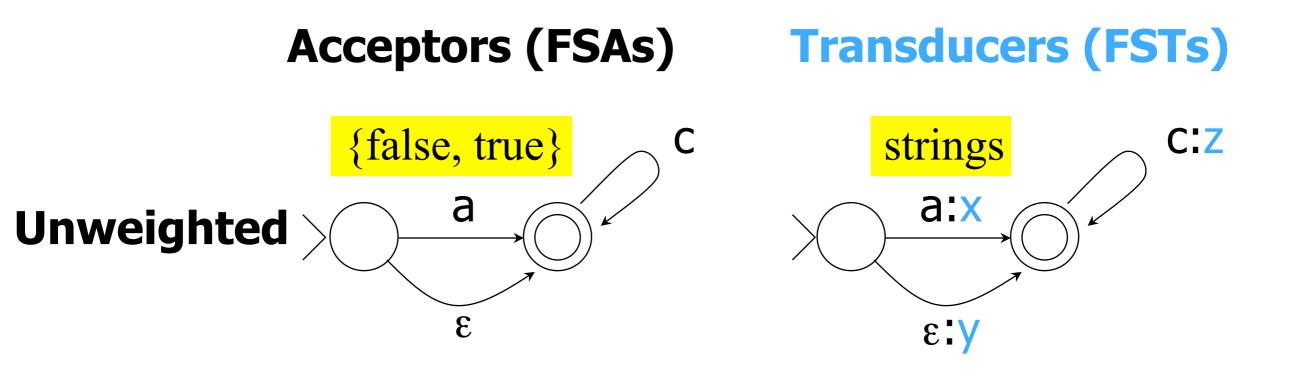


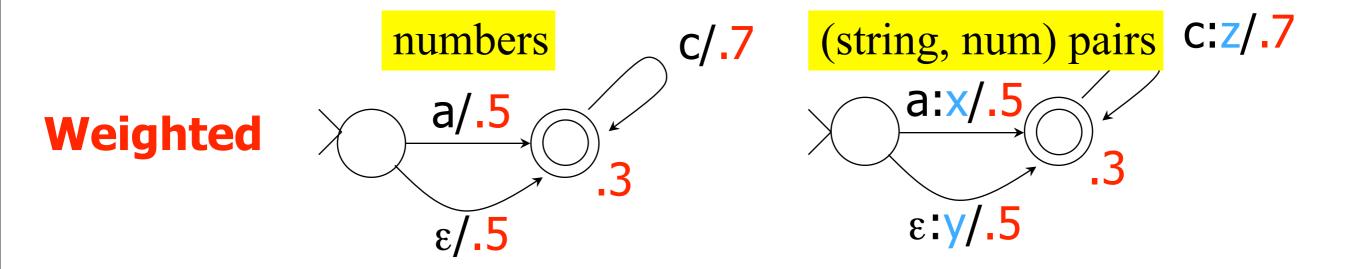












Acceptors (FSAs)

{false, true}

Unweighted

numbers

(string, num) pairs

**Transducers (FSTs)** 

strings

Weighted

Acceptors (FSAs)

#### {false, true}

**Unweighted** Grammatical?

Transducers (FSTs)

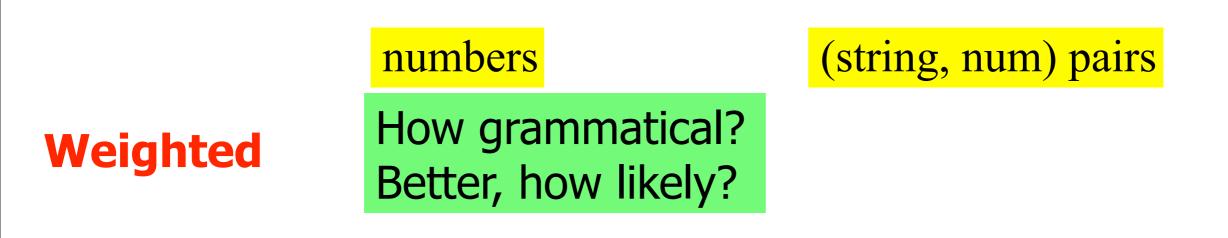
strings

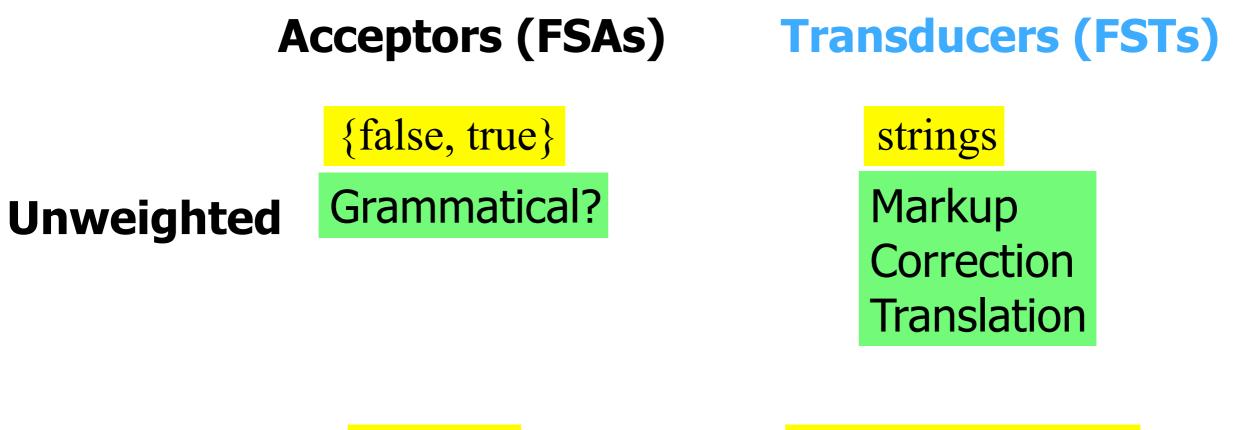
numbers

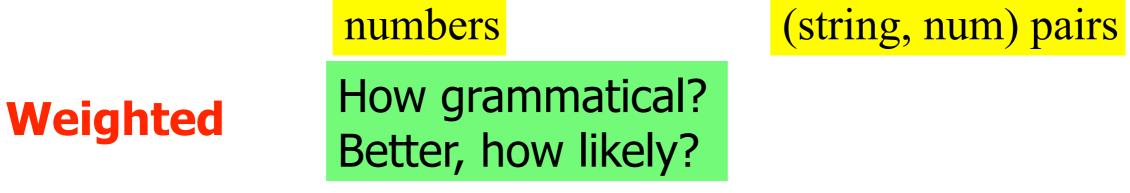
(string, num) pairs

Weighted

Acceptors (FSAs)Transducers (FSTs){false, true}stringsUnweightedGrammatical?







Acceptors (FSAs)

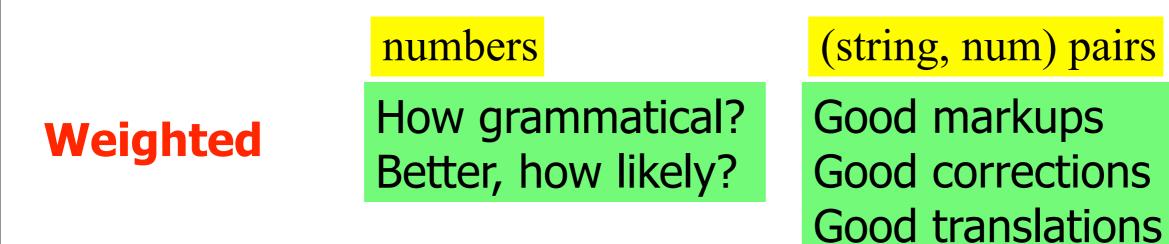
{false, true}

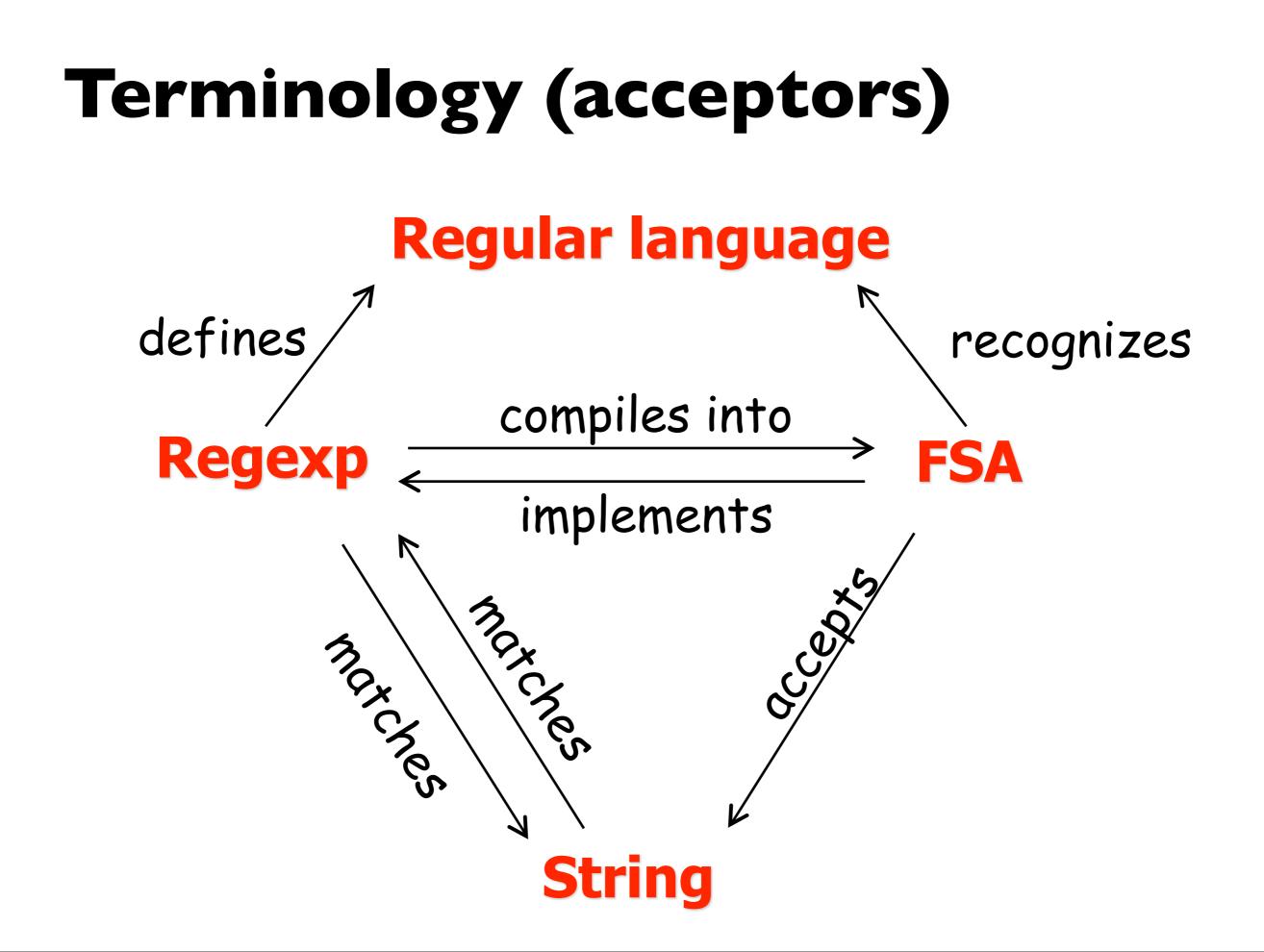
**Unweighted** Grammatical?

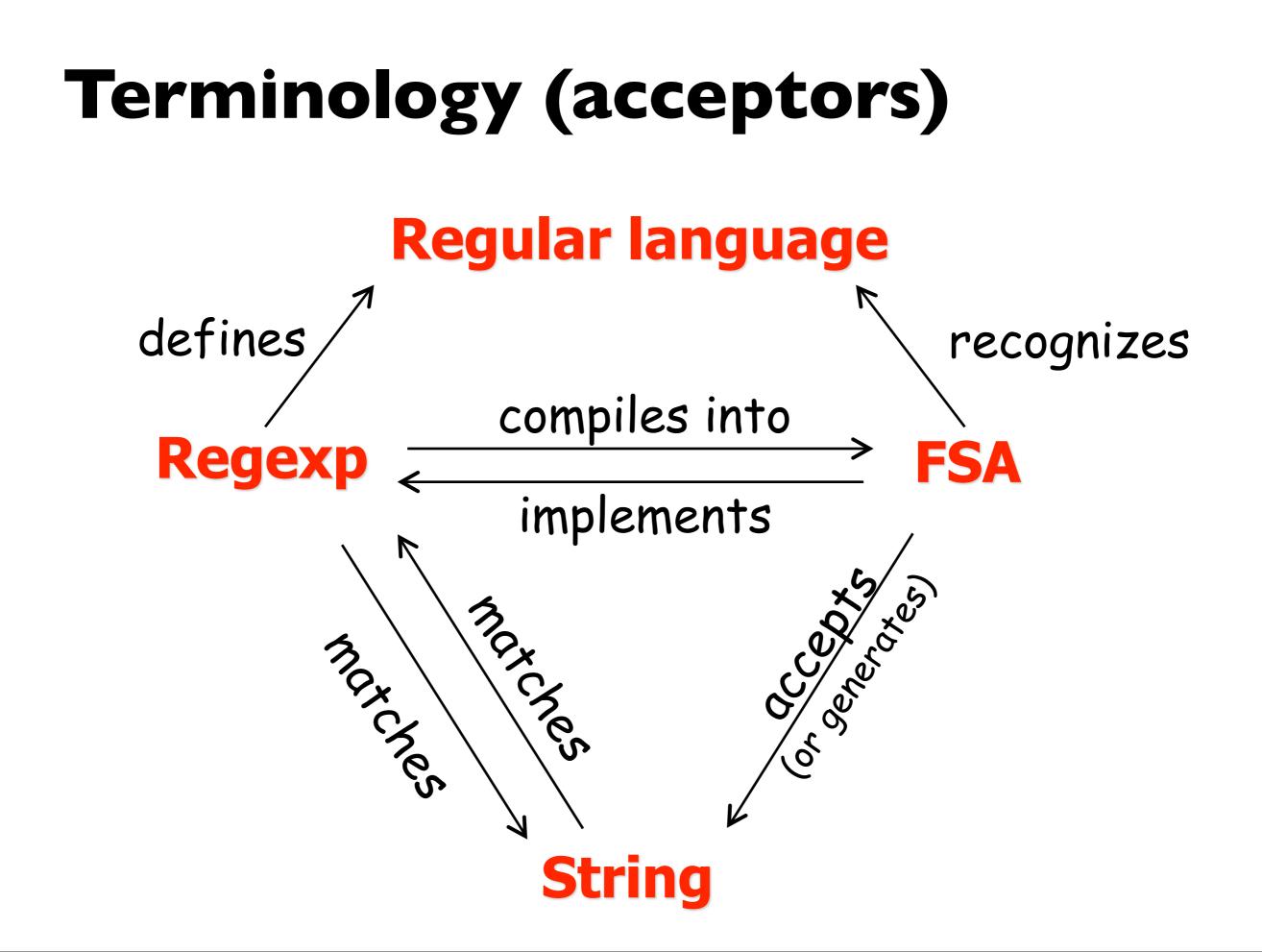
Transducers (FSTs)

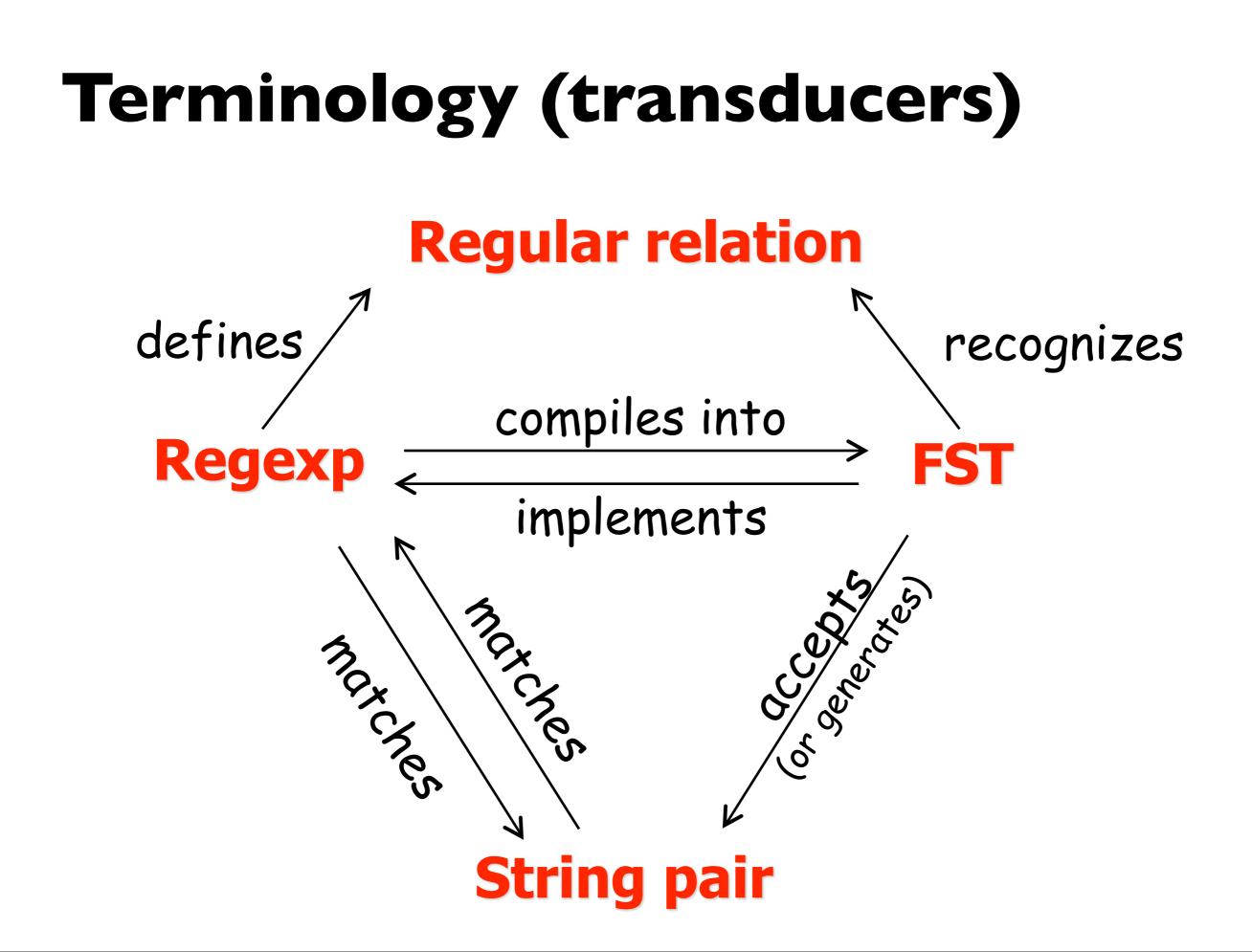
strings

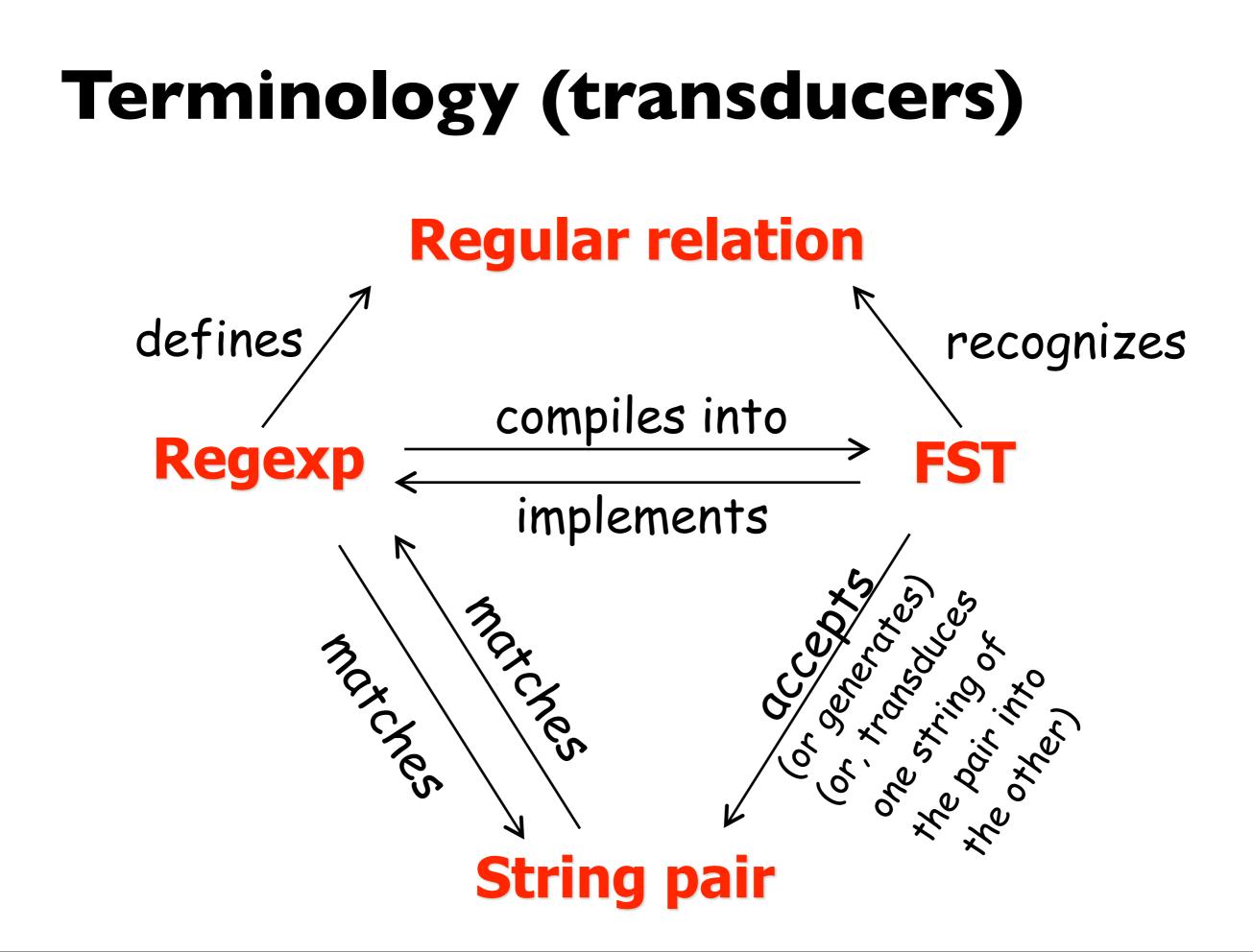
Markup Correction Translation

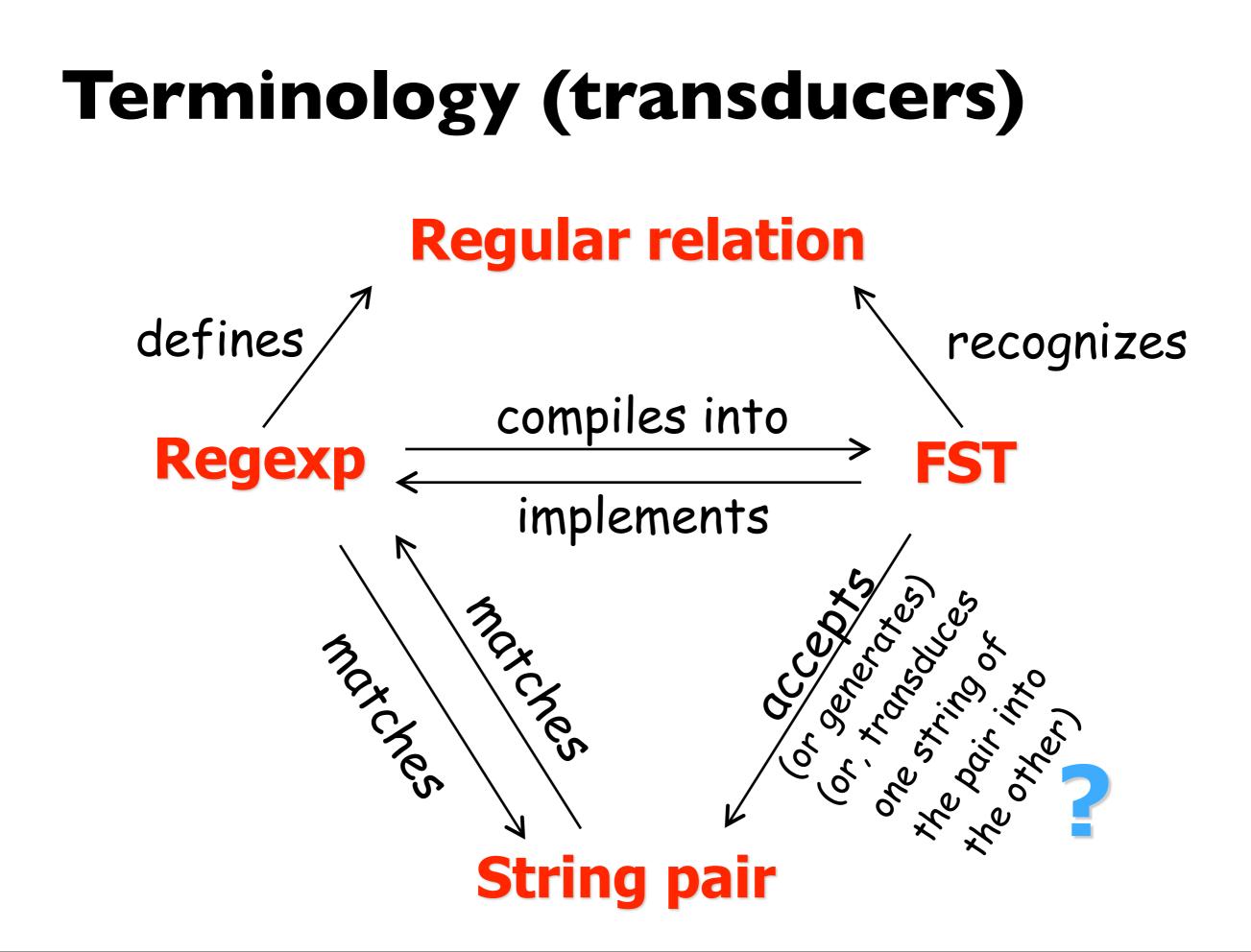






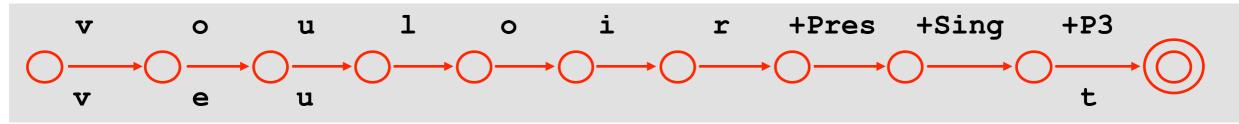






#### **Perspectives on a Transducer**

- Given 0 strings, generate a new string pair (by picking a path)
- Given one string (upper or lower), transduce it to the other kind
- Given two strings (upper & lower), decide whether to accept the pair

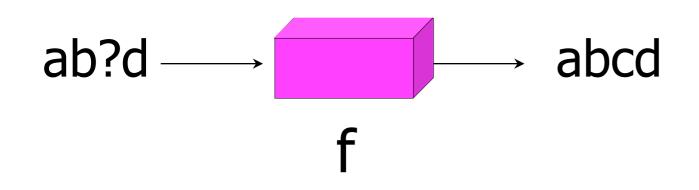


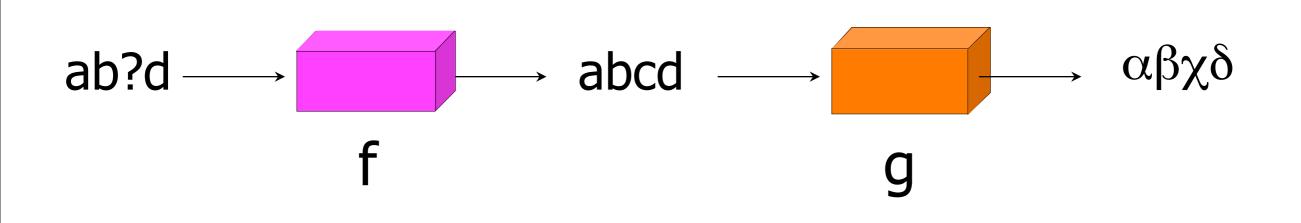
FST just defines the regular relation (mathematical object: set of pairs).

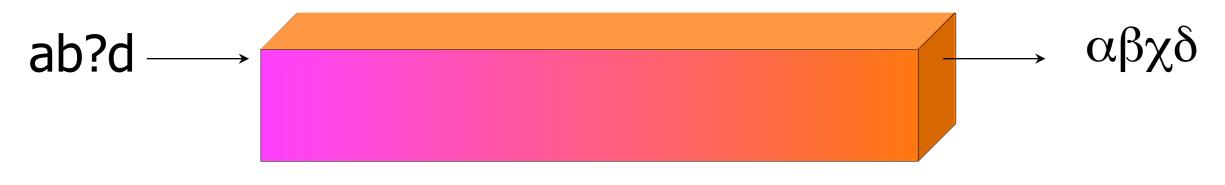
What's "input" and "output" depends on what one <u>asks</u> about the relation.

The 0, 1, or 2 given string(s) constrain which paths you can use.

ab?d →

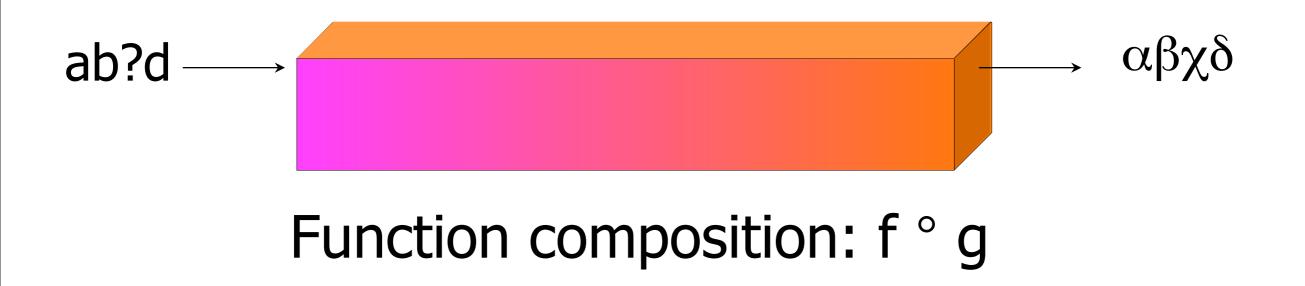




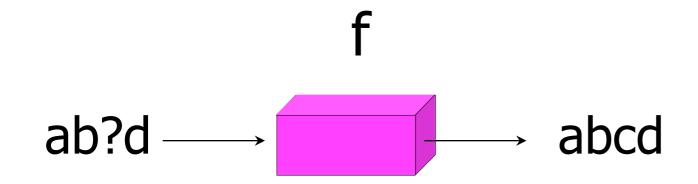


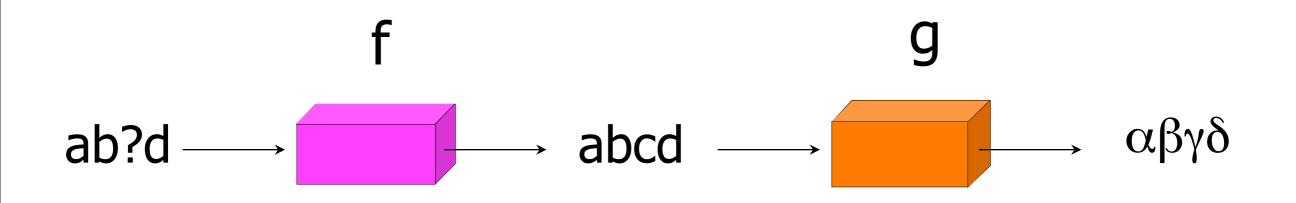
#### Function composition: f ° g

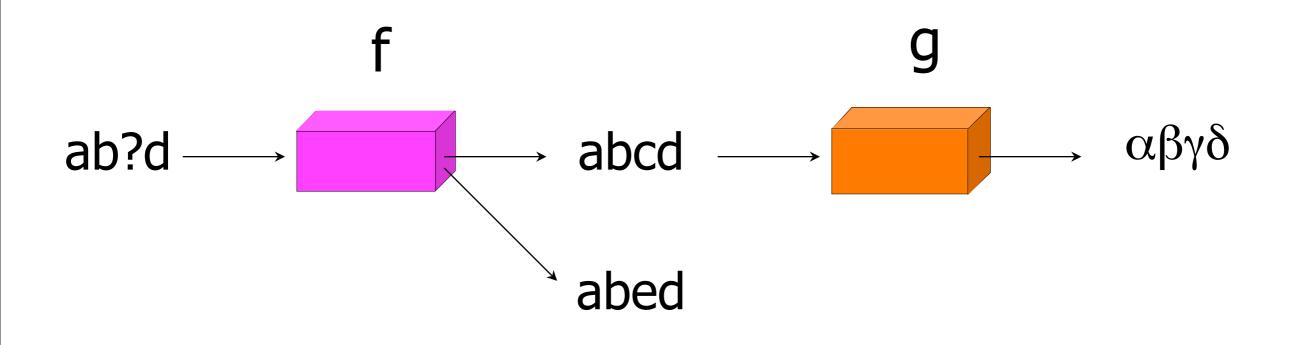


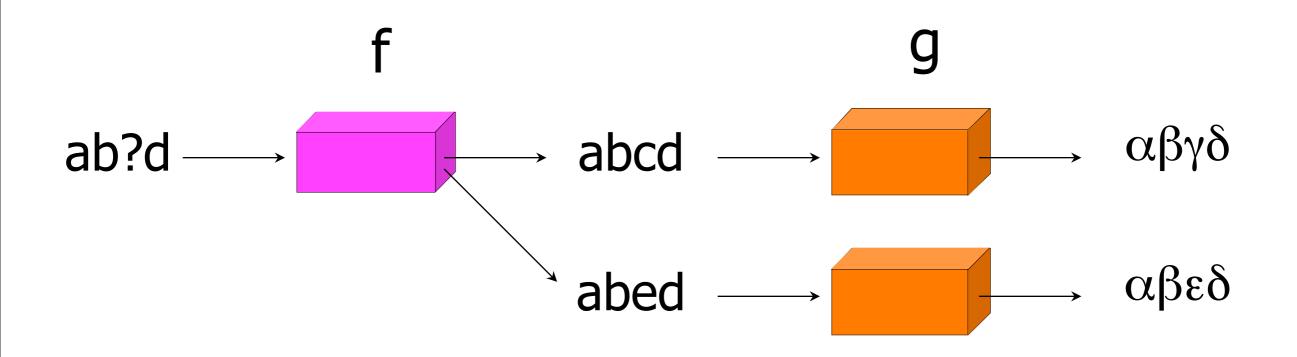


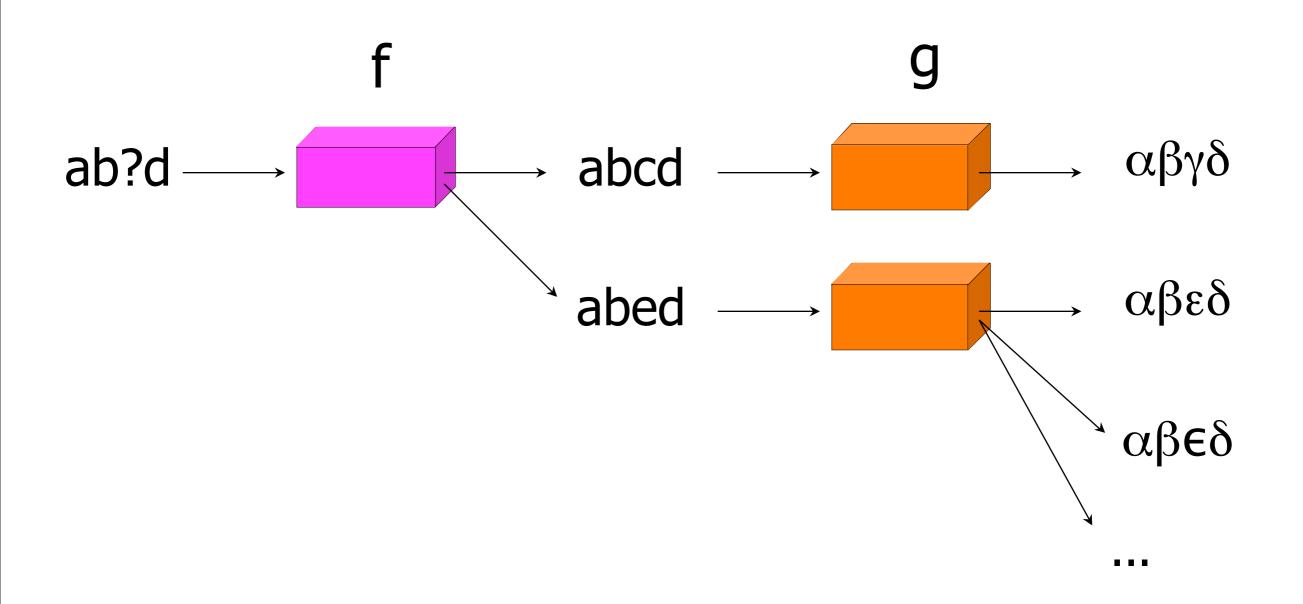
[first f, then g – intuitive notation, but opposite of the traditional math notation]

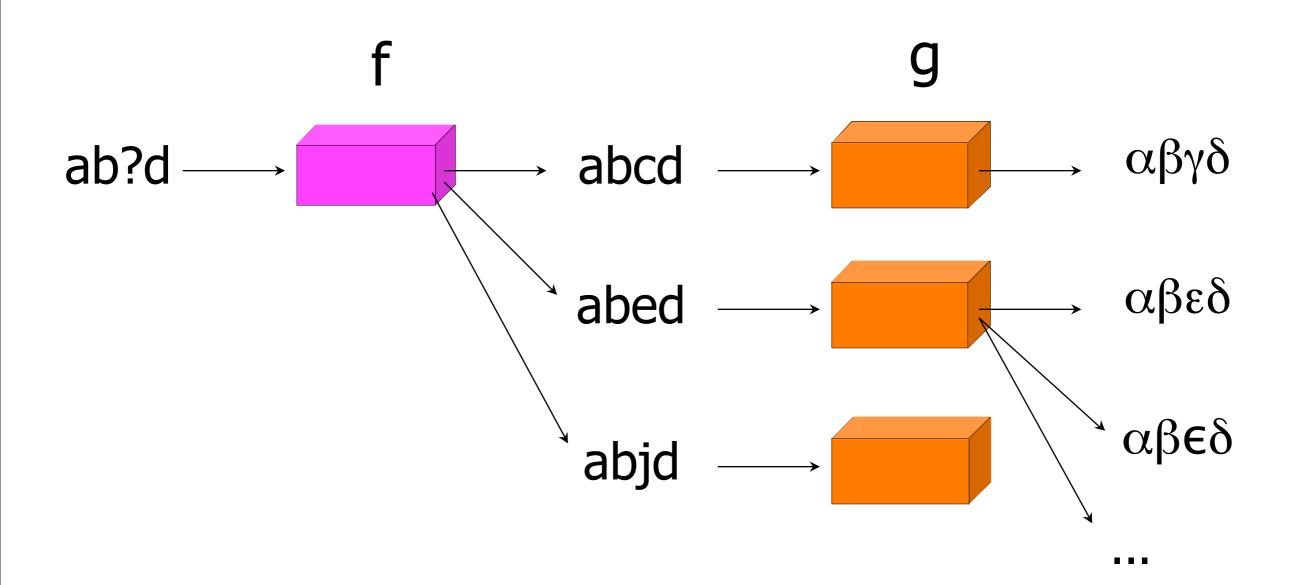


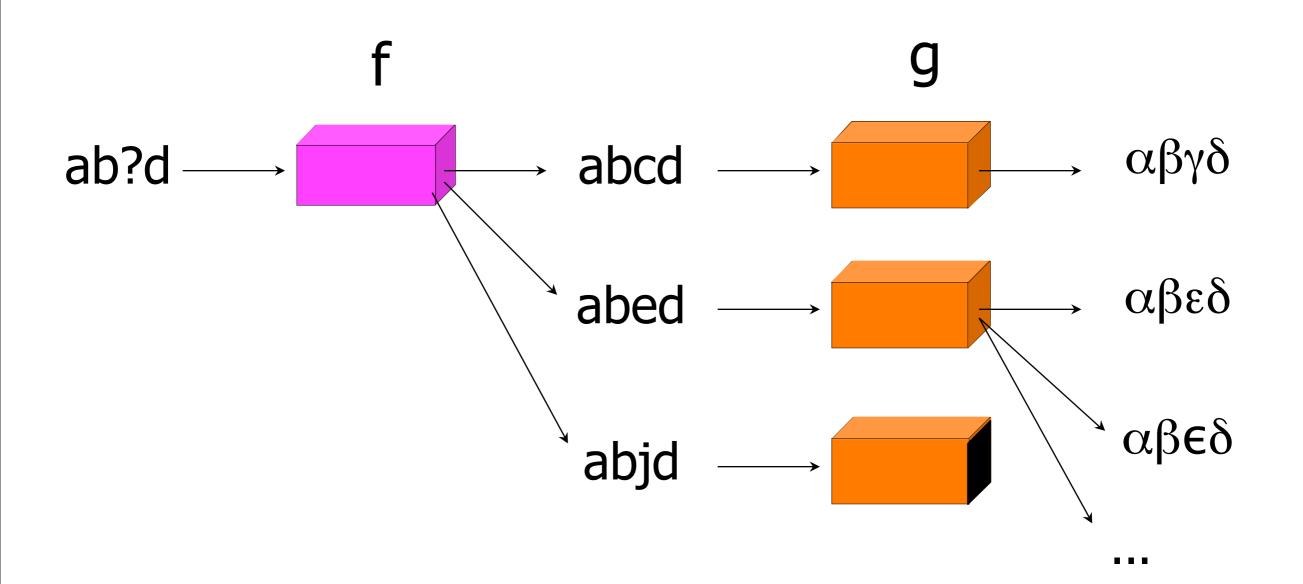


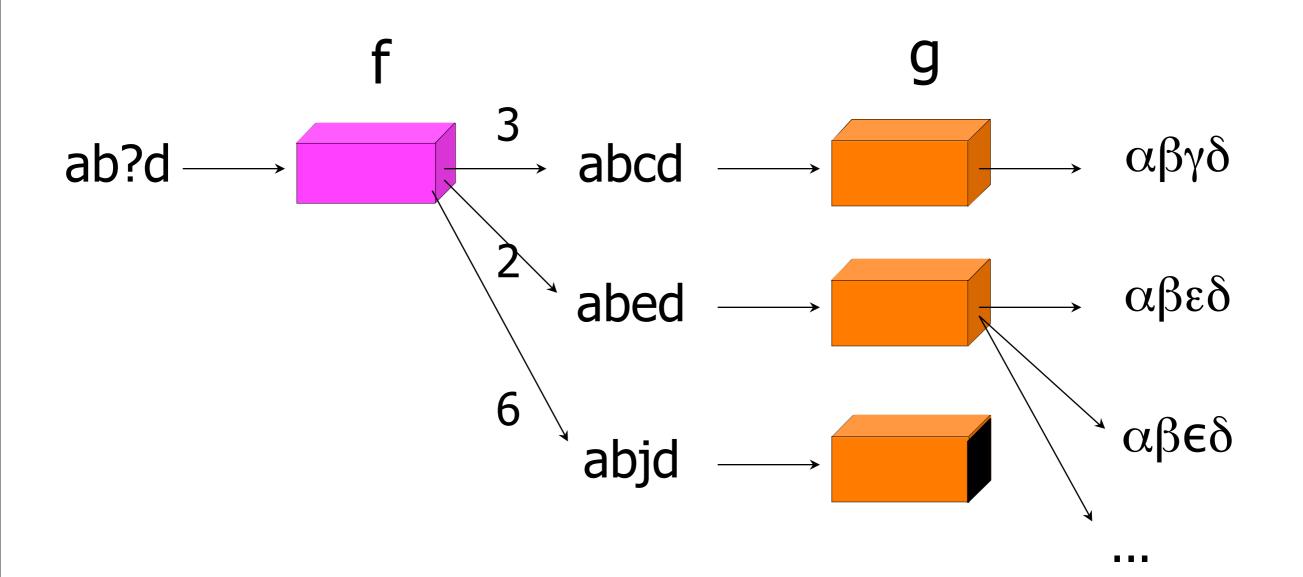


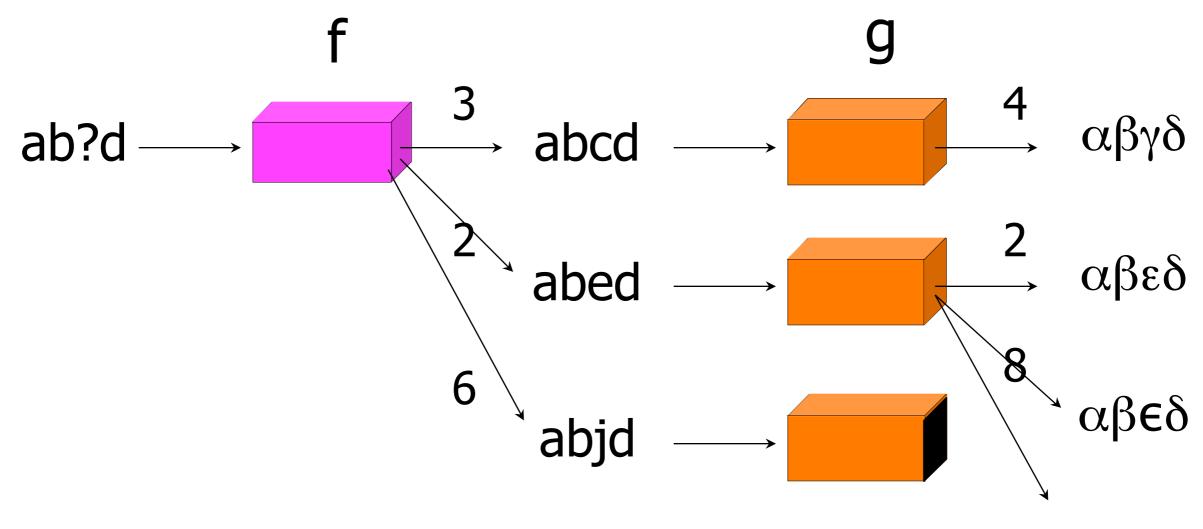


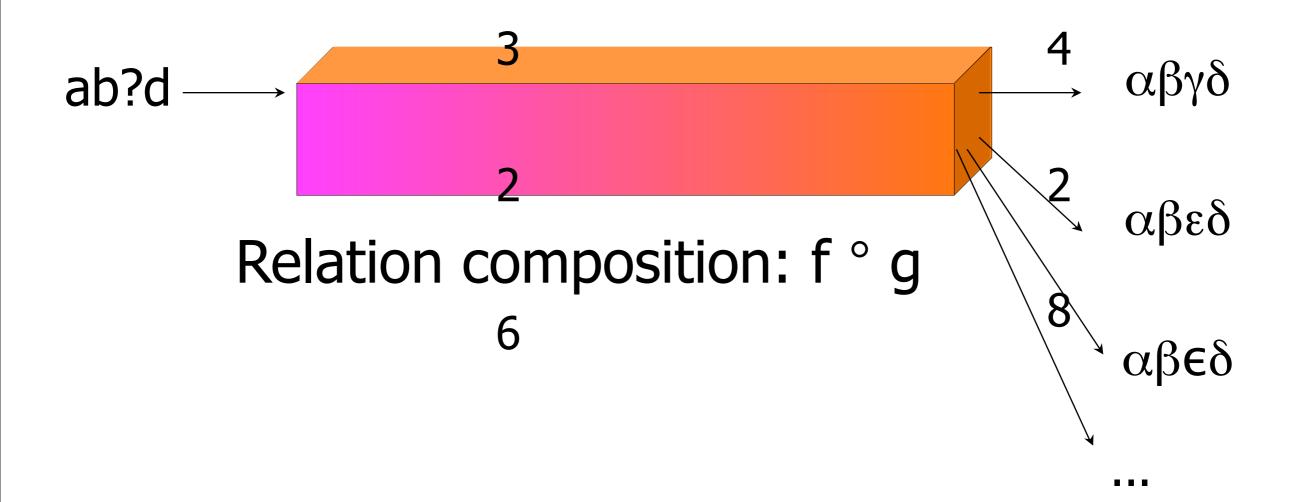














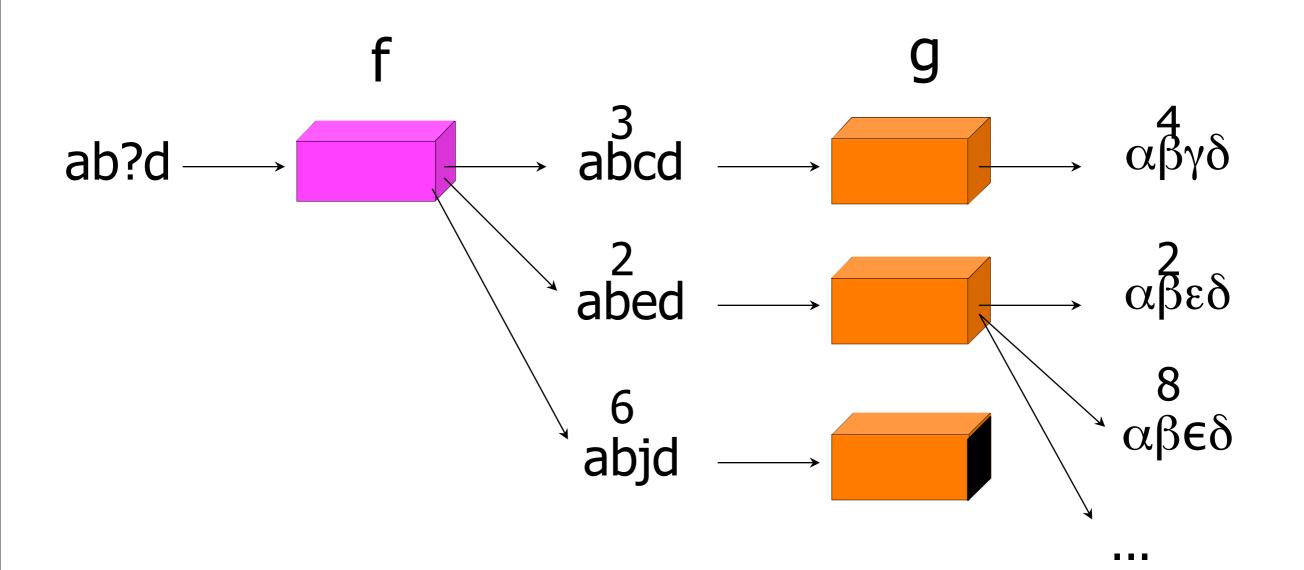




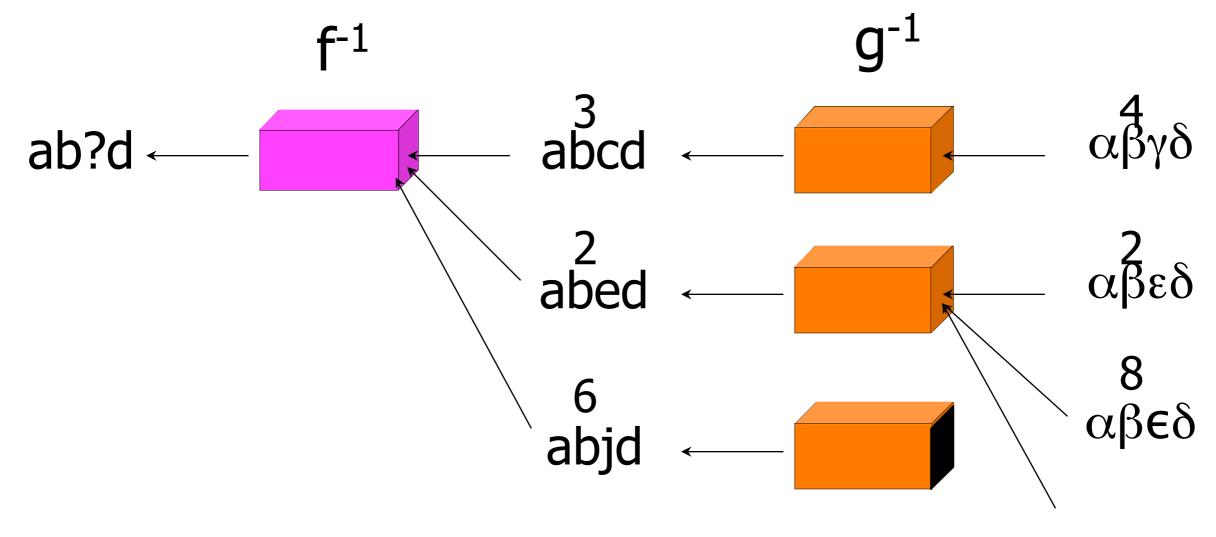


Often in NLP, all of the functions or relations involved can be described as finite-state machines, and manipulated using standard algorithms.

# **Inverting Relations**

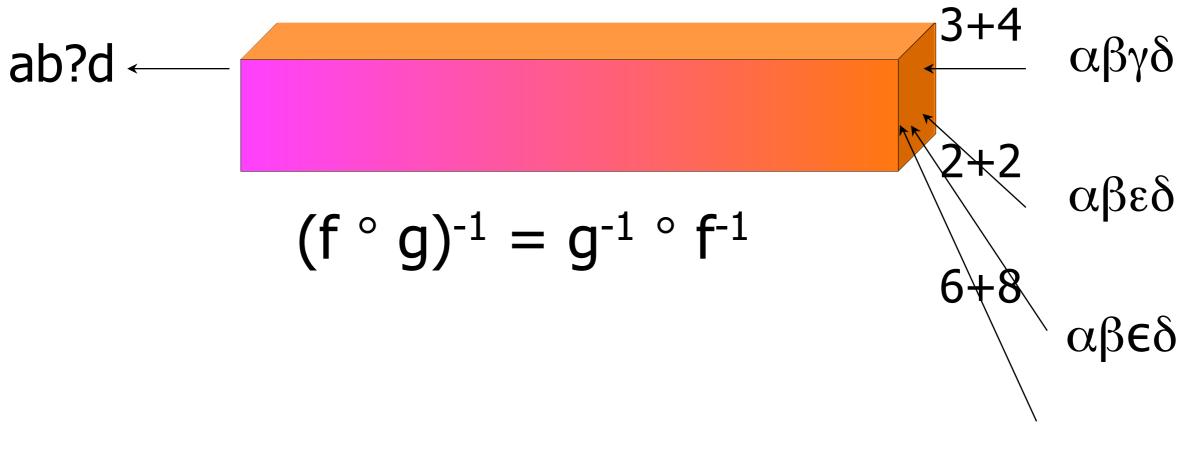


# **Inverting Relations**



. . .

# **Inverting Relations**

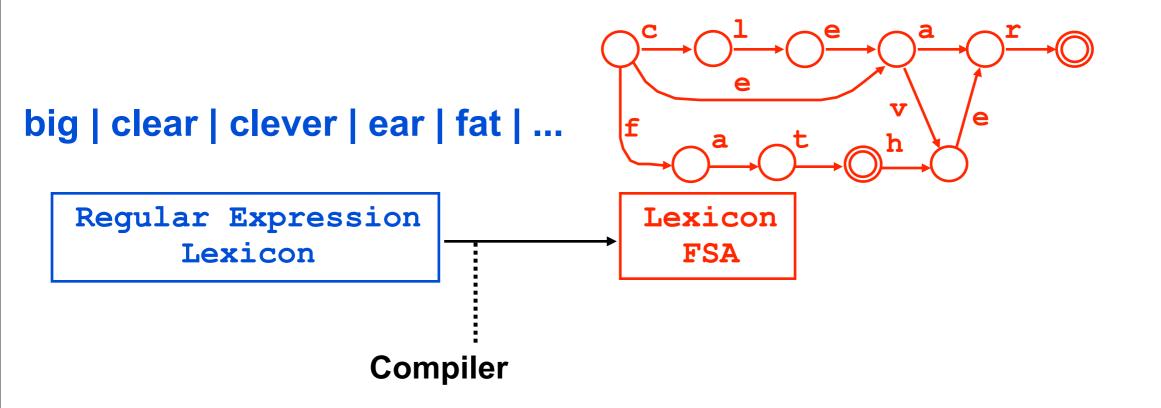


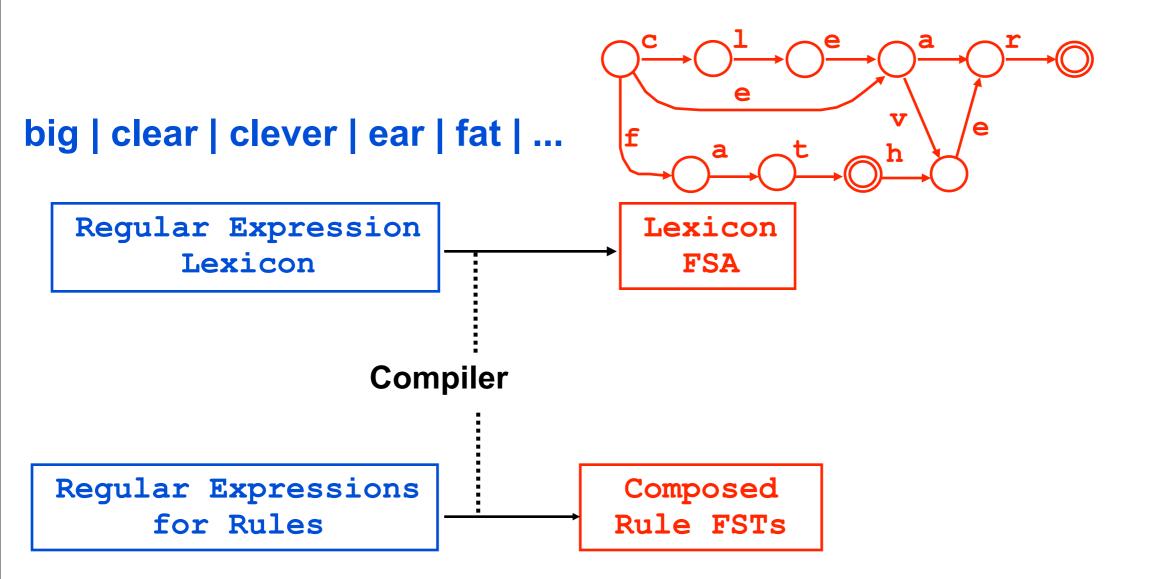
. . .

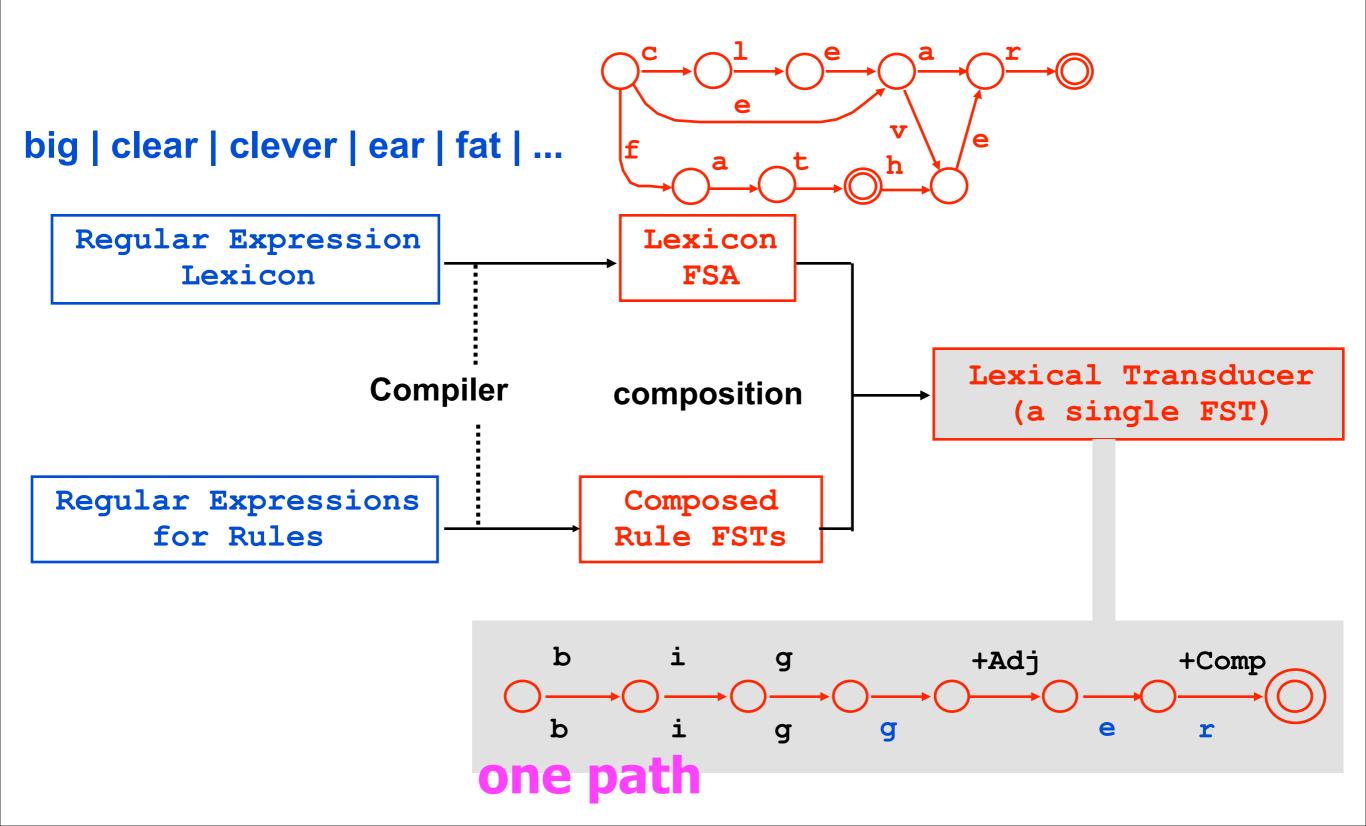
# **Building a lexical transducer**

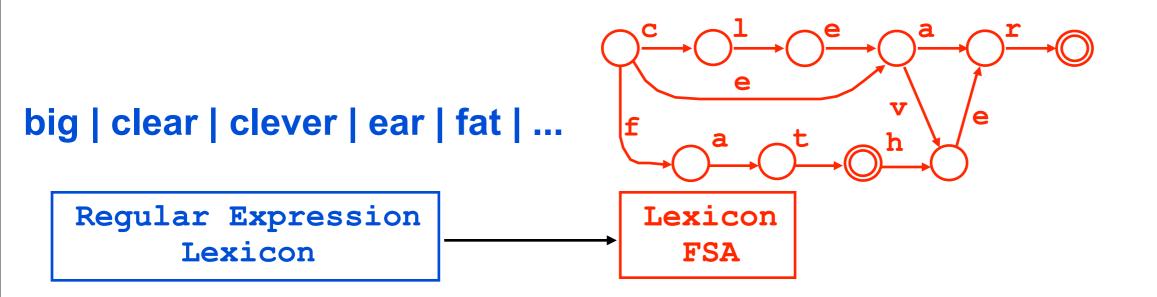
big | clear | clever | ear | fat | ...

Regular Expression Lexicon

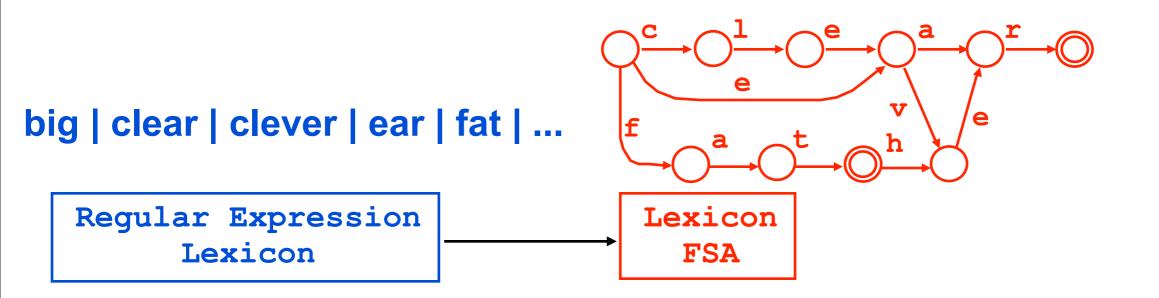




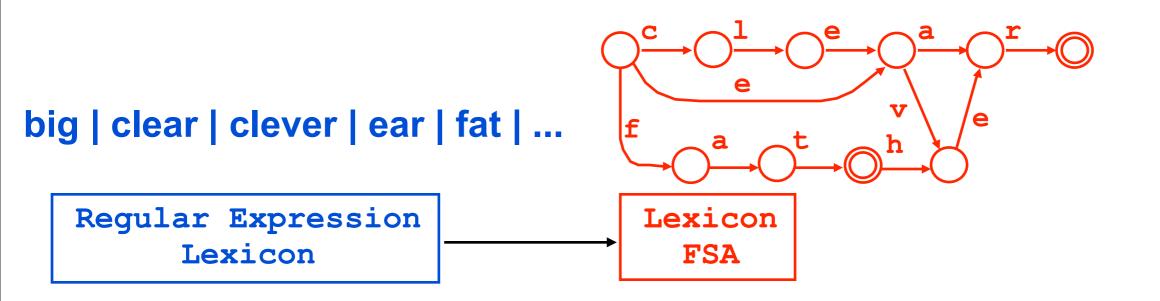




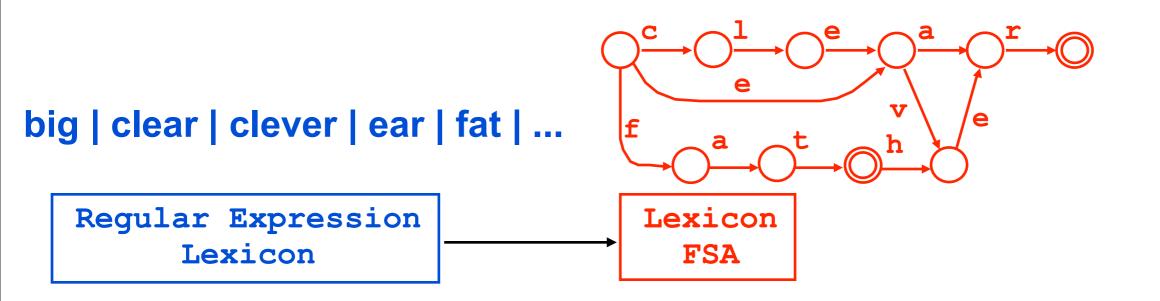
# **Building a lexical transducer**



 Actually, the lexicon must contain elements like big +Adj +Comp

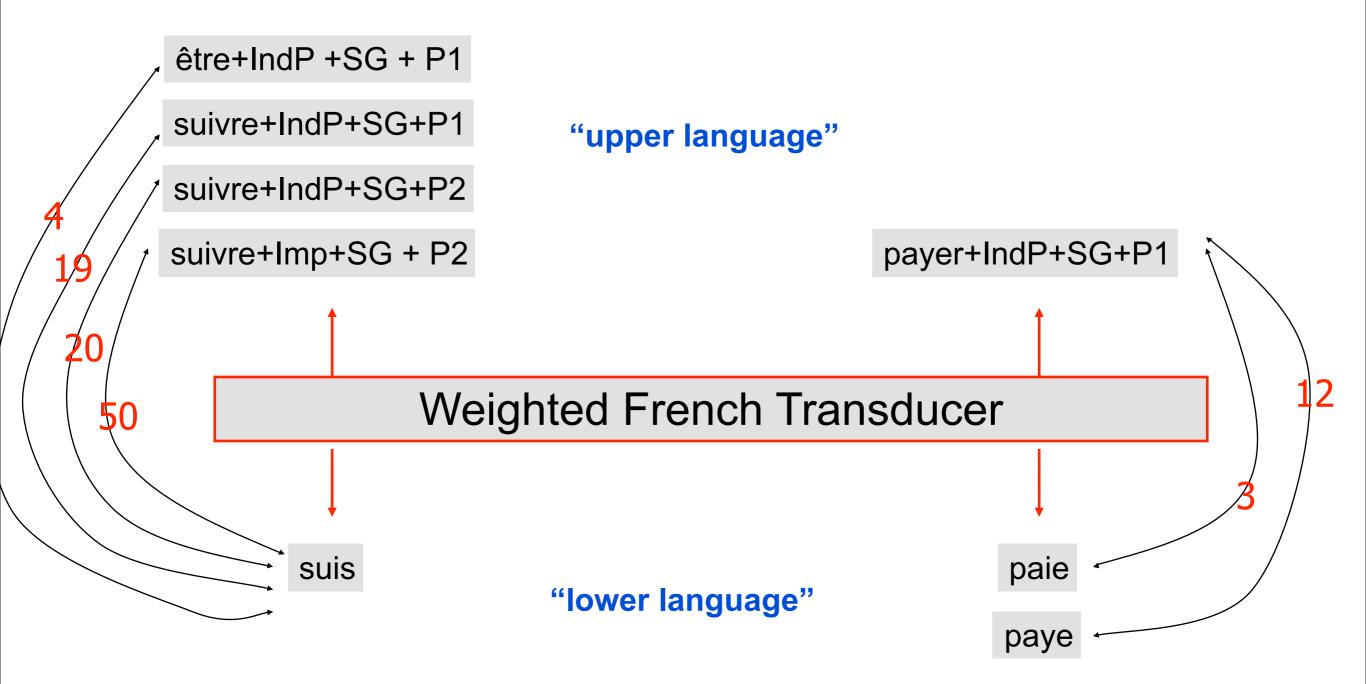


- Actually, the lexicon must contain elements like big +Adj +Comp
- So write it as a more complicated expression: (big | clear | clever | fat | ...) +Adj (ε | +Comp | +Sup) ← adjectives | (ear | father | ...) +Noun (+Sing | +PI) ← nouns | ...



- Actually, the lexicon must contain elements like big +Adj +Comp
- So write it as a more complicated expression: (big | clear | clever | fat | ...) +Adj (ε | +Comp | +Sup) ← adjectives | (ear | father | ...) +Noun (+Sing | +PI) ← nouns | ...
- Q: Why do we need a lexicon at all?

# Weighted version of transducer: Assigns a weight to each string pair



# Constructing Regular Languages

# **Xerox Regex Notation (Paper)**

- concatenationEF\* +iterationE\*, E+IunionE | F
- intersection
- complementation, minus
- .x. crossproduct
- .o. composition
- . u upper (input) language
- . lower (output) language
- E\*, E+ E & F ~E, \x, F-E E .x. F **E** .o. **F** E.u "domain" E. "range"

# Common Regular Expression Operators (in XFST notation)

#### concatenation



#### $\mathsf{EF} = \{\mathsf{ef}: \mathsf{e} \in \mathsf{E}, \mathsf{f} \in \mathsf{F}\}$

ef denotes the concatenation of 2 strings.

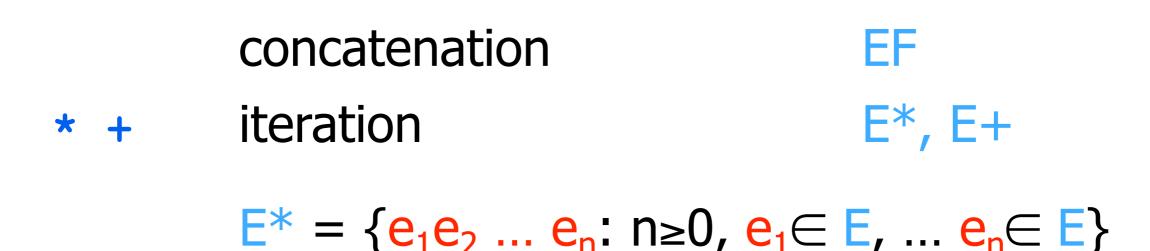
EF denotes the concatenation of 2 languages.

- To pick a string in EF, pick  $e \in E$  and  $f \in F$  and concatenate them.
- To find out whether  $w \in EF$ , look for at least one way to split w into two "halves," w = ef, such that  $e \in E$  and  $f \in F$ .

#### A **language** is a set of strings.

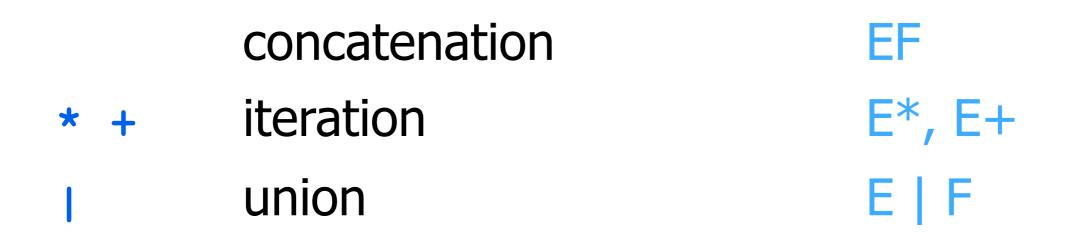
- It is a **regular language** if there exists an FSA that accepts all the strings in the language, and no other strings.
- If E and F denote regular languages, than so does EF. (We will have to prove this by finding the FSA for EF!)

# Common Regular Expression Operators (in XFST notation)



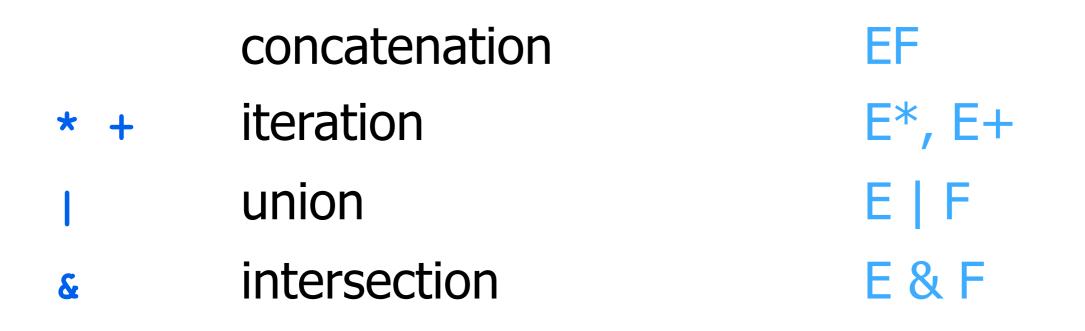
- To pick a string in E\*, pick any number of strings in E and concatenate them.
- To find out whether w ∈ E\*, look for at least one way to split w into 0 or more sections, e<sub>1</sub>e<sub>2</sub> ... e<sub>n</sub>, all of which are in E.

$$E + = \{e_1 e_2 \dots e_n : n > 0, e_1 \in E, \dots e_n \in E\} = EE^*$$



### $E \mid F = \{w: w \in E \text{ or } w \in F\} = E \cup F$

- To pick a string in E | F, pick a string from either E or F.
- To find out whether  $w \in E \mid F$ , check whether  $w \in E$  or  $w \in F$ .



### $E \& F = \{w: w \in E \text{ and } w \in F\} = E \cap F$

To pick a string in E & F, pick a string from E that is also in F.

• To find out whether  $w \in E \& F$ , check whether  $w \in E$  and  $w \in F$ .

	concatenation	EF
* +	iteration	E*, E+
	union	E F
&	intersection	E & F
~ \ -	complementation, minus	~E, \x, F-E

 $\sim E = \{e: e \notin E\} = \Sigma^* - E$  $E - F = \{e: e \in E \text{ and } e \notin F\} = E \& \sim F$  $E = \Sigma - E \quad (any single character not in E)$  $\Sigma \text{ is set of all letters; so } \Sigma^* \text{ is set of all strings; } ?* in XFST$ 

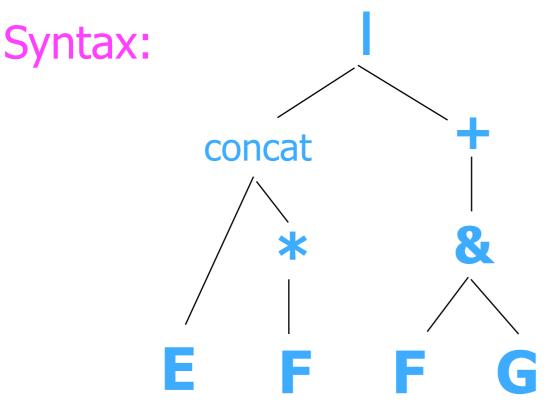
# **Regular Expressions**

### A language is a set of strings.

It is a **regular language** if there exists an FSA that accepts all the strings in the language, and no other strings.

If E and F denote regular languages, than so do EF, etc.

Regular expression: EF\* | (F & G)+



#### Semantics:

Denotes a regular language. As usual, can build semantics compositionally bottom-up. E, F, G must be regular languages. As a base case, e denotes {e} (a language containing a single string), so ef\*|(f&g)+ is regular.

# Regular Expressions for Regular Relations

A language is a set of strings.

It is a **regular language** if there exists an FSA that accepts all the strings in the language, and no other strings.

If E and F denote regular languages, than so do EF, etc.

A **relation** is a set of pairs – here, pairs of strings.

It is a **regular relation** if here exists an FST that accepts all the pairs in the language, and no other pairs.

If E and F denote regular <u>relations</u>, then so do EF, etc.

 $EF = \{(ef,e'f'): (e,e') \in E, (f,f') \in F\}$ 

Can you guess the definitions for E\*, E+, E | F, E & F when E and F are regular relations?

Surprise: E & F isn't necessarily regular in the case of relations; so not supported.

	concatenation	EF
* +	iteration	E*, E+
I	union	E F
&	intersection	E & F
~ \ -	complementation, minus	~E, \x, F-E
.x.	crossproduct	E .x. F

### $\mathsf{E} . \mathsf{x}. \mathsf{F} = \{(\mathsf{e},\mathsf{f}): \mathsf{e} \in \mathsf{E}, \mathsf{f} \in \mathsf{F}\}\$

Combines two regular <u>languages</u> into a regular <u>relation</u>.

	concatenation	EF
* +	iteration	E*, E+
I	union	E F
&	intersection	E & F
~ \ -	complementation, minus	~E, \x, F-E
. <b>x</b> .	crossproduct	E .x. F
.0.	composition	E .o. F

### E .o. $F = \{(e,f): \exists m. (e,m) \in E, (m,f) \in F\}$

- Composes two regular <u>relations</u> into a regular <u>relation</u>.
- As we've seen, this generalizes ordinary function composition.

	concatenation	EF
* +	iteration	E*, E+
I	union	E F
&	intersection	E & F
~ \ -	complementation, minus	~E, \x, F-E
.x.	crossproduct	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
	E.u = {e: ∃m. (e,m) ∈ E}	

	concatenation	EF
* +	iteration	E*, E+
	union	E F
<b>&amp;</b>	intersection	E & F
~ \ -	complementation, minus	~E, \x, F-E
. X .	crossproduct	E .x. F
. 0 .	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

# Finite-State Programming

# Finite-state "programming"

Function

Source code Object code

Compiler Optimization of object code Function on (set of) strings

Regular expression Finite state machine

Regexp compiler Determinization, minimization, pruning

# Finite-state "programming"

Function composition	(Weighted) composition
Higher-order function	Operator
Function inversion (available in Prolog)	Function inversion
Structured programming	Ops + small regexps

# Finite-state "programming"

Parallelism Nondeterminism Stochasticity

Apply to set of strings

Nondeterminism

Prob.-weighted arcs

### Some Xerox Extensions

- \$ containment
- => restriction
- -> @-> replacement

Make it easier to describe complex languages and relations without extending the formal power of finite-state systems.

### \$[ab\*c]

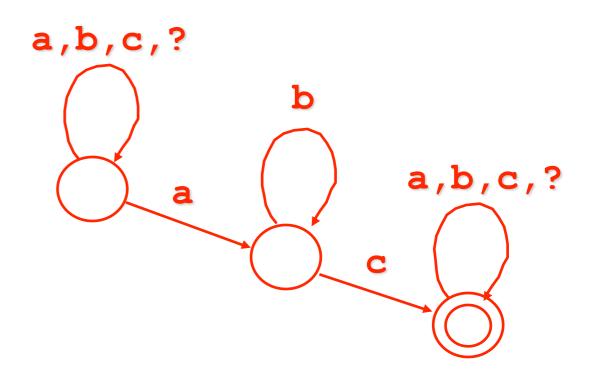
"Must contain a substring that matches **ab\*c**."

> Accepts **xxxacyy** Rejects **bcba**

\$[ab\*c]

"Must contain a substring that matches ab\*c."

> Accepts **xxxacyy** Rejects **bcba**



### \$[ab\*c]

"Must contain a substring that matches ab\*c."

> Accepts **xxxacyy** Rejects **bcba**

> > ?\* [ab\*c] ?\*

a,b,c,?

a

b

C

a,b,c,?

Equivalent expression

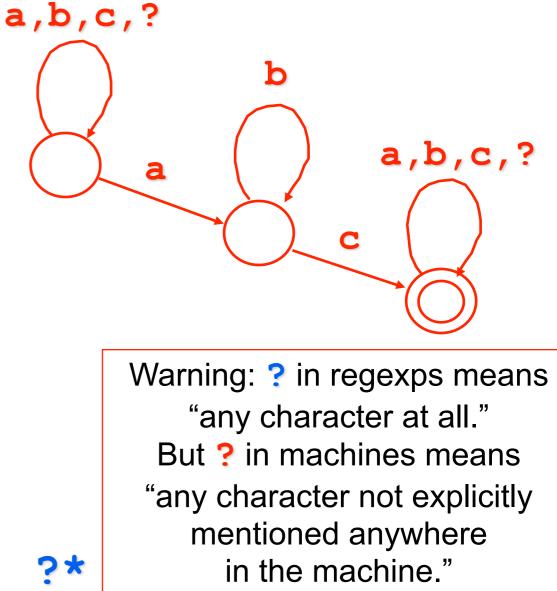
### \$[ab\*c]

"Must contain a substring that matches ab\*c."

> Accepts **xxxacyy** Rejects **bcba**

> > ?\* [ab\*c] ?\*

Equivalent expression



### Restriction

### Restriction

#### a => b \_ c

"Any a must be preceded by b and followed by c."

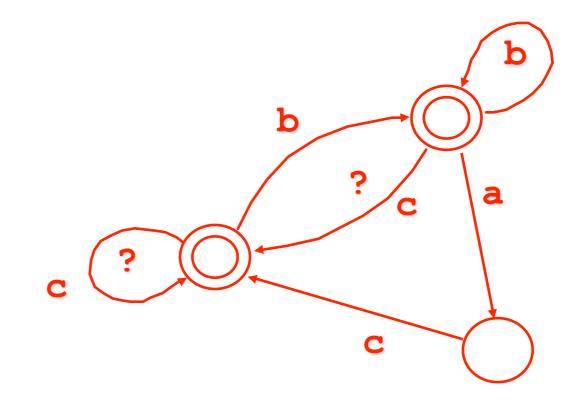
> Accepts **bacbbacde** Rejects **bac<u>a</u>**

### Restriction

a => b \_ c

"Any a must be preceded by b and followed by c."

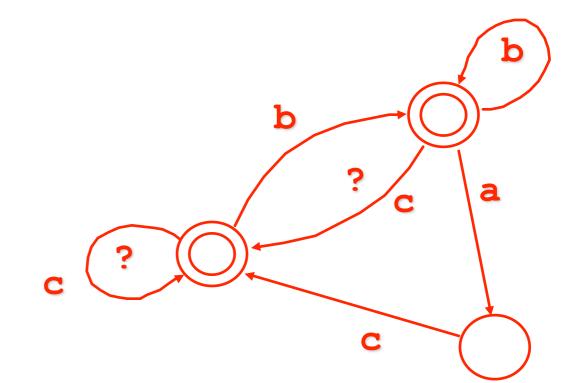
> Accepts **bacbbacde** Rejects **bac<u>a</u>**



### Restriction

a => b \_ c

"Any a must be preceded by b and followed by c."



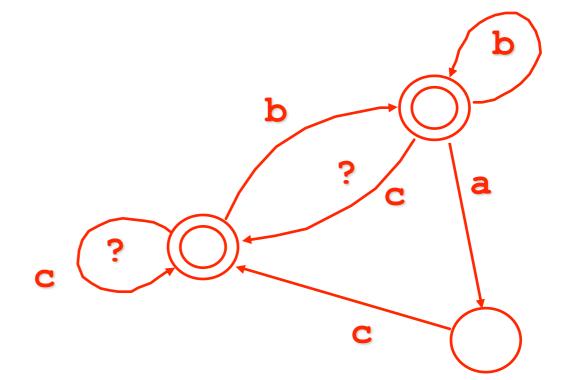
Accepts **bacbbacde** Rejects **bac<u>a</u>** 

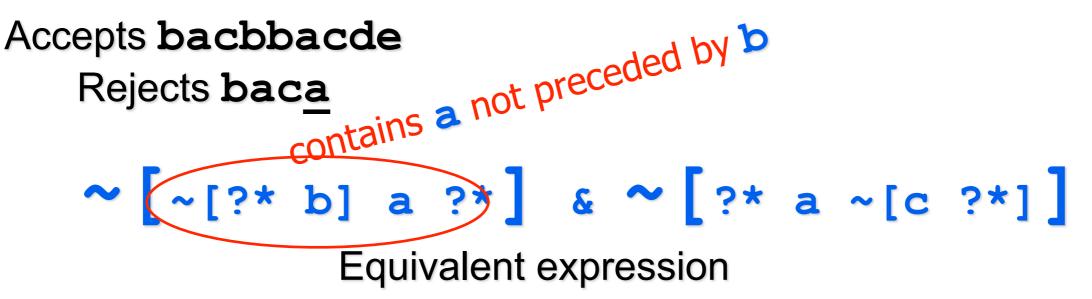
> ~ [~[?\* b] a ?\*] & ~ [?\* a ~[c ?\*]] Equivalent expression

### Restriction

a => b \_ c

"Any a must be preceded by b and followed by c."

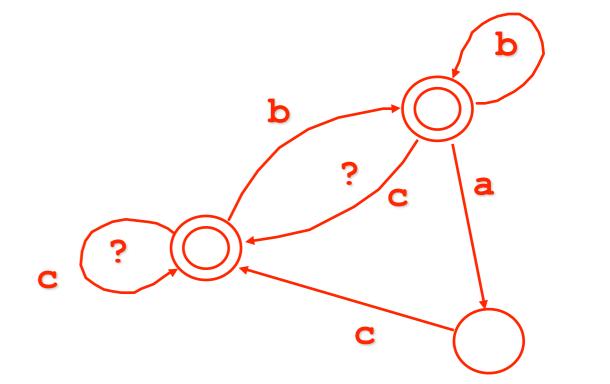


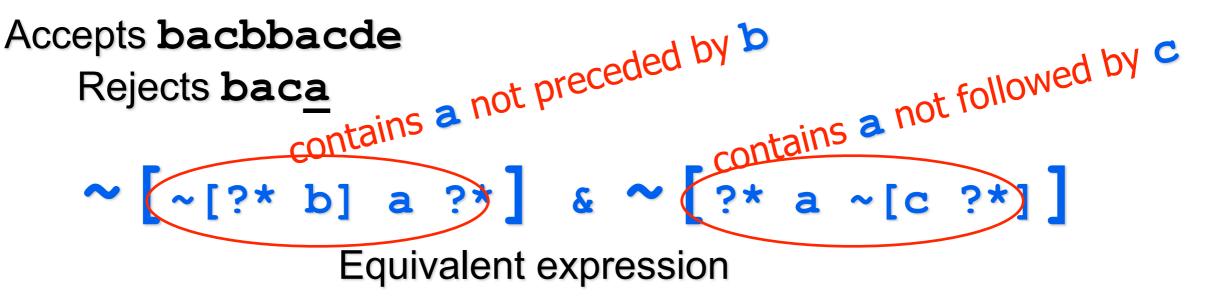


### Restriction

a => b \_ c

"Any a must be preceded by b and followed by c."





### Replacement

### Replacement

#### a b -> b a

"Replace 'ab' by 'ba'."

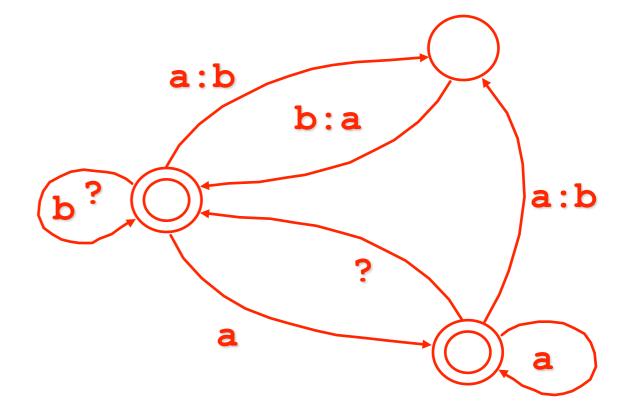
Transduces <u>abcdbaba</u> to <u>bacdbba</u>a

### Replacement

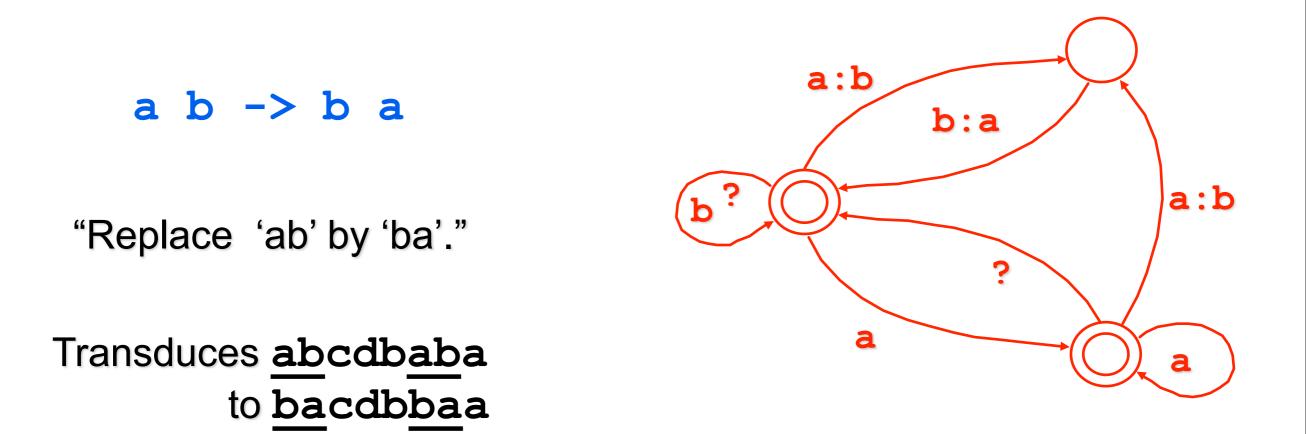


"Replace 'ab' by 'ba'."

Transduces <u>abcdbaba</u> to <u>bacdbba</u>a



### Replacement



[~\$[a b] [[a b] .x. [b a]]] \* ~\$[a b] Equivalent expression

#### a b -> b a | x

"Replace 'ab' by 'ba' or 'x', nondeterministically."

Transduces <u>abcdbaba</u> to {<u>bacdbbaa</u>, <u>bacdbxa</u>, <u>xcdbbaa</u>, <u>xcdbxa</u>}

#### [ a b -> b a | x ] .o. [ x => \_ c ]

"Replace 'ab' by 'ba' or 'x', nondeterministically."

Transduces <u>abcdbaba</u> to {<u>bacdbbaa</u>, <u>bacdbxa</u>, <u>x</u>cdb<u>ba</u>a, <u>x</u>cdb<u>x</u>a}

#### [ a b -> b a | x ] .o. [ x => \_ c ]

"Replace 'ab' by 'ba' or 'x', nondeterministically."

Transduces <u>abcdbaba</u> to {<u>bacdbbaa</u>, <u>bacdbxa</u>, <u>xcdbbaa</u>, <u>xcdbxa</u>}

#### a b | b | b a | a b a -> x

applied to "aba" Four overlapping substrings match; we haven't told it which one to replace so it chooses nondeterministically

_	ba aba —	a ]	b a	a b	a
	x a X	a	X	X	a

### More Replace Operators

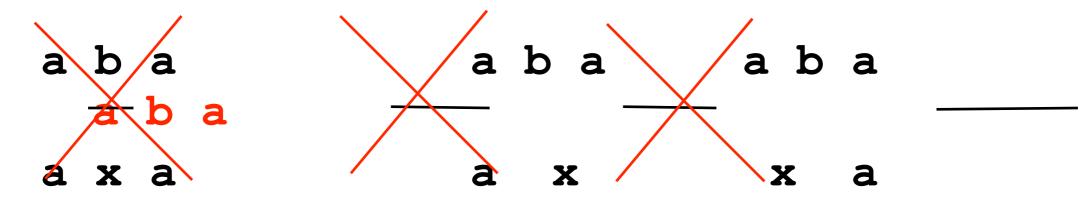
Optional replacement: a b (->) b a

- Directed replacement
  - guarantees a unique result by constraining the factorization of the input string by
    - Direction of the match (rightward or leftward)
    - Length (longest or shortest)

#### **@-> Left-to-right, Longest-match Replacement**

#### a b | b | b a | a b a @-> x

applied to "aba"

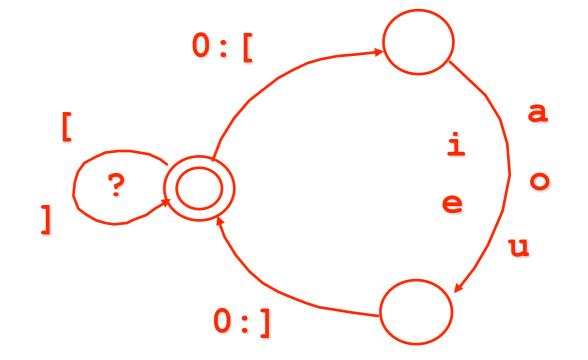


@-> left-to-right, longest match (cf. perl s///)

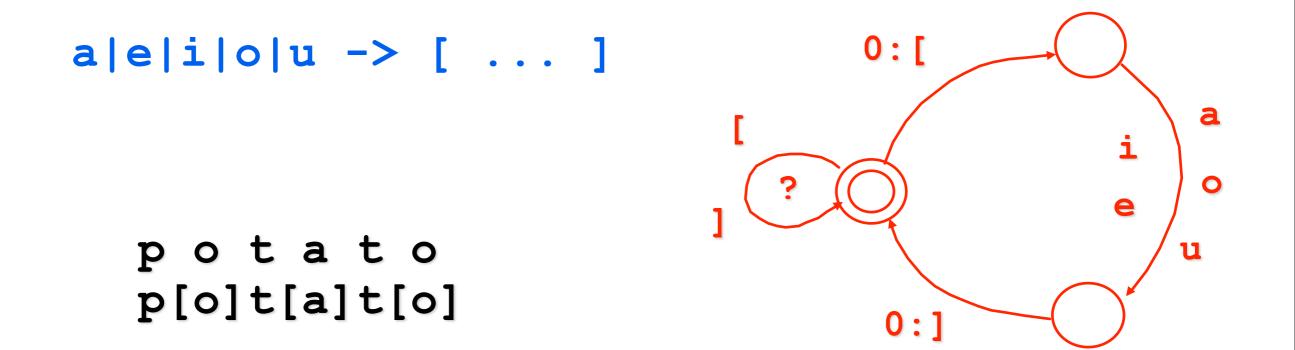
- @> left-to-right, shortest match
- ->@ right-to-left, longest match
- >@ right-to-left, shortest match

# Using "..." for marking

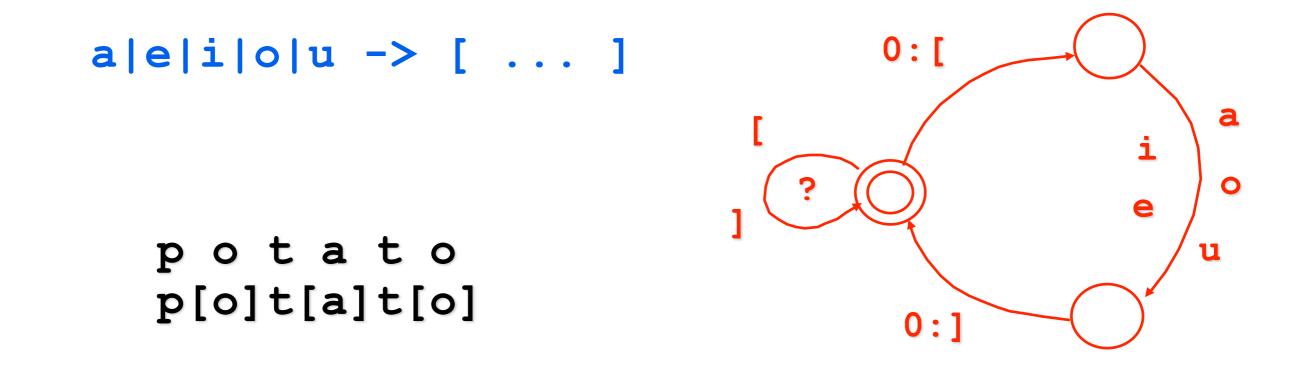
#### a|e|i|o|u -> [ ... ]



# Using "..." for marking

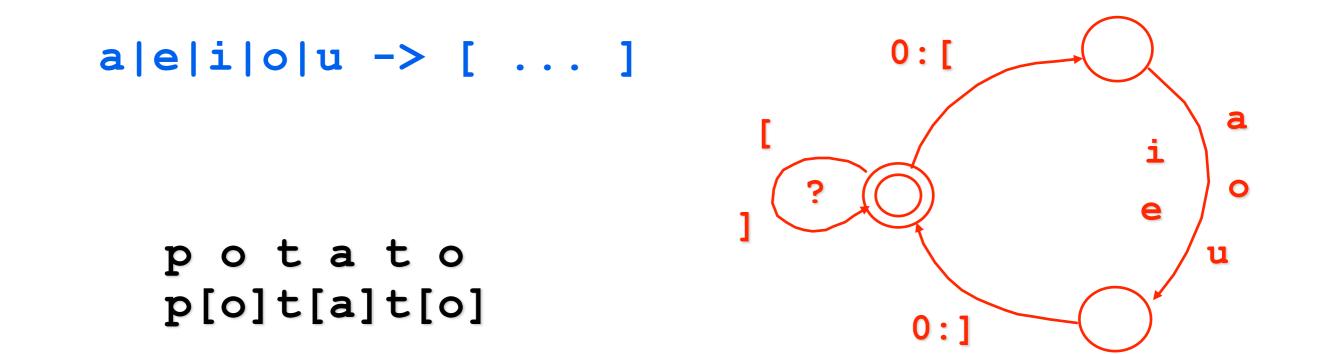


## Using "..." for marking



Note: actually have to write as -> %[ ... %] or -> "[" ... "]" since [] are parens in the regexp language

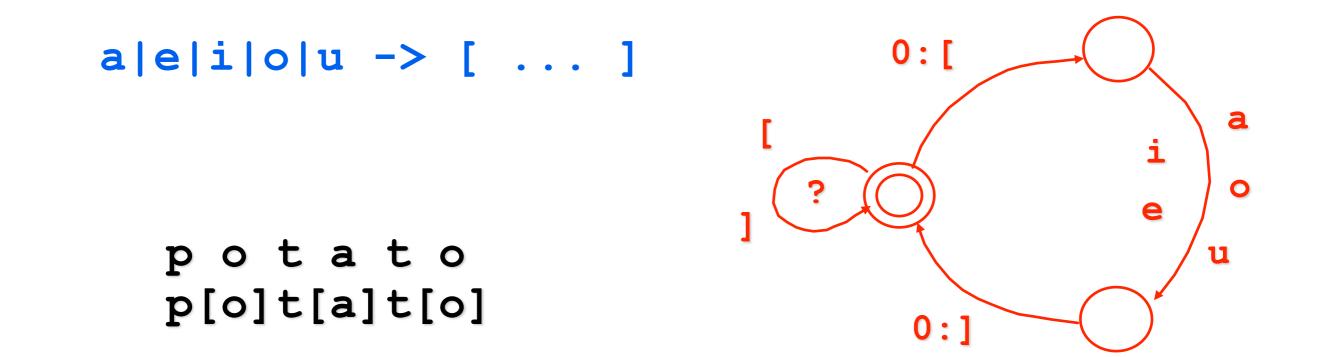
# Using "..." for marking



Which way does the FST transduce potatoe?

potato e vs. potato e
p[o]t[a]t[o][e] vs. p[o]t[a]t[o e]

# Using "..." for marking



Which way does the FST transduce potatoe?

potato e vs. potato e
p[o]t[a]t[o][e] vs. p[o]t[a]t[o e]

How would you change it to get the other answer?

define C [ b | c | d | f ... define V [ a | e | i | o | u | y | ä | ...

#### define C [ b | c | d | f ... define V [ a | e | i | o | u | y | ä | ...

#### $[C*V+C*] @-> \dots "-" || _ [CV]$

"Insert a hyphen after the longest instance of the C\* V+ C\* pattern in front of a C V pattern."

#### $[C*V+C*] @-> \dots "-" || [CV]$

"Insert a hyphen after the longest instance of the C\* V+ C\* pattern in front of a C V pattern."

define C [ b | c | d | f ... define V [ a | e | i | o | u | y | ä | ...

#### $[C*V+C*]@-> \dots "-"|| [CV]$

"Insert a hyphen after the longest instance of the C\* V+ C\* pattern in front of a C V pattern." why?

# **Conditional Replacement**

# **Conditional Replacement**



The relation that replaces **A** by **B** between **L** and **R** leaving everything else unchanged.

# **Conditional Replacement**



The relation that replaces **A** by **B** between **L** and **R** leaving everything else unchanged.

Sources of complexity:

- Replacements and contexts may overlap
- Alternative ways of interpreting "between left and right."

# Hand-Coded Example: slide courtesy of L. Karttunen Parsing Dates

Today is [Tuesday, July 25, 2000].

# Hand-Coded Example: slide courtesy of L. Karttunen Parsing Dates

Today is [Tuesday, July 25, 2000].

**Best result** 

Today is Tuesday, [July 25, 2000]. Today is [Tuesday, July 25], 2000. Today is Tuesday, [July 25], 2000. Today is [Tuesday], July 25, 2000.

**Bad results** 

# Hand-Coded Example: slide courtesy of L. Karttunen Parsing Dates

Today is [Tuesday, July 25, 2000].

**Best result** 

Today is Tuesday, [July 25, 2000]. Today is [Tuesday, July 25], 2000. Today is Tuesday, [July 25], 2000. Today is [Tuesday], July 25, 2000.

**Bad results** 

Need left-to-right, longest-match constraints.

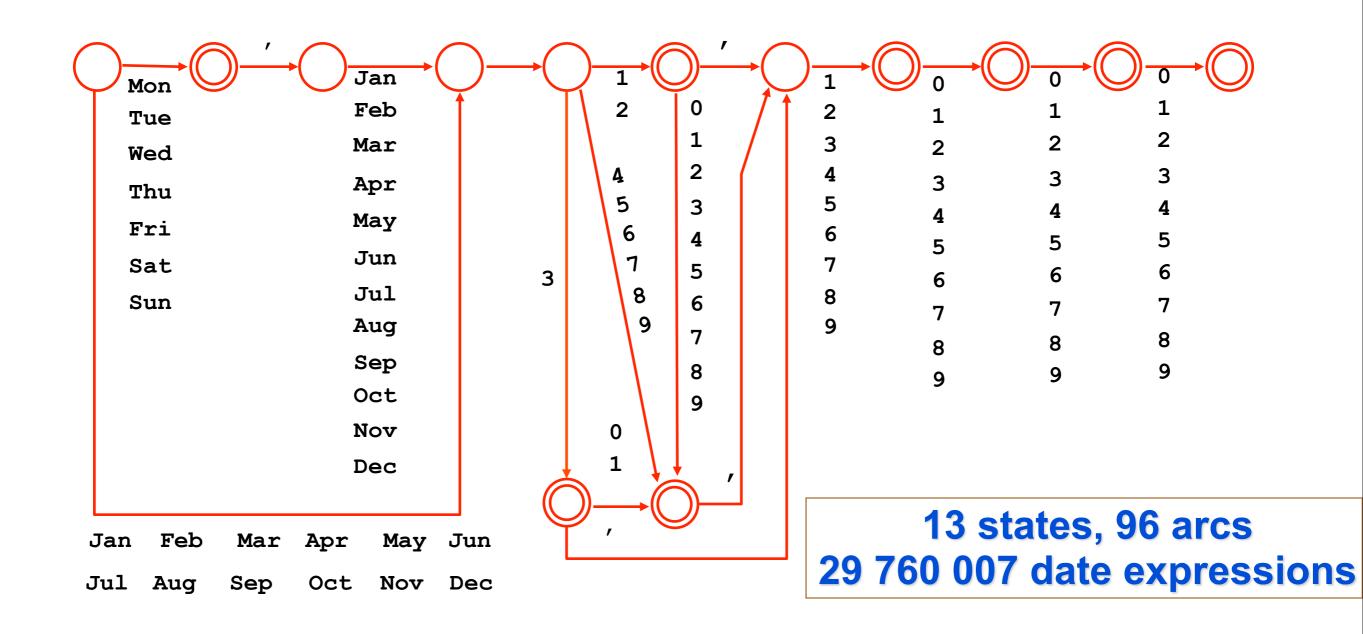
### Source code: Language of Dates

### Source code: Language of Dates

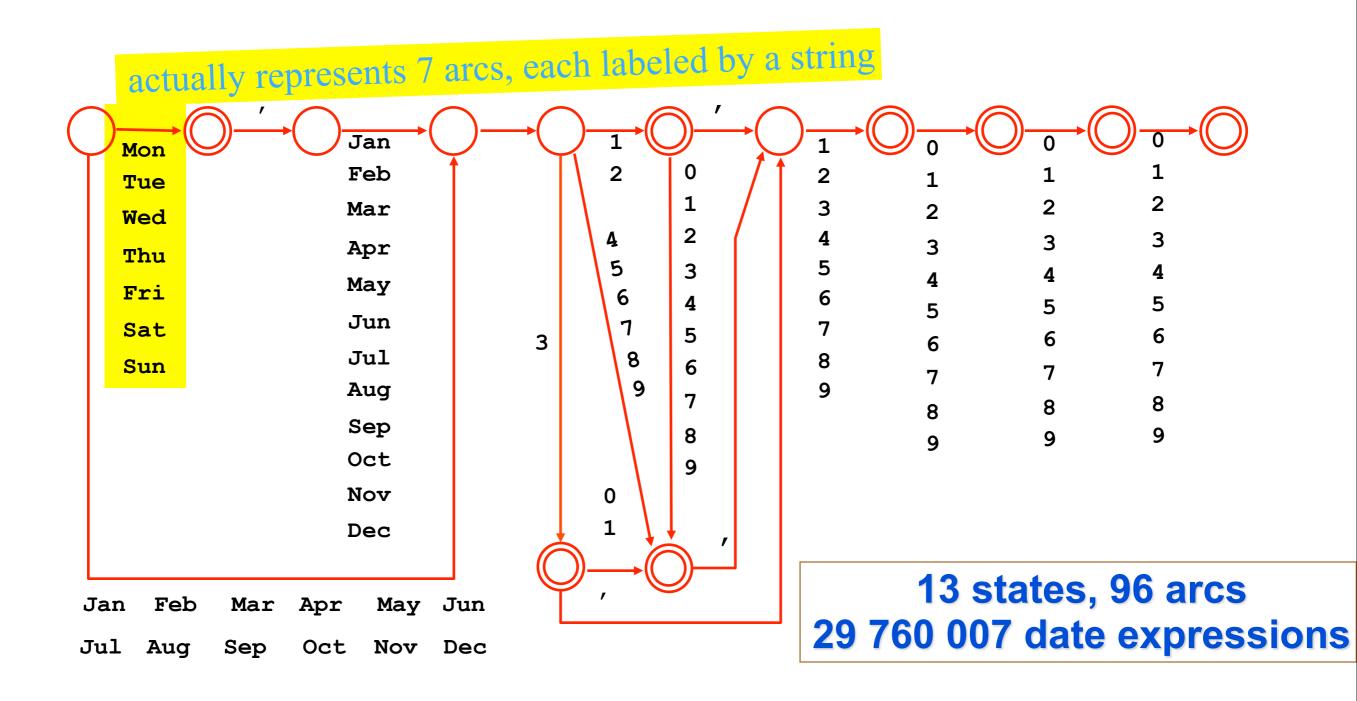
### Source code: Language of Dates

Day = Monday | Tuesday | ... | Sunday Month = January | February | ... | December Date = 1 | 2 | 3 | ... | 3 1 Year = %0To9 (%0To9 (%0To9 (%0To9))) - %0?\* from 1 to 9999

# Object code:slide courtesy of L. KarttunenAll Dates from I/I/I to I2/3I/9999



# Object code:slide courtesy of L. KarttunenAll Dates from I/I/I to I2/3I/9999



### **Parser for Dates**

### **Parser for Dates**

AllDates @-> "[DT " ... "]"

### **Parser for Dates**

AllDates @-> "[DT " ... "]"

Compiles into an unambiguous transducer (23 states, 332 arcs).

### **Parser for Dates**

AllDates @-> "[DT " ... "]"

Compiles into an unambiguous transducer (23 states, 332 arcs).

Today is [DT Tuesday, July 25, 2000] because yesterday was [DT Monday] and it was [DT July 24] so tomorrow must be [DT Wednesday, July 26] and not [DT July 27] as it says on the program.

### **Parser for Dates**

AllDates @-> "[DT " ... "]"

Compiles into an unambiguous transducer (23 states, 332 arcs).

Xerox left-to-right replacement operator

Today is [DT Tuesday, July 25, 2000] because yesterday was [DT Monday] and it was [DT July 24] so tomorrow must be [DT Wednesday, July 26] and not [DT July 27] as it says on the program.

## **Problem of Reference**

#### **Valid dates**

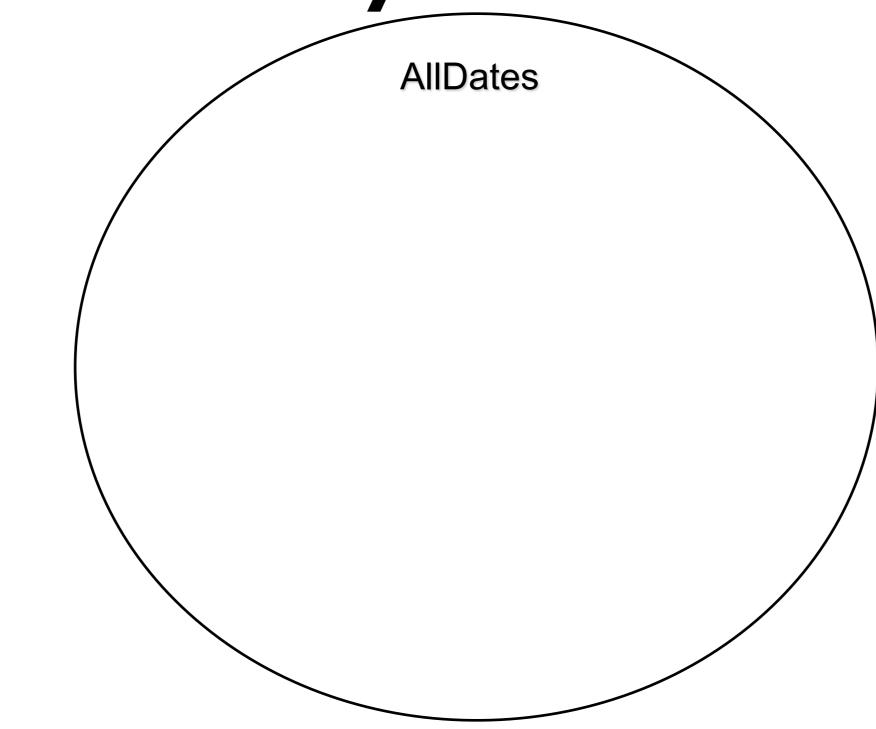
Tuesday, July 25, 2000

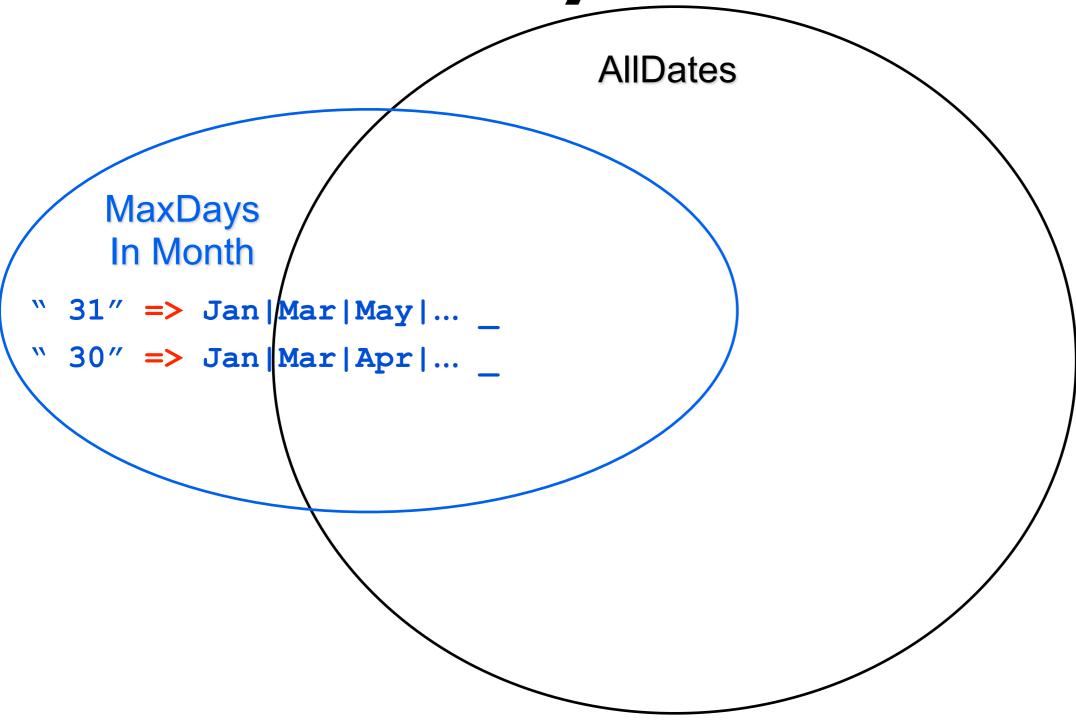
Tuesday, February 29, 2000

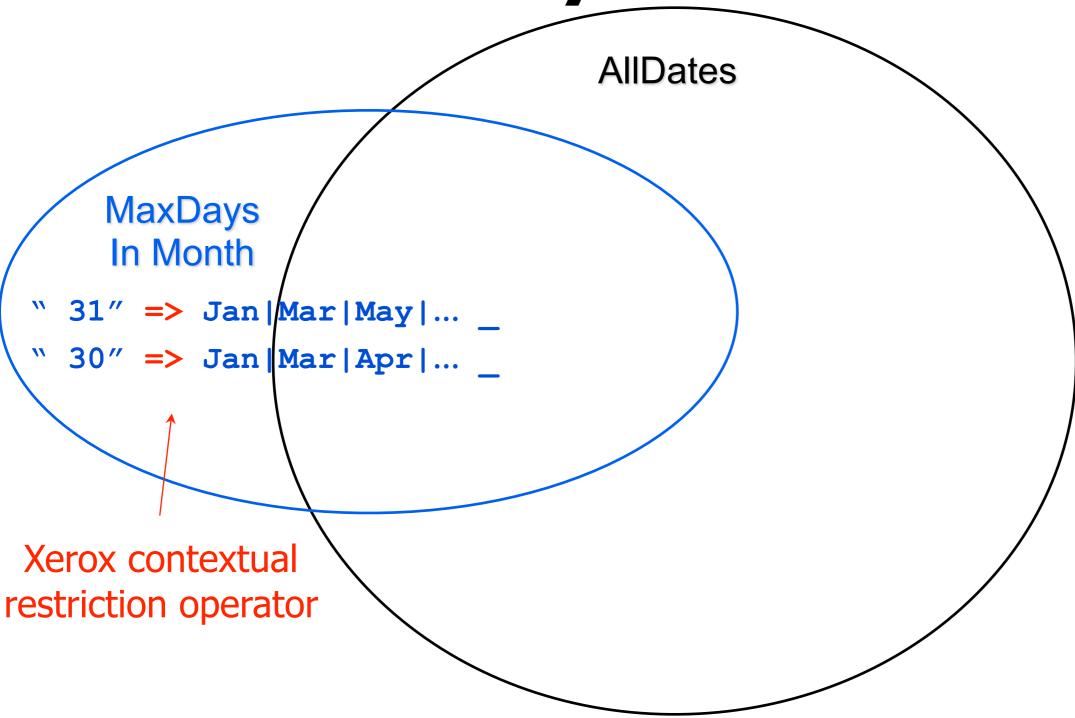
Monday, September 16, 1996

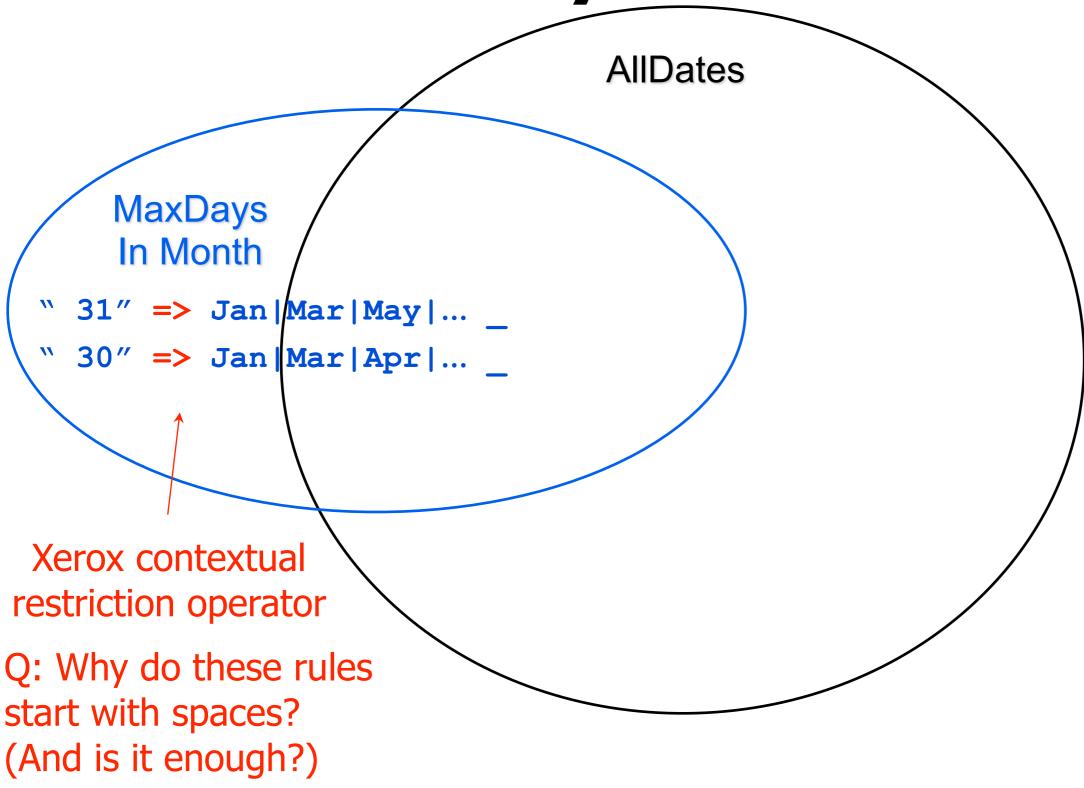
#### **Invalid dates**

Wednesday, April 31, 1996 Thursday, February 29, 1900 Tuesday, July 26, 2000

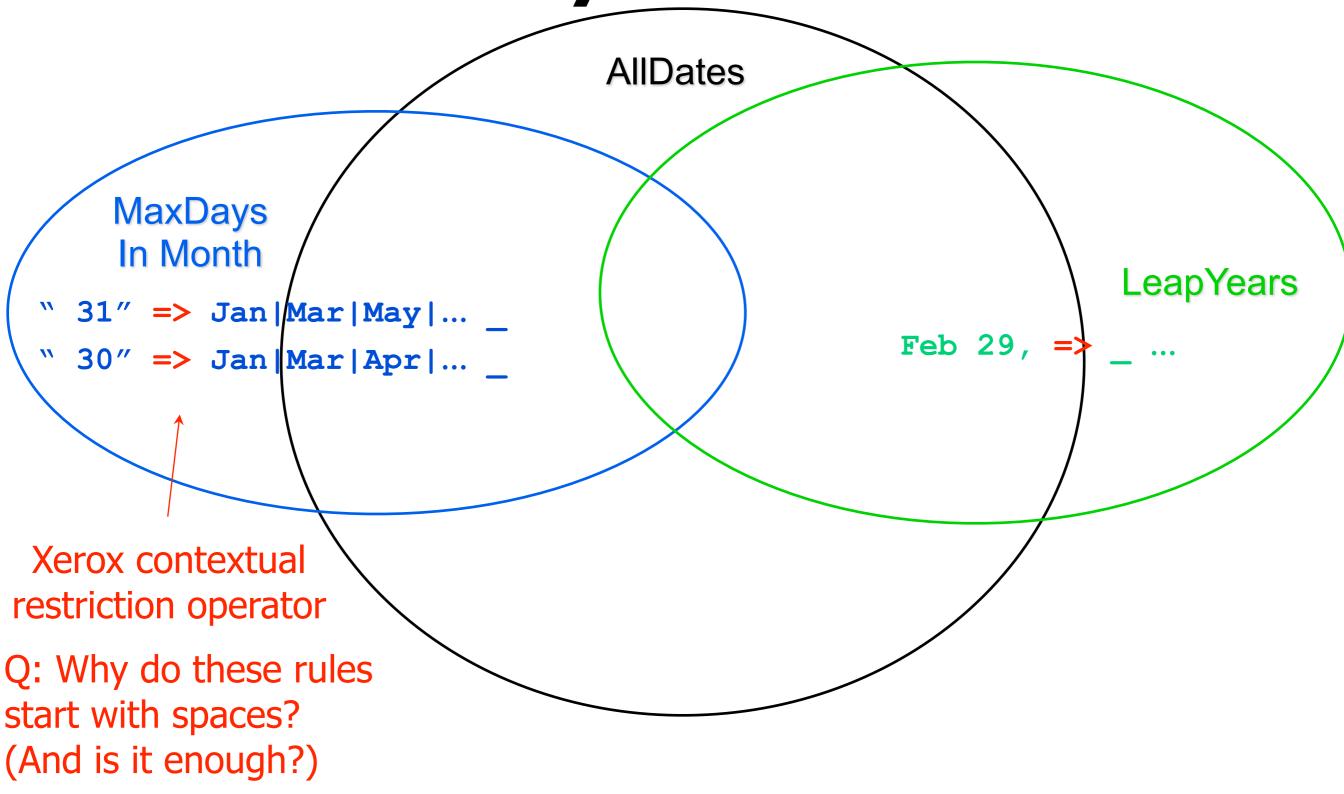




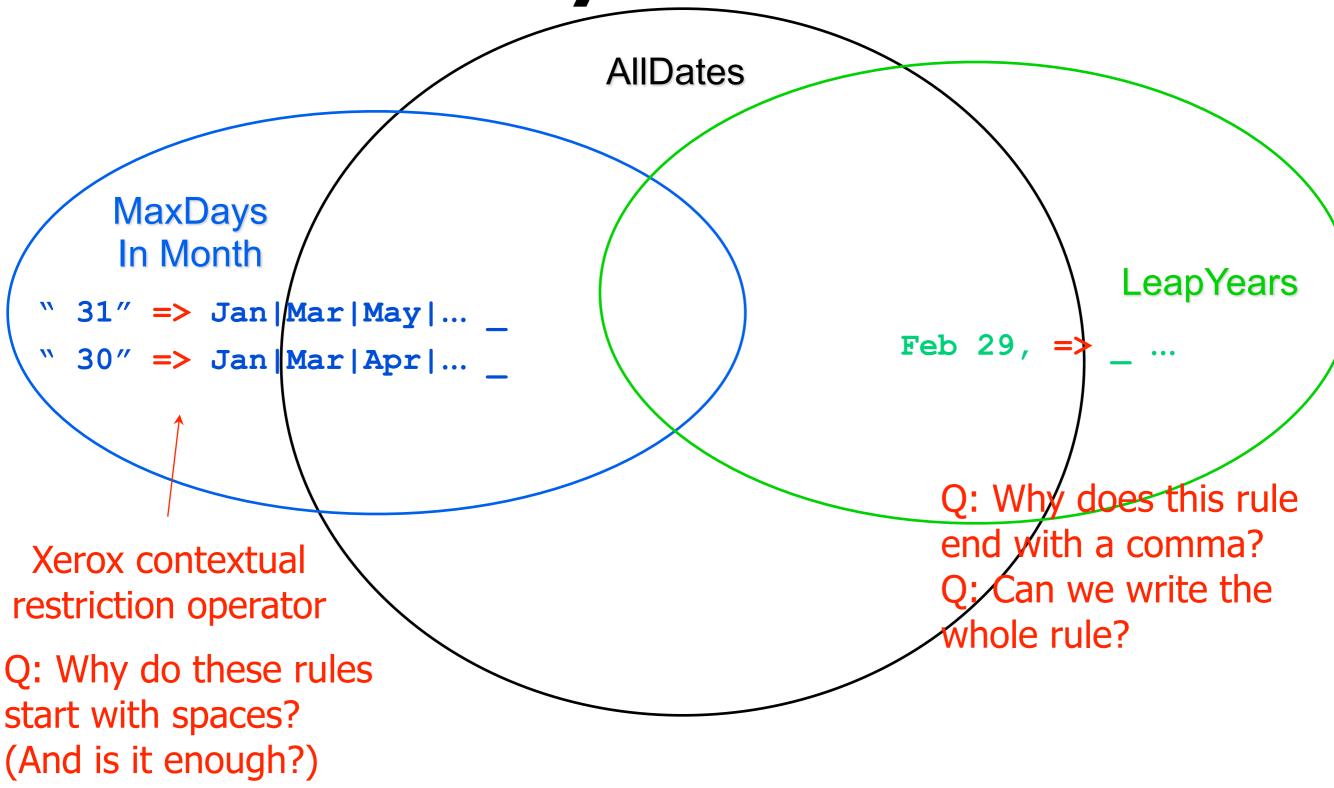


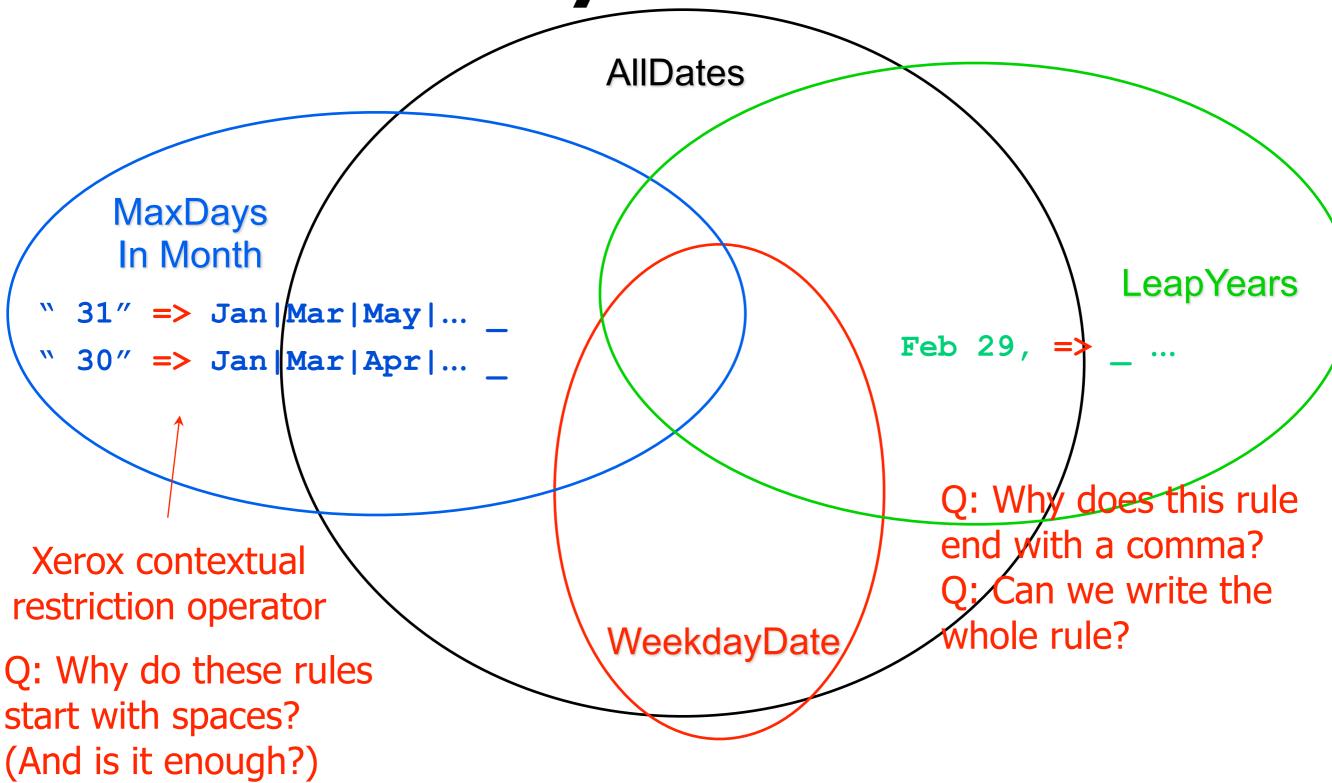


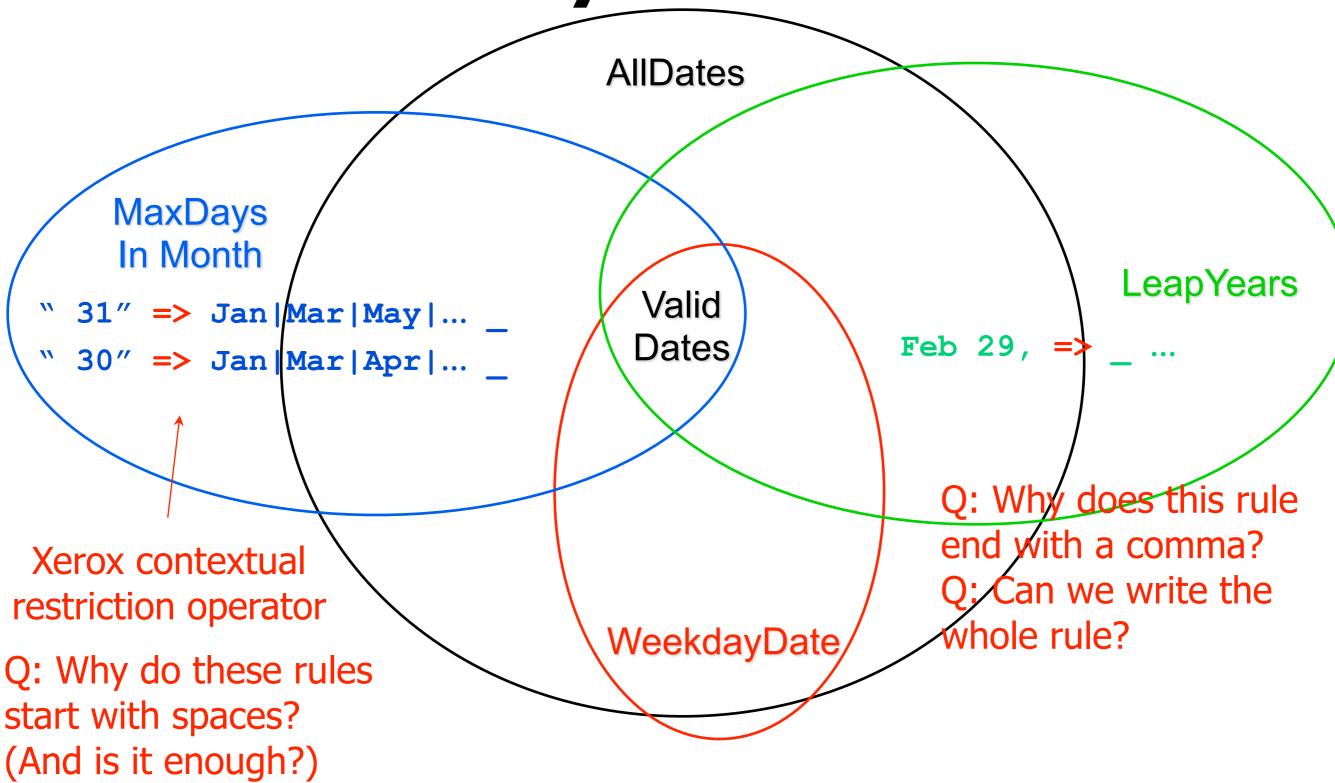
slide courtesy of L. Karttunen (modified)



slide courtesy of L. Karttunen (modified)







slide courtesy of L. Karttunen

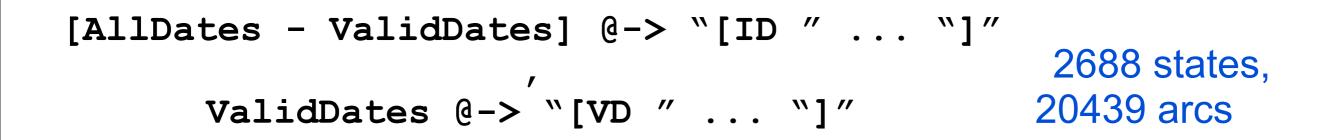
# **Defining Valid Dates**

AllDates: 13 states, 96 arcs 29 760 007 date expressions

= ValidDates

ValidDates: 805 states, 6472 arcs 7 307 053 date expressions

#### **Parser for Valid and Invalid Dates**



Today is [VD Tuesday, July 25, 2000], not [ID Tuesday, July 26, 2000]. valid date invalid date

#### **Parser for Valid and Invalid Dates**

Today is [VD Tuesday, July 25, 2000], not [ID Tuesday, July 26, 2000]. valid date

invalid date

Markup

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance
  - Phonology, morphology: series of little fixups? constraints?

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance
  - Phonology, morphology: series of little fixups? constraints?
  - Speech

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance
  - Phonology, morphology: series of little fixups? constraints?
  - Speech
  - Transliteration / back-transliteration

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance
  - Phonology, morphology: series of little fixups? constraints?
  - Speech
  - Transliteration / back-transliteration
  - Machine translation?

- Markup
  - Dates, names, places, noun phrases; spelling/grammar errors?
  - Hyphenation
  - Informative templates for information extraction (FASTUS)
  - Word segmentation (use probabilities!)
  - Part-of-speech tagging (use probabilities maybe!)
- Translation
  - Spelling correction / edit distance
  - Phonology, morphology: series of little fixups? constraints?
  - Speech
  - Transliteration / back-transliteration
  - Machine translation?
- Learning ...

#### **FASTUS – Information Extraction** Appelt et al, 1992-?

Input: Bridgestone Sports Co. said Friday it has set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be shipped to Japan. The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with ...

Output:	
Relationship:	TIE-UP
Entities:	"Bridgestone Sports Co."
	"A local concern"
	"A Japanese trading house"
Joint Venture Company:	"Bridgestone Sports Taiwan Co."
Amount:	NT\$2000000

#### **FASTUS: Successive Markups**

(details on subsequent slides)

Tokenization .0. **Multiwords** .0. Basic phrases (noun groups, verb groups ...) .0. Complex phrases .0. Semantic Patterns .0. Merging different references

### **FASTUS: Tokenization**

- Spaces, hyphens, etc.
- wouldn't  $\rightarrow$  would not
- their  $\rightarrow$  them 's
- company.  $\rightarrow$  company . but Co.  $\rightarrow$  Co.

# **FASTUS: Multiwords**

- •"set up"
- "joint venture"
- "San Francisco Symphony Orchestra,"
   "Canadian Opera Company"
- use a specialized regexp to match musical groups.
- -... what kind of regexp would match company names?

#### **FASTUS : Basic phrases**

Output looks like this (no nested brackets!): ... [NG it] [VG had set\_up] [NP a joint\_venture] [Prep in] ...

Company Name:	Bridgestone Sports Co.
Verb Group:	said
Noun Group:	Friday
Noun Group:	it
Verb Group:	had set up
Noun Group:	a joint venture
Preposition:	in
Location:	Taiwan
Preposition:	with
Noun Group:	a local concern

## **FASTUS: Noun Groups**

Build FSA to recognize phrases like approximately 5 kg more than 30 people the newly elected president the largest leftist political force a government and commercial project Use the FSA for left-to-right longest-match markup

What does FSA look like? See next slide ...

## **FASTUS: Noun Groups**

Described with a kind of non-recursive CFG ... (a regexp can include names that stand for other regexps)

NG → Pronoun | Time-NP | Date-NP NG → (Det) (Adjs) HeadNouns

Adjs → sequence of adjectives maybe with commas, conjunctions, adverbs

 $\mathsf{Det} \to \mathsf{DetNP} \mid \mathsf{DetNonNP}$ 

DetNP → detailed expression to match "the only five, another three, this, many, hers, all, the most ..."

. . .

#### **FASTUS: Semantic patterns**

```
BusinessRelationship =
NounGroup(Company/ies) VerbGroup(Set-up)
NounGroup(JointVenture) with NounGroup(Company/ies)
...
```

ProductionActivity =
 VerbGroup(Produce) NounGroup(Product)

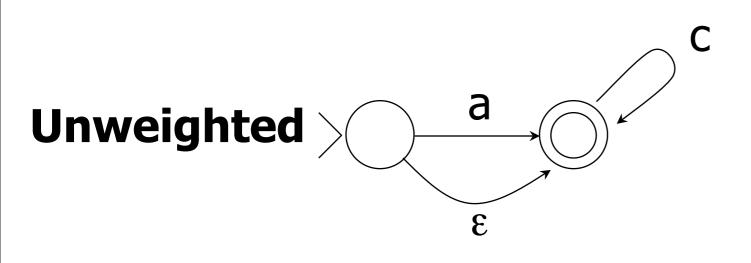
NounGroup(Company/ies) → NounGroup & ... is made easy by the processing done at a previous level

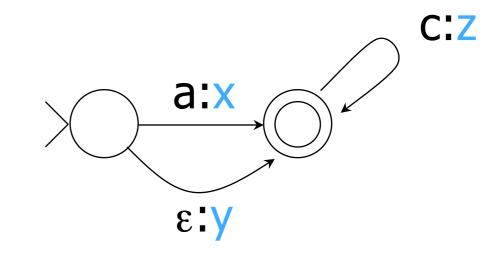
Use this for spotting references to put in the database.

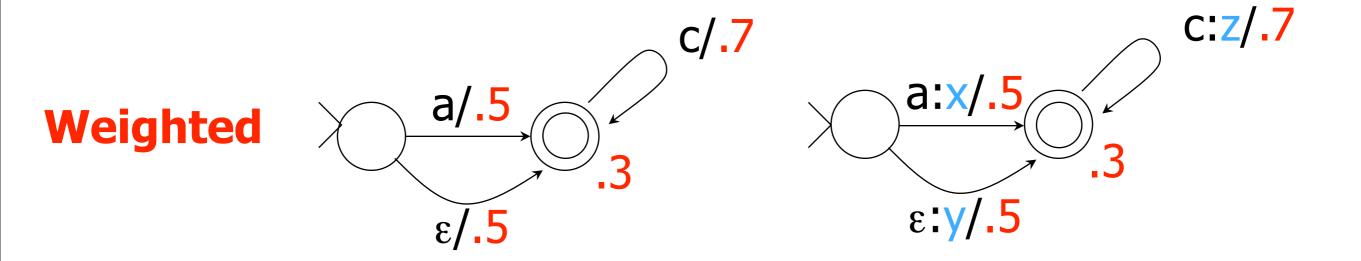
# Weighted FSMs

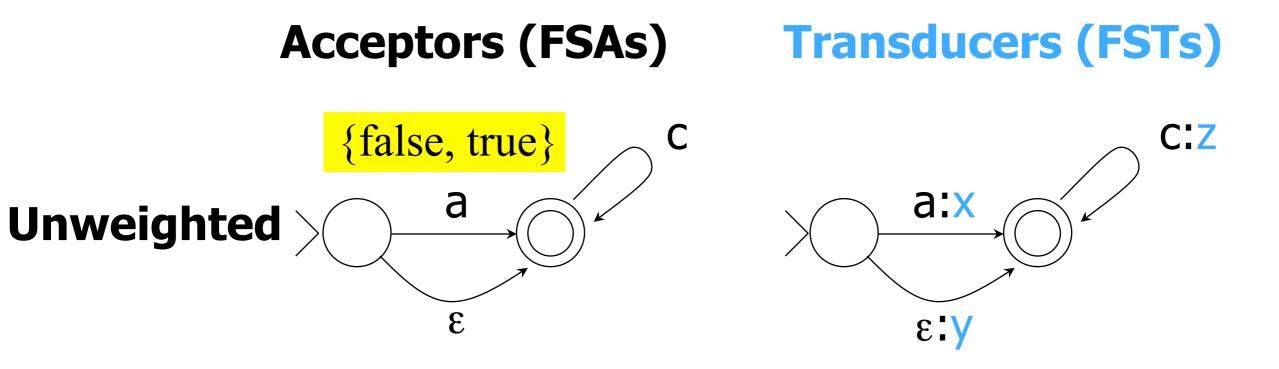
Acceptors (FSAs)

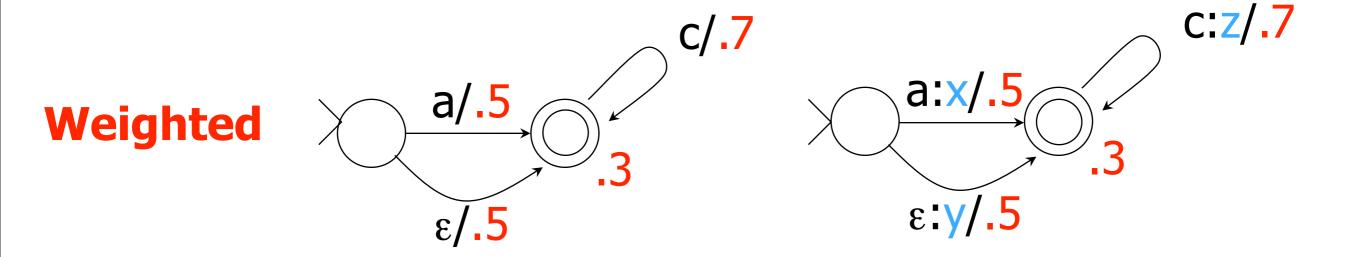
**Transducers (FSTs)** 

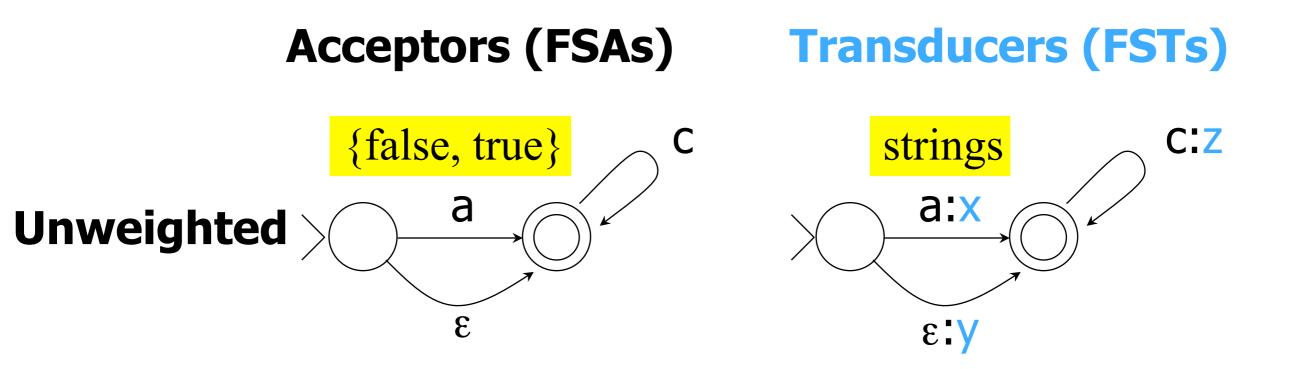


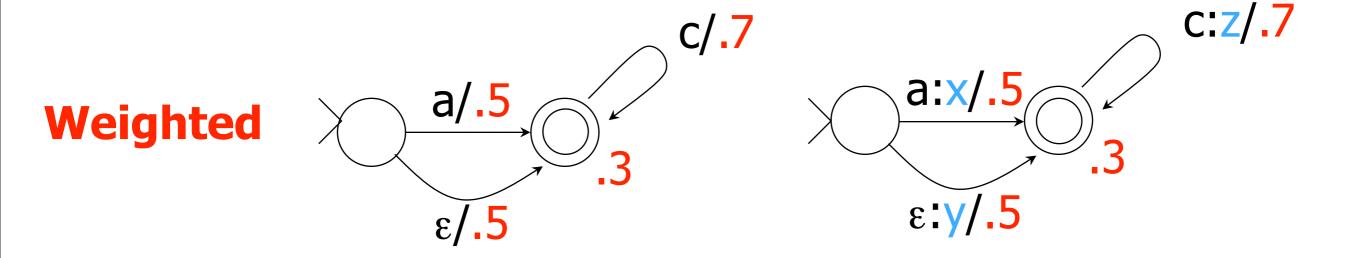


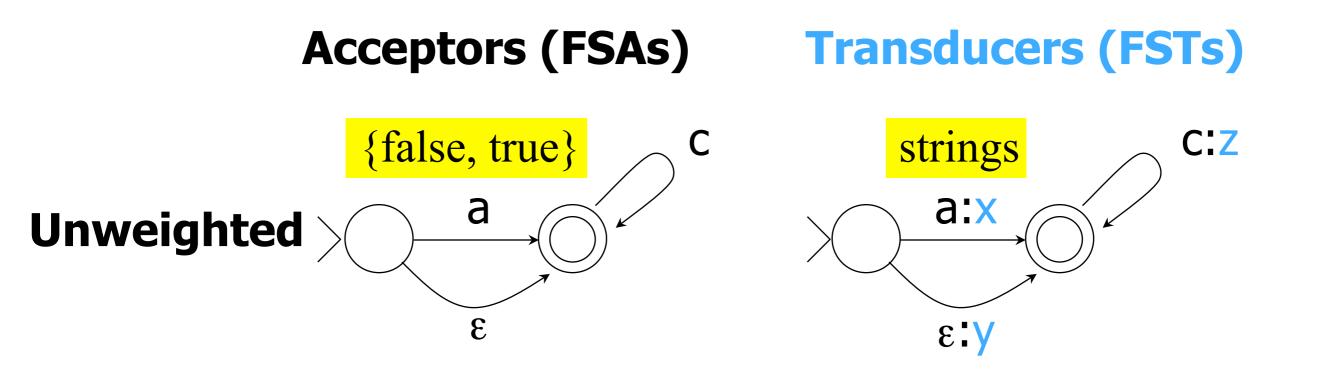


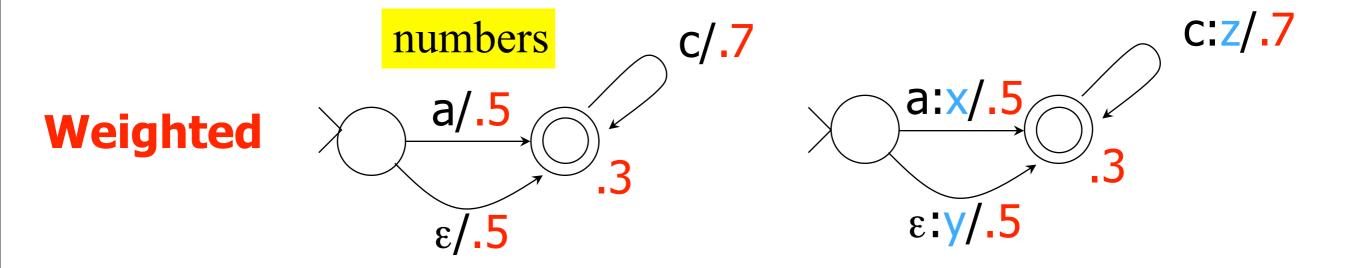


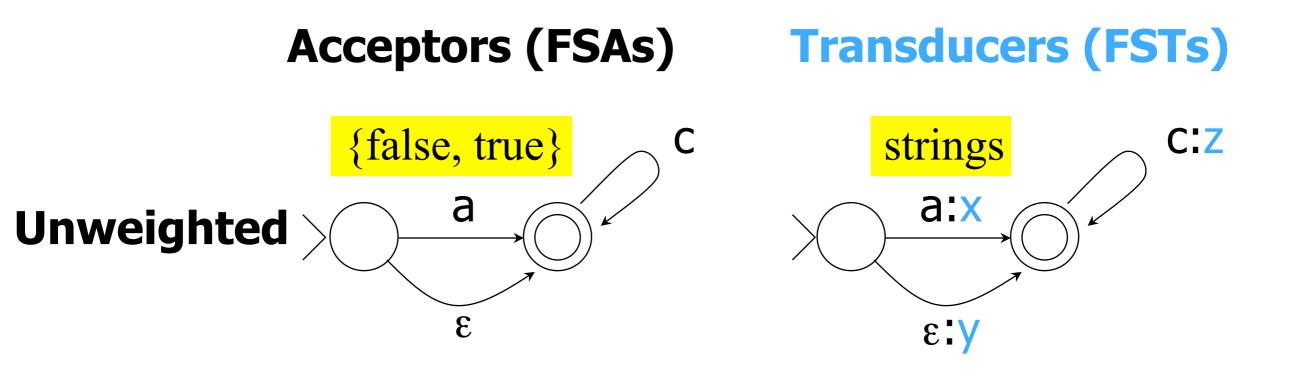


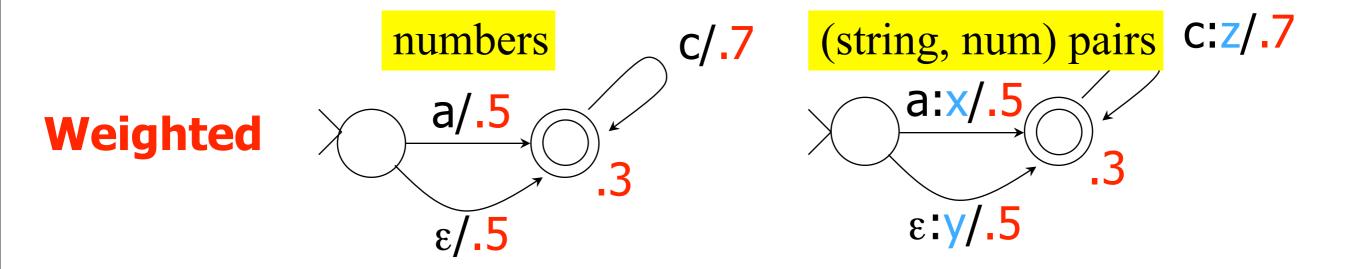












# Weighted Relations

- If we have a language [or relation], we can ask it: Do you contain this string [or string pair]?
- If we have a weighted language [or relation], we ask: What weight do you assign to this string [or string pair]?
- Pick a semiring: all our weights will be in that semiring.
  What?!

# Semirings

	Set	$\oplus$	$\otimes$	0	Ι
Prob	R	+	X	0	I
Max	R	max	X	0	Ι
Log	R∪{±∞}	log+	+	- 00	0
"Tropical"	R∪{±∞}	max	+	- 00	0
Shortest path	R∪{±∞}	min	+	$\infty$	0
Boolean	{0,I}	V	Λ	F	Т
String	Σ	longest common prefix	concat	$\infty$	3

# Weighted Relations

- If we have a language [or relation], we can ask it: Do you contain this string [or string pair]?
- If we have a weighted language [or relation], we ask: What weight do you assign to this string [or string pair]?
- Pick a semiring: all our weights will be in that semiring.
  - **Don't panic!** We will cover this again when we get to HMMs and parsing.
  - The unweighted case is the boolean semring {true, false}.
  - If a string is not in the language, it has weight 0.
  - If an FST or regular expression can choose among multiple ways to match, use 
     to combine the weights of the different choices.
  - If an FST or regular expression matches by matching multiple substrings, use is to combine those different matches.

#### Which Semiring Operators are Needed?

concatenationEF\* +iterationE\*, E+IunionE | F

- complementation, minus
- Intersection
- .x. crossproduct
- .o. composition
- . u upper (input) language
- . lower (output) language

~E, \x, E-F E & F E .x. F **E** .o. **F** E.u "domain" E. "range"

#### Which Semiring Operators are Needed?

	concatenation	EF	
* +	iteration	E*, E+	
	union $\oplus$ to sum over 2 choices	E F	
~ \ -	complementation, minus	~E, \x, E-F	
&	intersection	E & F	
. X .	crossproduct	E .x. F	
.0.	composition	E .o. F	
.u	upper (input) language	E.u "domain"	
.1	lower (output) language	E. "range"	

	concatenation	EF
* +	iteration	E*, E+
I	union $\oplus$ to sum over 2 choices	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
. x .	crossproduct against E and F	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

# Common Regular Expression Operators (in XFST notation)



### $\mathsf{E} \mid \mathsf{F} = \{\mathsf{w}: \mathsf{w} \in \mathsf{E} \text{ or } \mathsf{w} \in \mathsf{F}\} = \mathsf{E} \cup \mathsf{F}$

 Weighted case: Let's write E(w) to denote the weight of w in the weighted language E.

$$(\mathsf{E}|\mathsf{F})(\mathsf{w}) = \mathsf{E}(\mathsf{w}) \oplus \mathsf{F}(\mathsf{w})$$

	concatenation	EF
* +	iteration	E*, E+
	union $\oplus$ to sum over 2 choices	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
. X .	crossproduct	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

	concatenation	EF
* +	iteration	E*, E+
I	union $\oplus$ to sum over 2 choices	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
. x .	crossproduct against E and F	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

	concatenation need both	EF
* +	iteration	E*, E+
I	union	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
.x.	crossproduct J against E and F	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

concatenation iteration

+

 $\begin{array}{c} & \mathsf{EF} \\ \text{need both} \oplus \text{and} \\ \otimes & \mathsf{E*, E+} \end{array}$ 

# $\mathsf{EF} = \{\mathsf{ef}: \mathsf{e} \in \mathsf{E}, \mathsf{f} \in \mathsf{F}\}$

 Weighted case must match two things (⊗), but there's a choice (⊕) about which two things:

$$(EF)(W) = \bigcup_{e,f \text{ such}} (E(e) \otimes F(f))$$

that w=ef

	concatenation need both	EF
* +	iteration	E*, E+
I	union	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
.x.	crossproduct J against E and F	E .x. F
.0.	composition	E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

	concatenation iteration ⊗	
* +	iteration	E*, E+
I	union ⊕ to sum over 2 choices	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
.x.	crossproduct against E and F	E .x. F
.0.	composition both	<sup>?)</sup> E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E. "range"

	concatenation need both	EF
* +	iteration	E*, E+
	union	E F
~ \ -	complementation, minus	~E, \x, E-F
&	intersection	E & F
.x.	crossproduct against E and F	E .x. F
.0.	composition both	<sup>)</sup> E .o. F
.u	upper (input) language	E.u "domain"
.1	lower (output) language	E.I "range"

[Red material shows differences from FSAs.]

- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).

- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).
  - Each edge is also labeled with a "lower string" (in  $\Delta^*$ ).
  - [Upper/lower are sometimes regarded as input/output.]

- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).
  - Each edge is also labeled with a "lower string" (in  $\Delta^*$ ).
  - [Upper/lower are sometimes regarded as input/output.]
  - Each edge and final state is also labeled with a semiring weight.

- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).
  - Each edge is also labeled with a "lower string" (in  $\Delta^*$ ).
  - [Upper/lower are sometimes regarded as input/output.]
  - Each edge and final state is also labeled with a semiring weight.
- More traditional definition specifies an FST via these:
  - a state set Q
  - initial state i
  - set of final states F
  - input alphabet  $\Sigma$  (also define  $\Sigma^*$ ,  $\Sigma^+$ ,  $\Sigma^?$ )
  - output alphabet  $\Delta$

- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).
  - Each edge is also labeled with a "lower string" (in  $\Delta^*$ ).
  - [Upper/lower are sometimes regarded as input/output.]
  - Each edge and final state is also labeled with a semiring weight.
- More traditional definition specifies an FST via these:
  - a state set Q
  - initial state i
  - set of final states F
  - input alphabet  $\Sigma$  (also define  $\Sigma^*$ ,  $\Sigma^+$ ,  $\Sigma$ ?)
  - output alphabet  $\Delta$
  - transition function d: Q x  $\Sigma$ ? --> 2<sup>Q</sup>

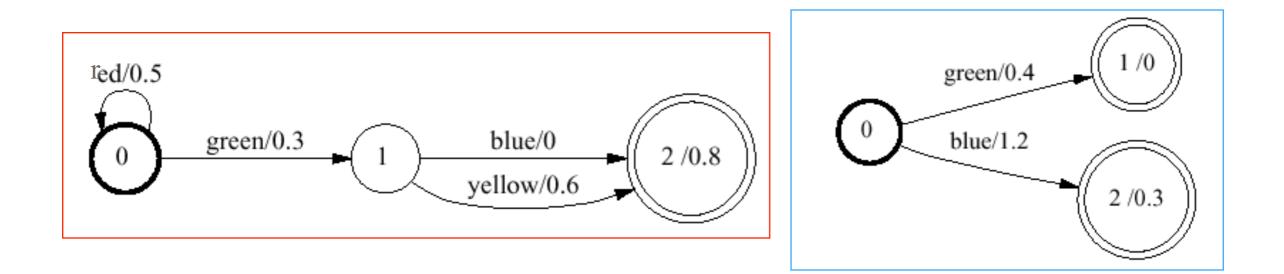
- [Red material shows differences from FSAs.]
- Simple view:
  - An FST is simply a finite directed graph, with some labels.
  - It has a designated initial state and a set of final states.
  - Each edge is labeled with an "upper string" (in  $\Sigma^*$ ).
  - Each edge is also labeled with a "lower string" (in  $\Delta^*$ ).
  - [Upper/lower are sometimes regarded as input/output.]
  - Each edge and final state is also labeled with a semiring weight.
- More traditional definition specifies an FST via these:
  - a state set Q
  - initial state i
  - set of final states F
  - input alphabet  $\Sigma$  (also define  $\Sigma^*$ ,  $\Sigma^+$ ,  $\Sigma^?$ )
  - output alphabet  $\Delta$
  - transition function d: Q x  $\Sigma$ ? --> 2<sup>Q</sup>
  - output function s: Q x  $\Sigma$ ? x Q -->  $\Delta$ ?

#### slide courtesy of L. Karttunen (modified)

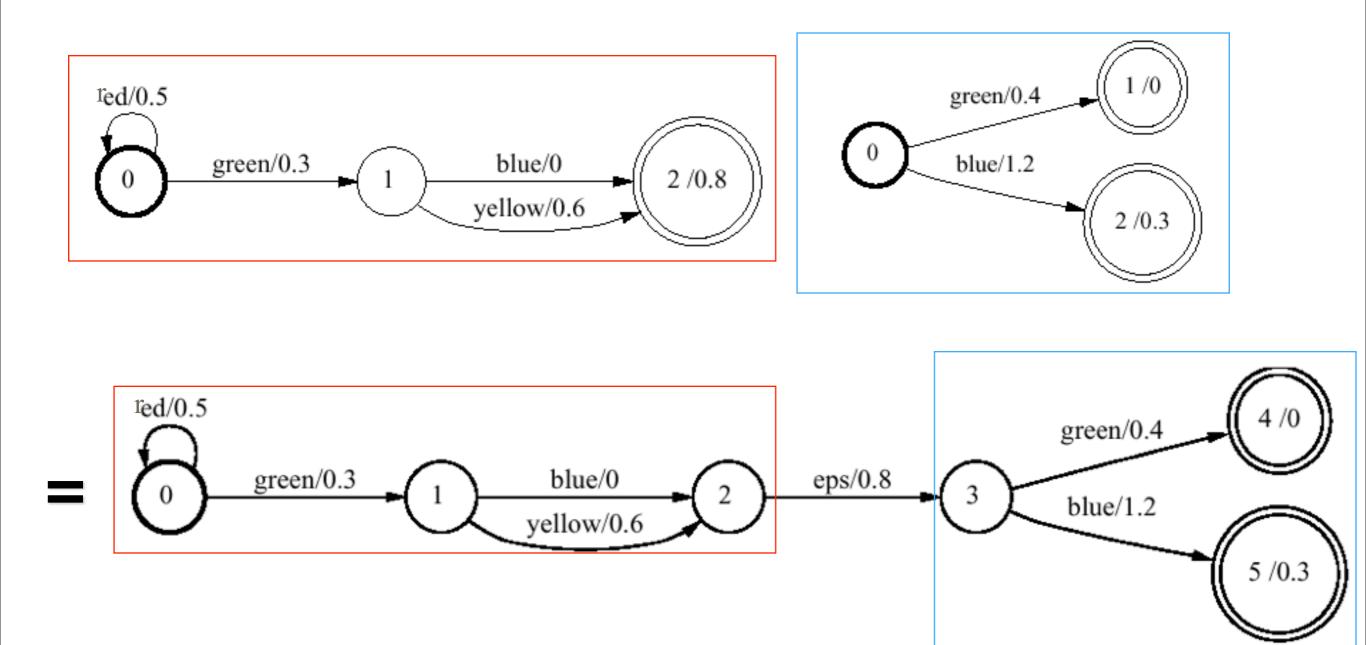
# How to implement?

- + iteration
   E\*, E+
   Union
   E | F
- complementation, minus
- intersection
- .x. crossproduct
- .o. composition
- . u upper (input) language
- . lower (output) language
- EF ~E, \x, E-F E & F E .x. F **E**.o. **F** E.u "domain" E. "range"

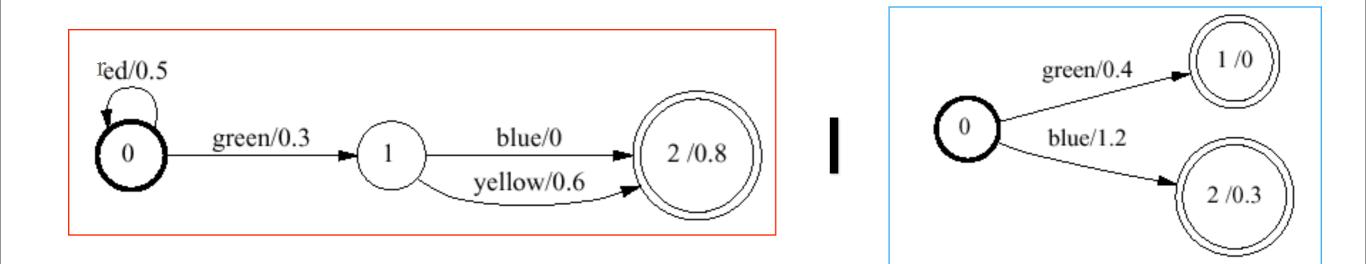
## Concatenation



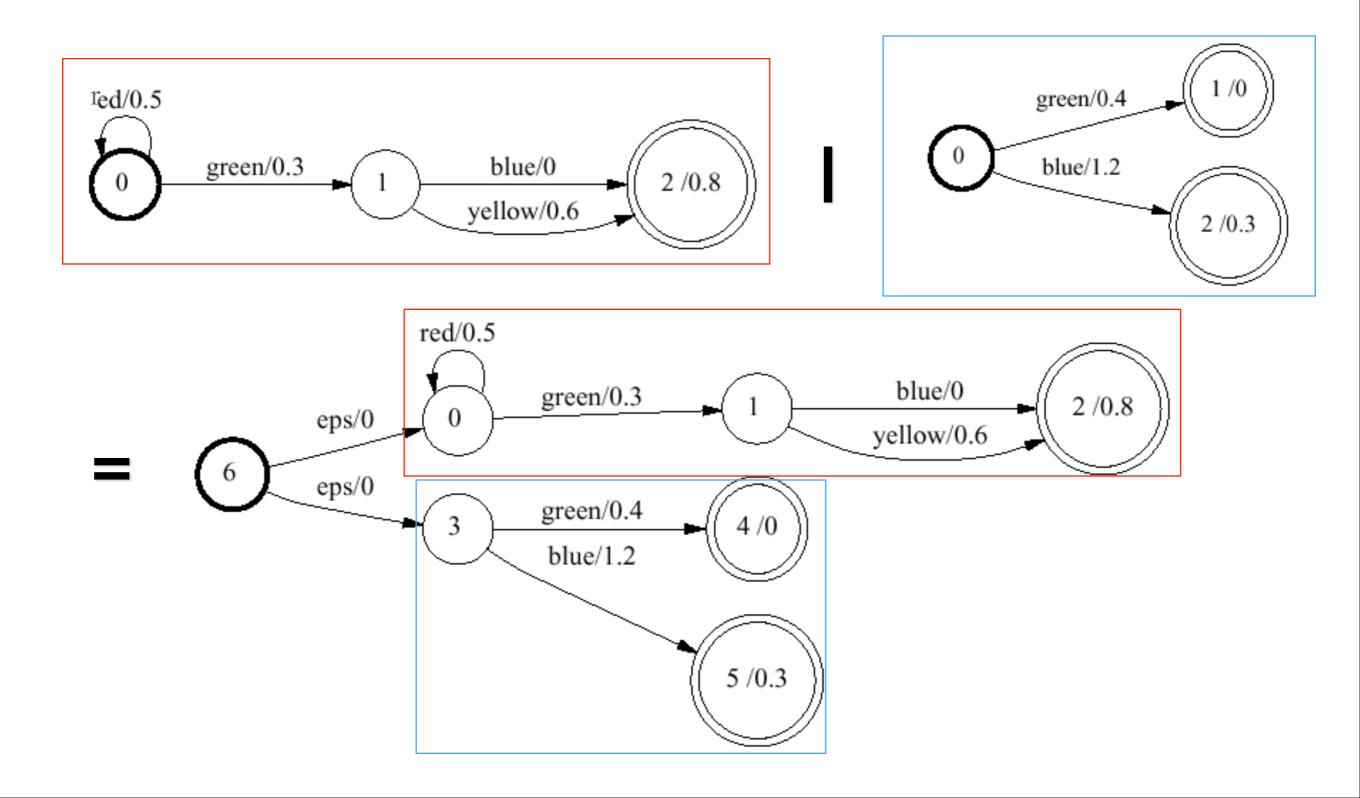
## Concatenation

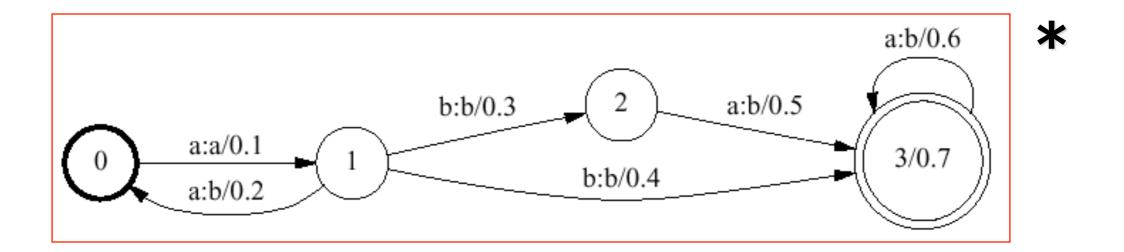


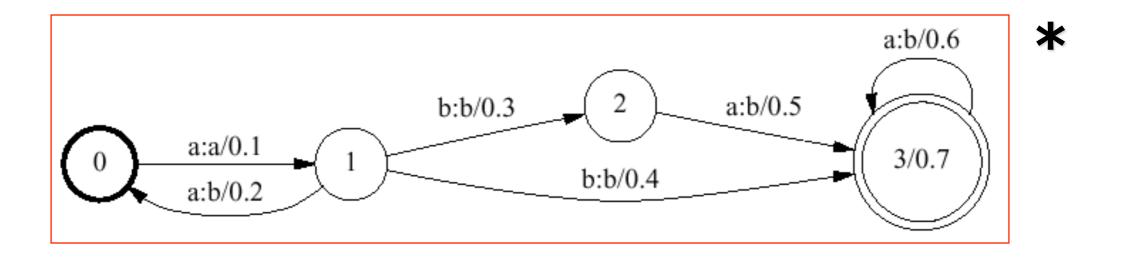
# Union

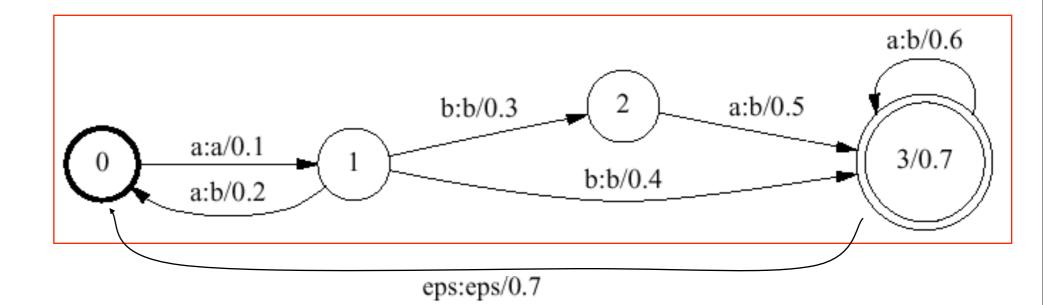


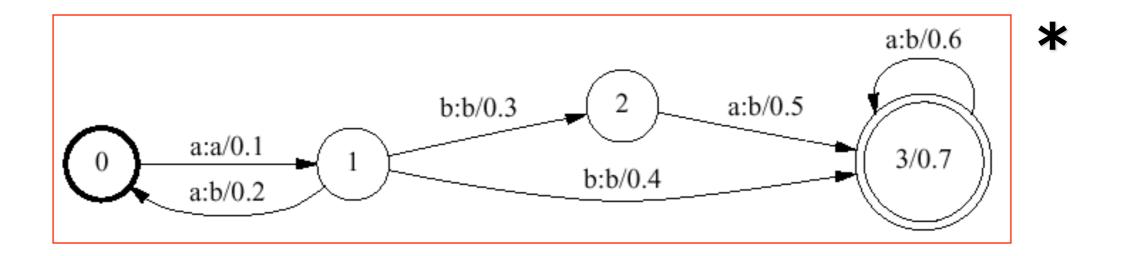
# Union

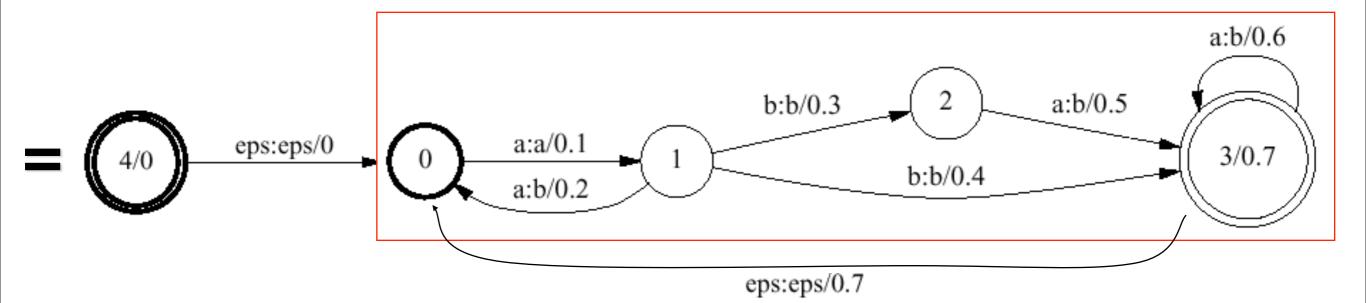


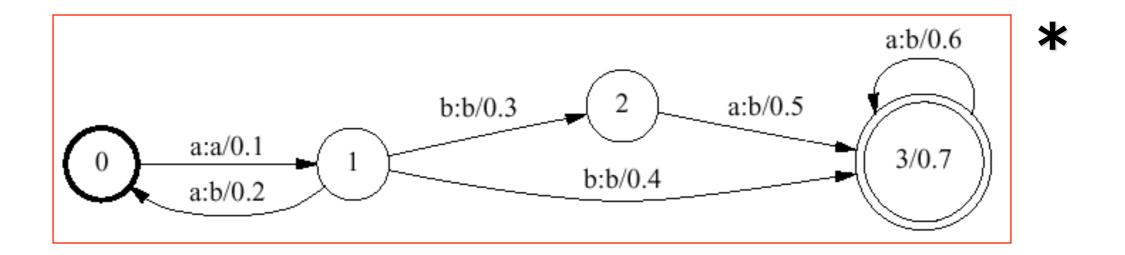


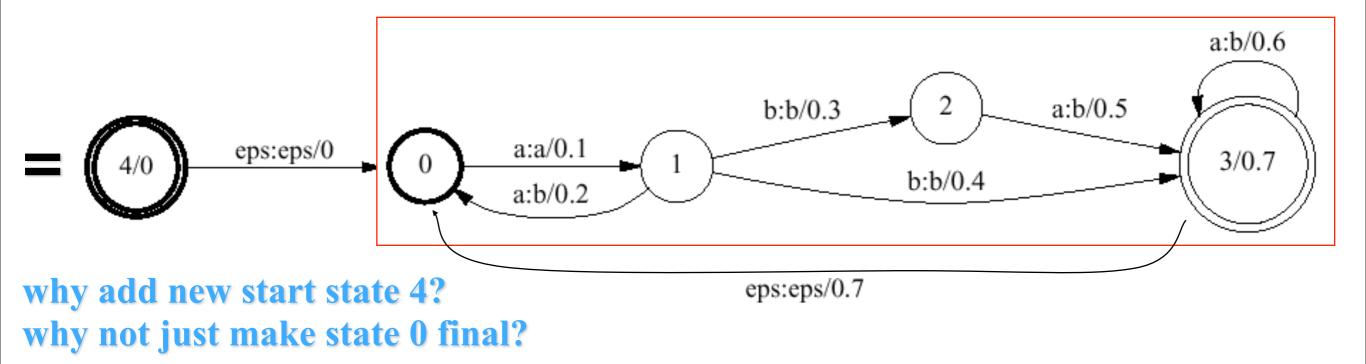




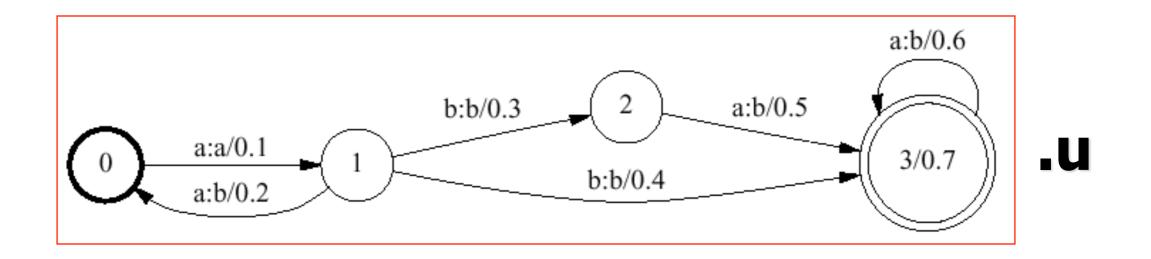


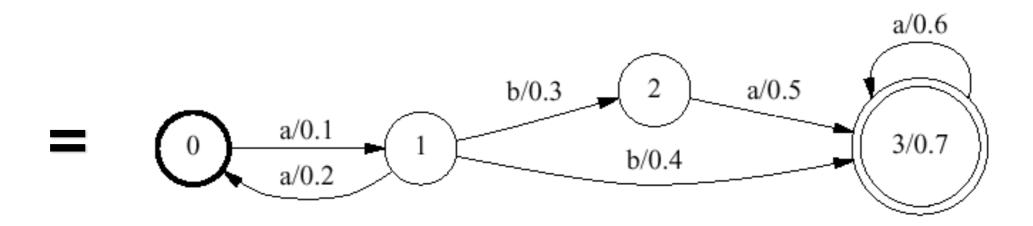




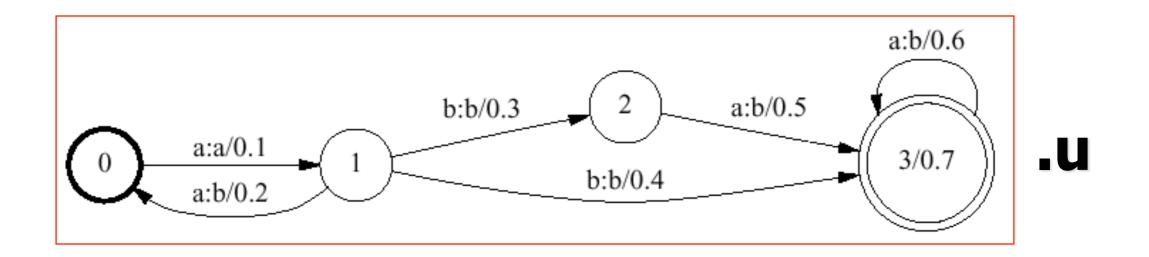


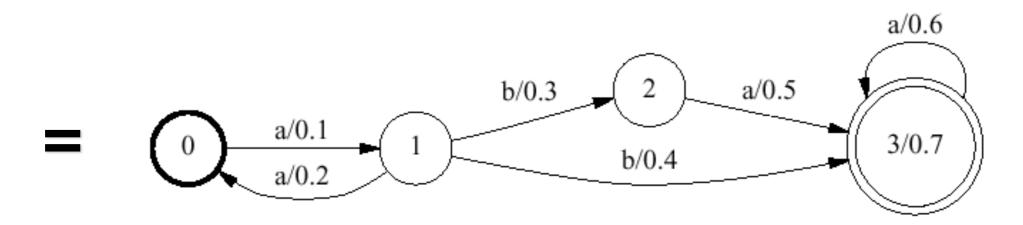
# Upper language (domain)





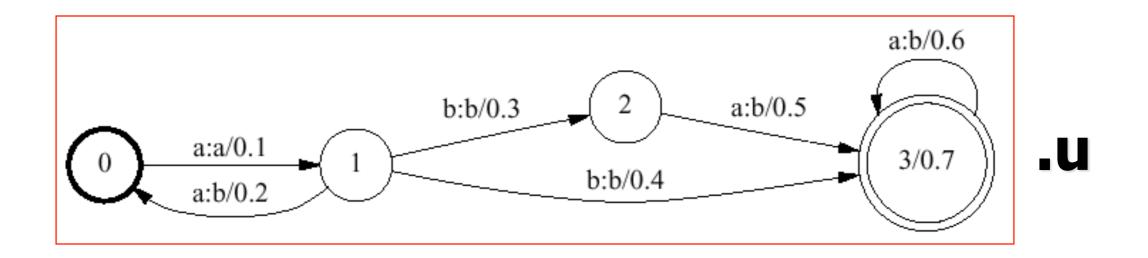
# Upper language (domain)

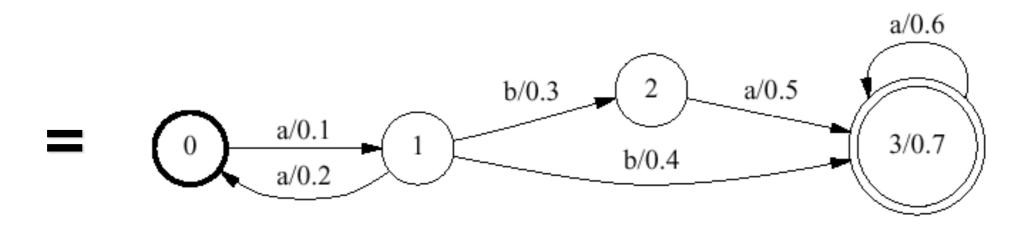




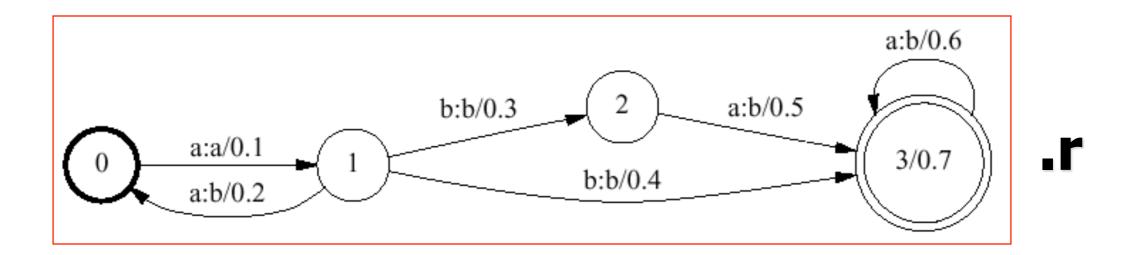
similarly construct lower language .I

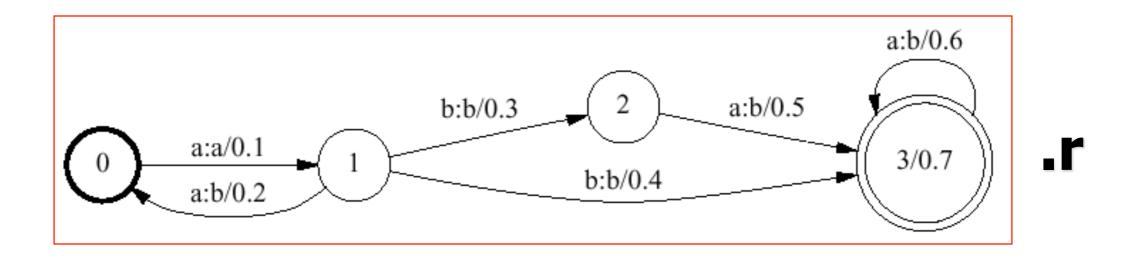
# Upper language (domain)

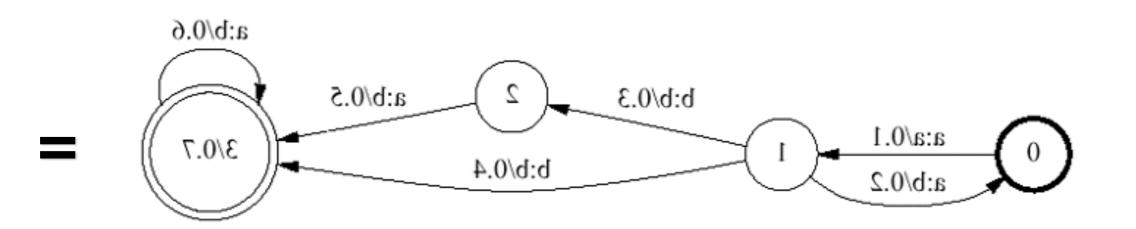


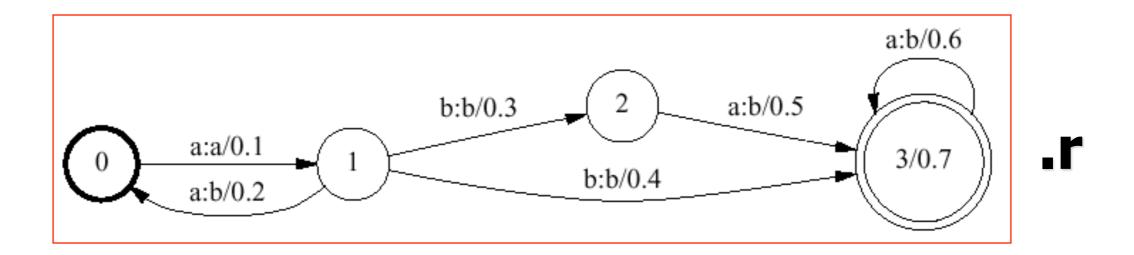


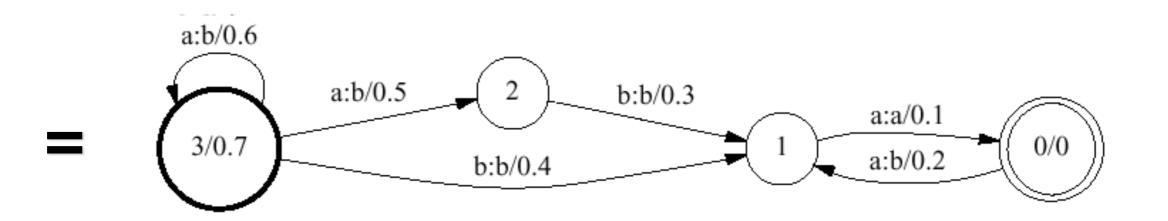
similarly construct lower language .l also called input & output languages



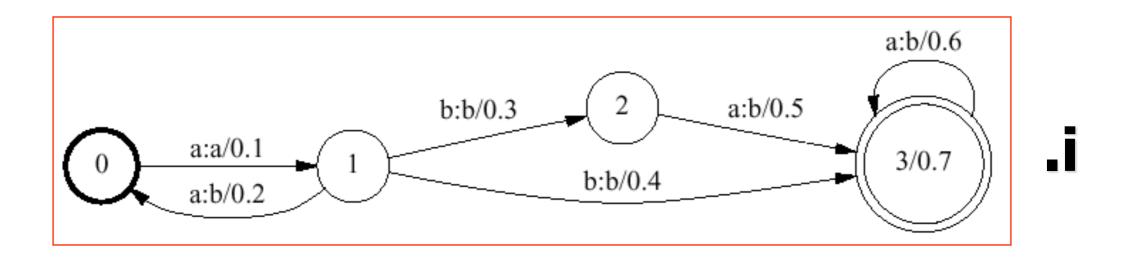




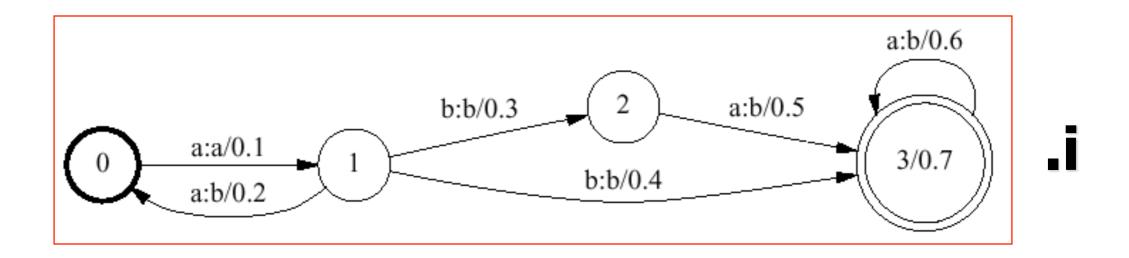


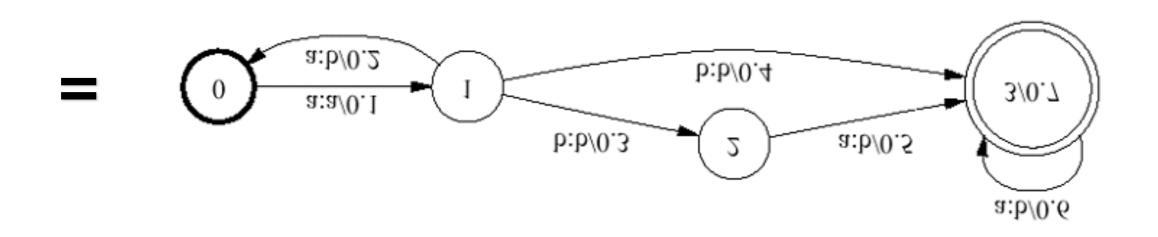


## Inversion

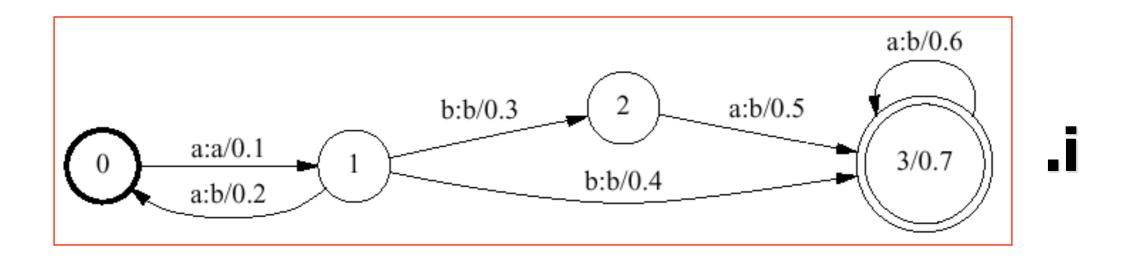


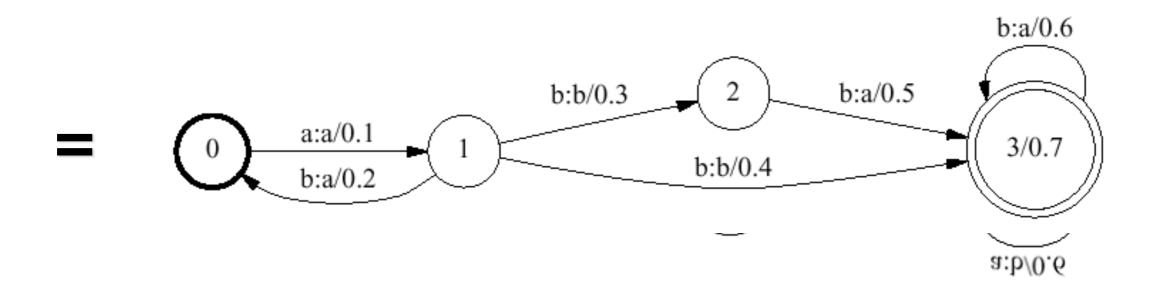
## Inversion





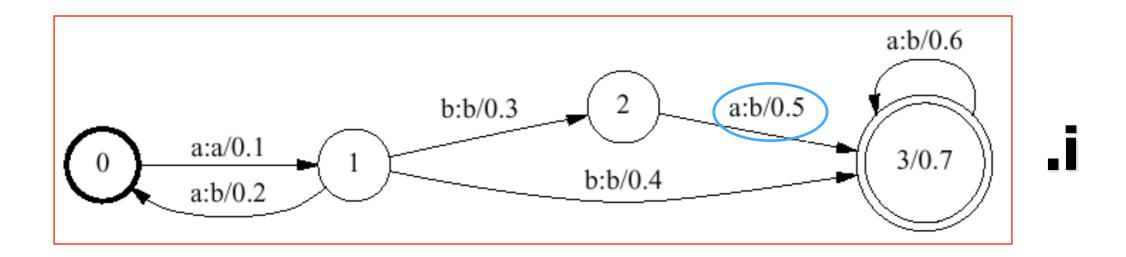
## Inversion

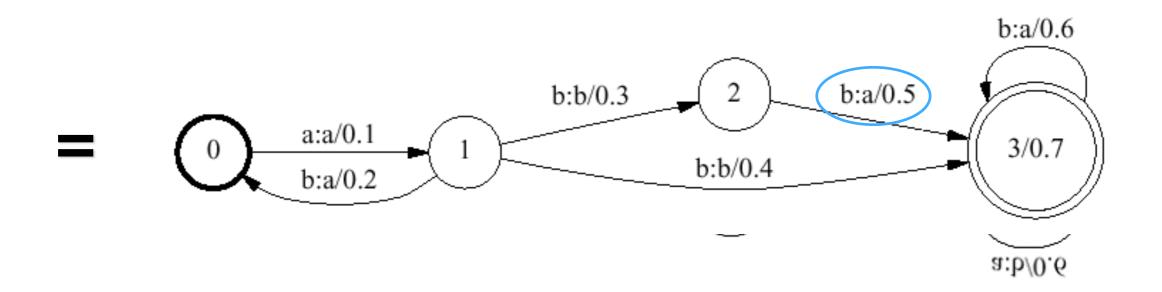




#### example courtesy of M. Mohri

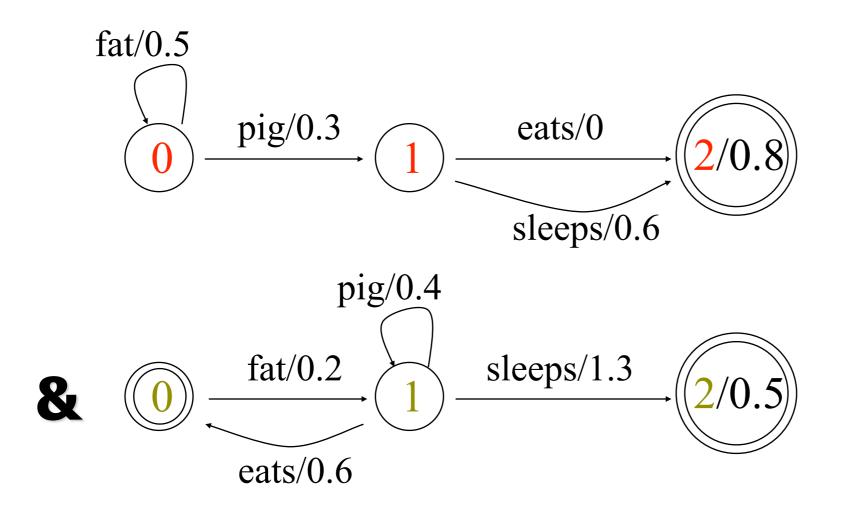
#### Inversion

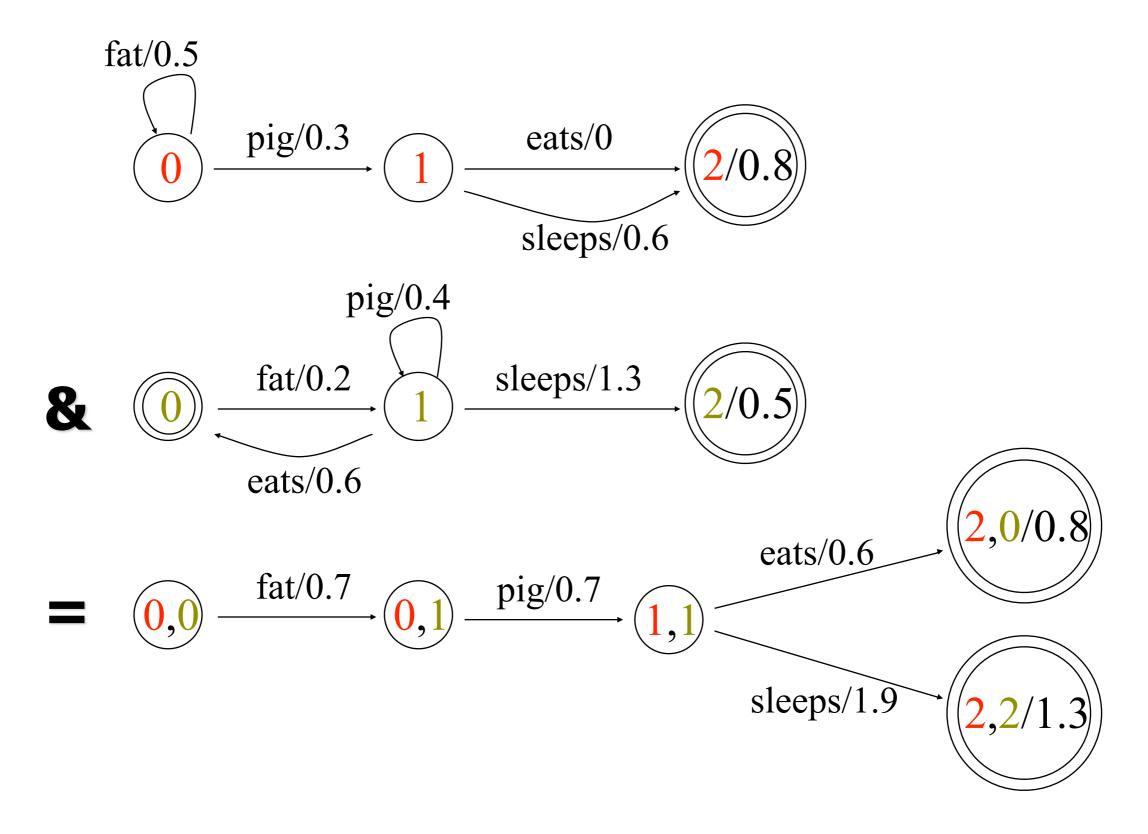


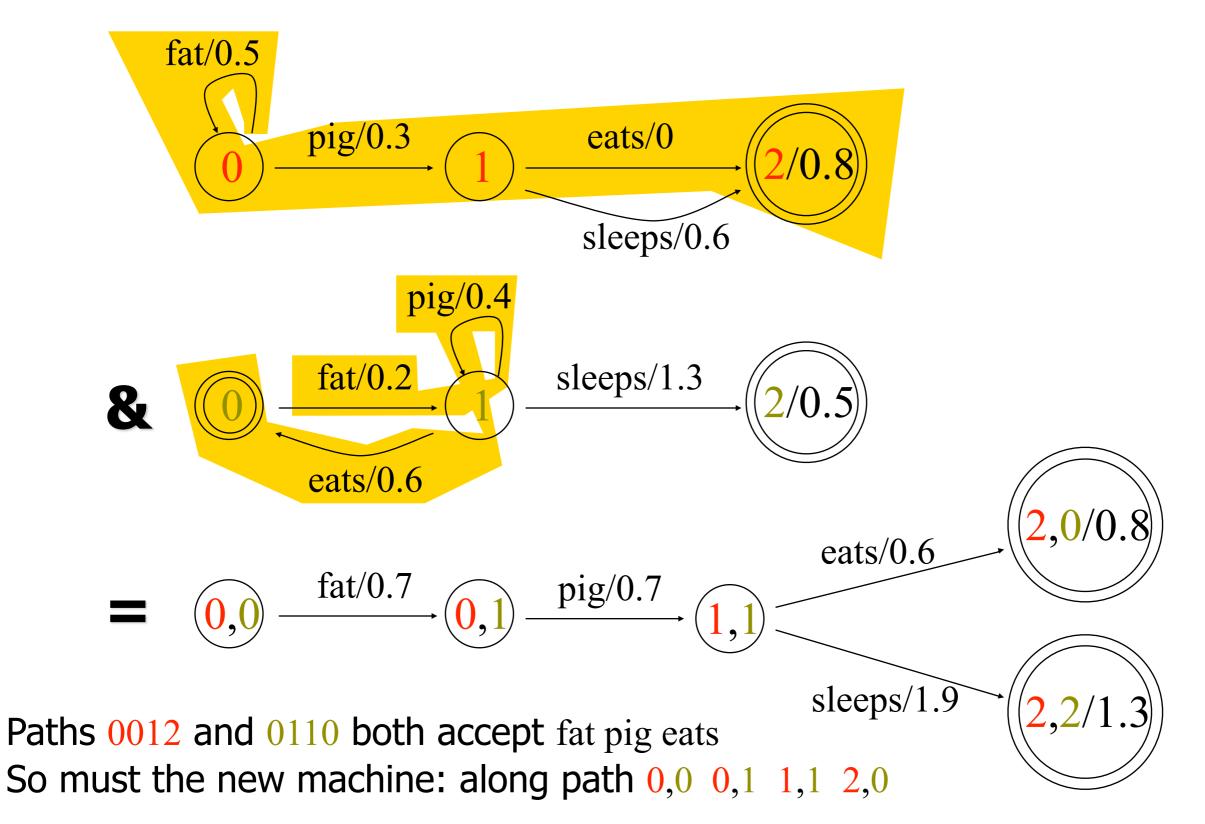


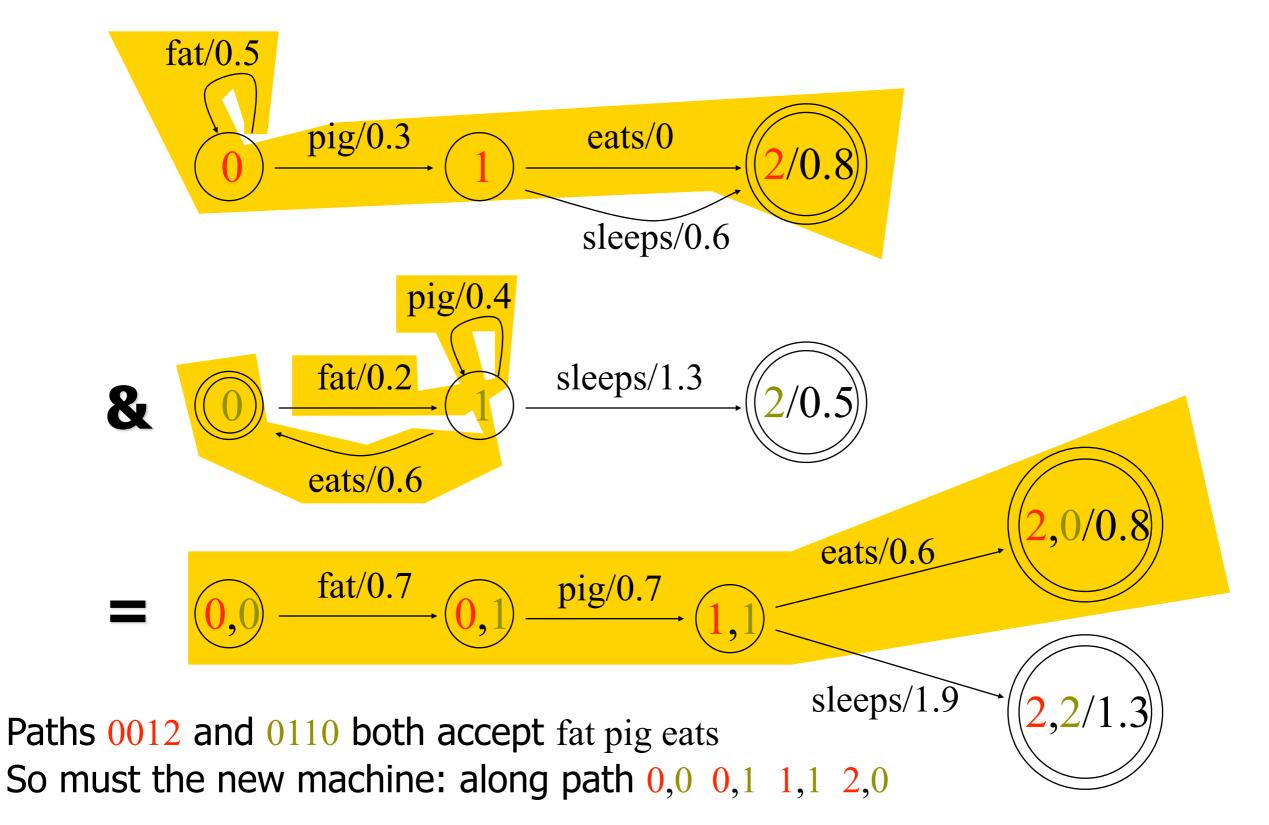
## Complementation

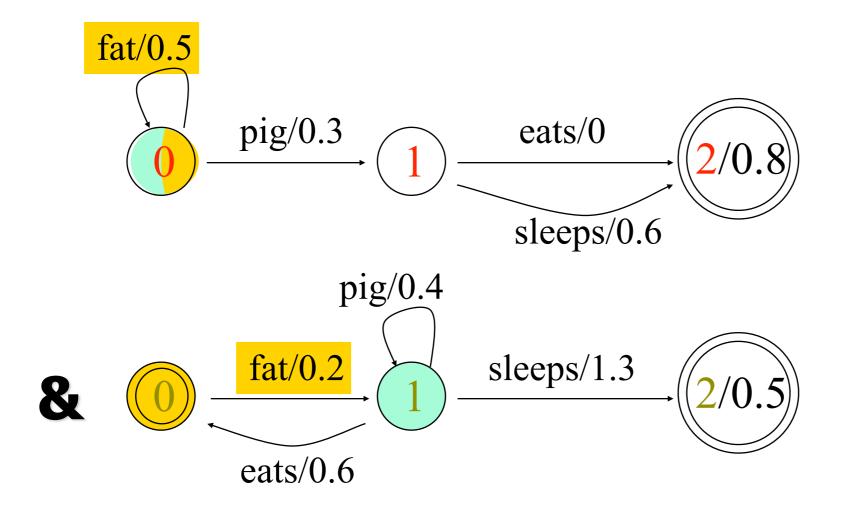
- Given a machine M, represent all strings not accepted by M
- Just change final states to non-final and vice-versa
- Works only if machine has been determinized and completed first





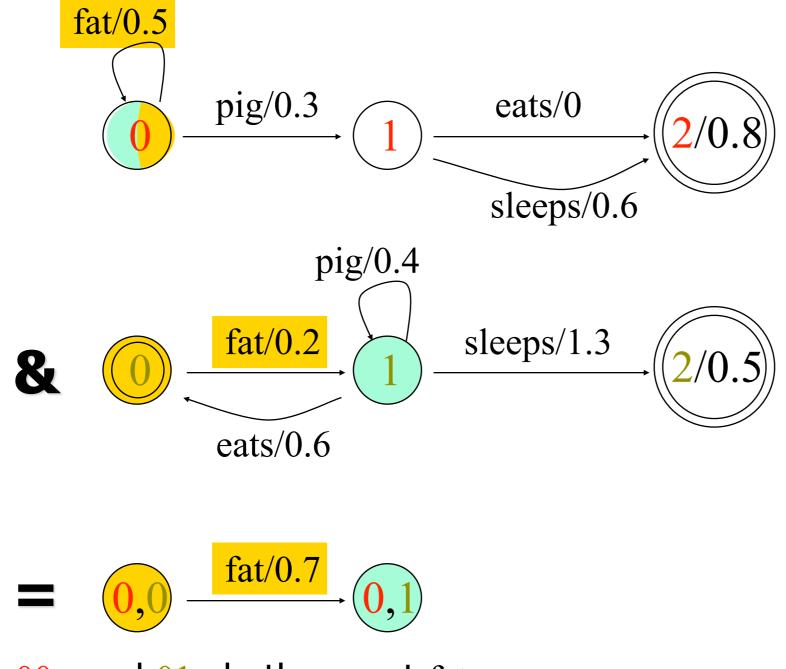




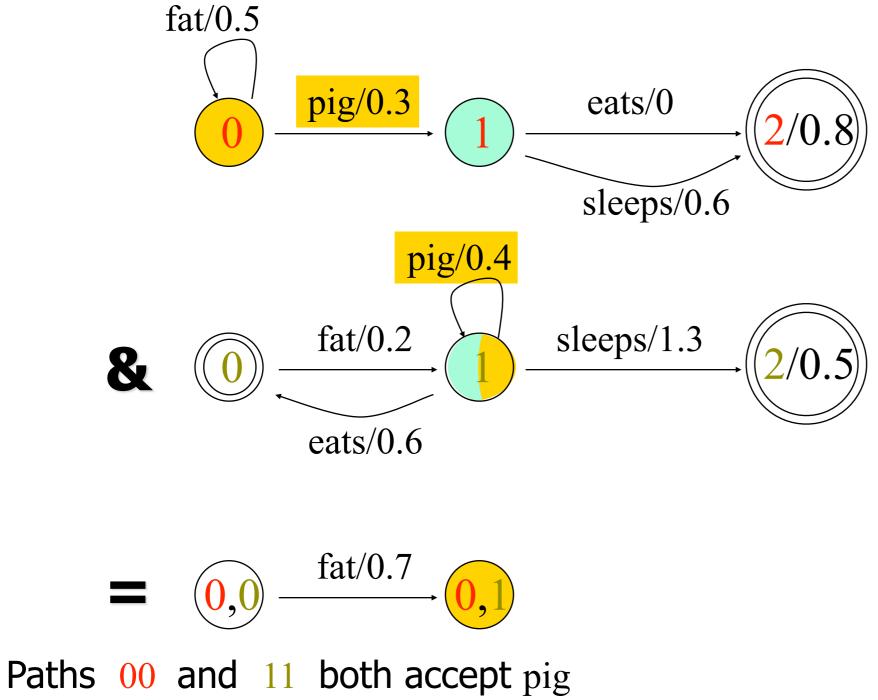




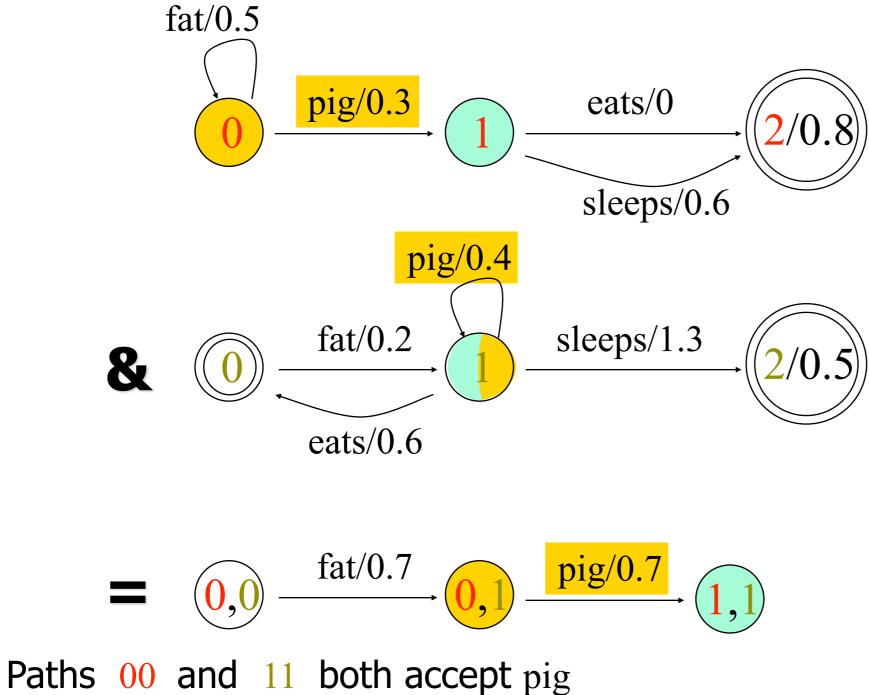
Paths 00 and 01 both accept fat So must the new machine: along path 0,0 0,1



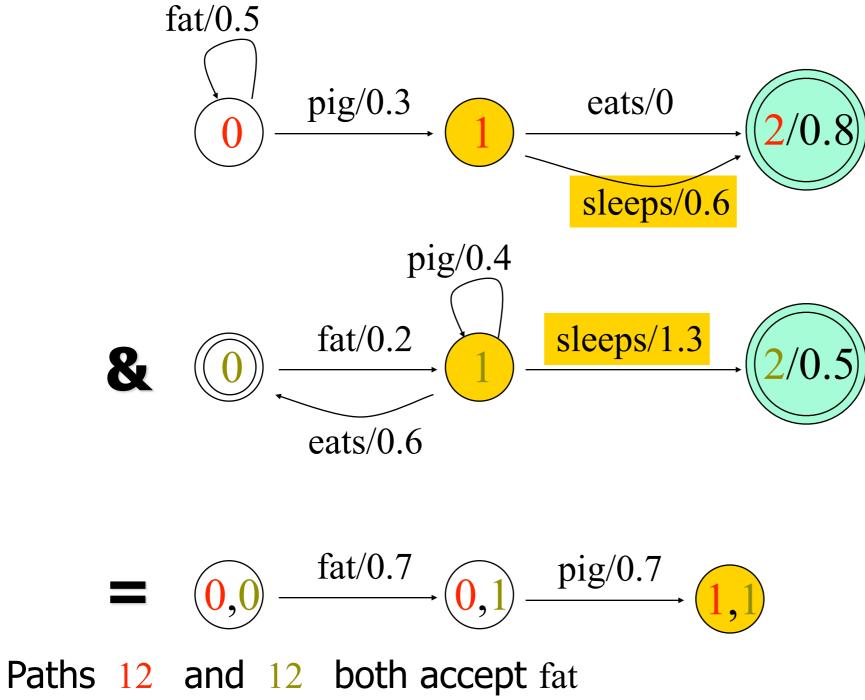
Paths 00 and 01 both accept fat So must the new machine: along path 0,0 0,1



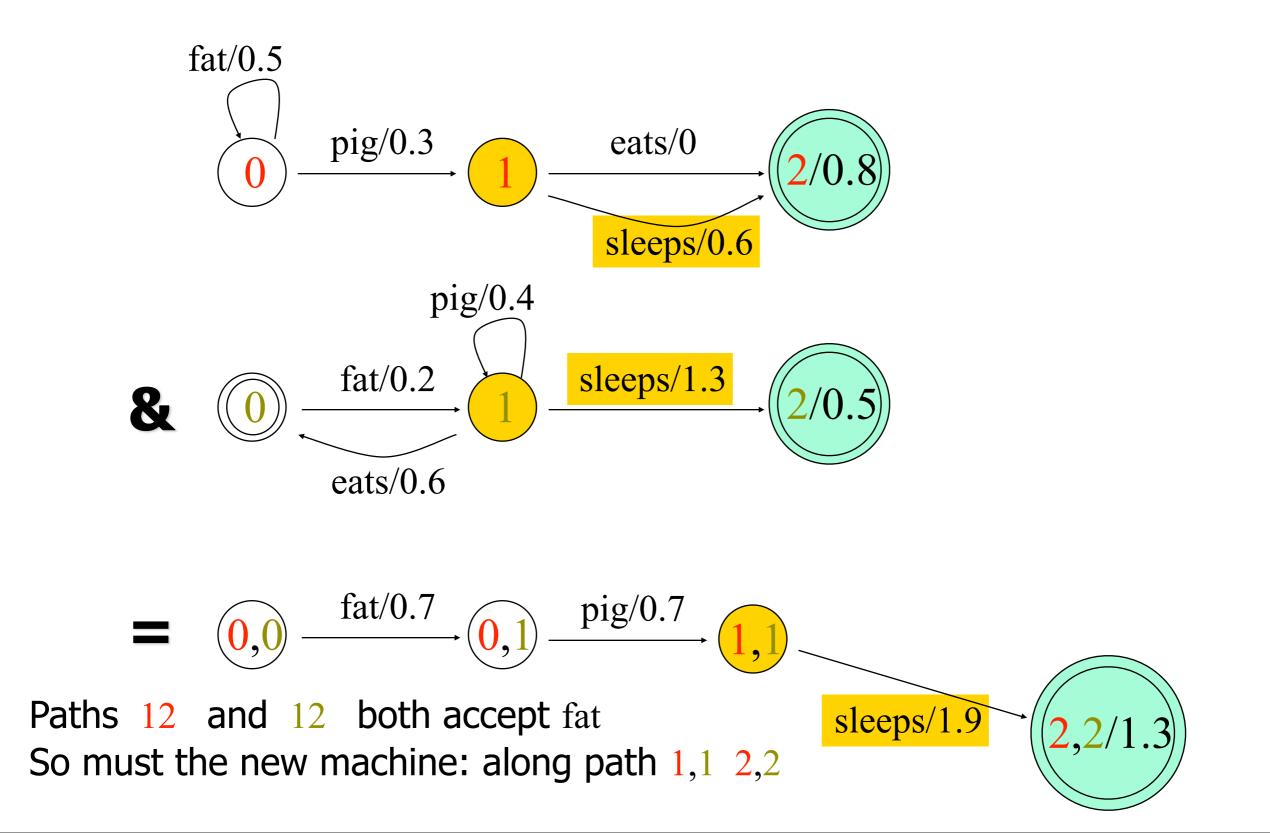
So must the new machine: along path 0,1 1,1

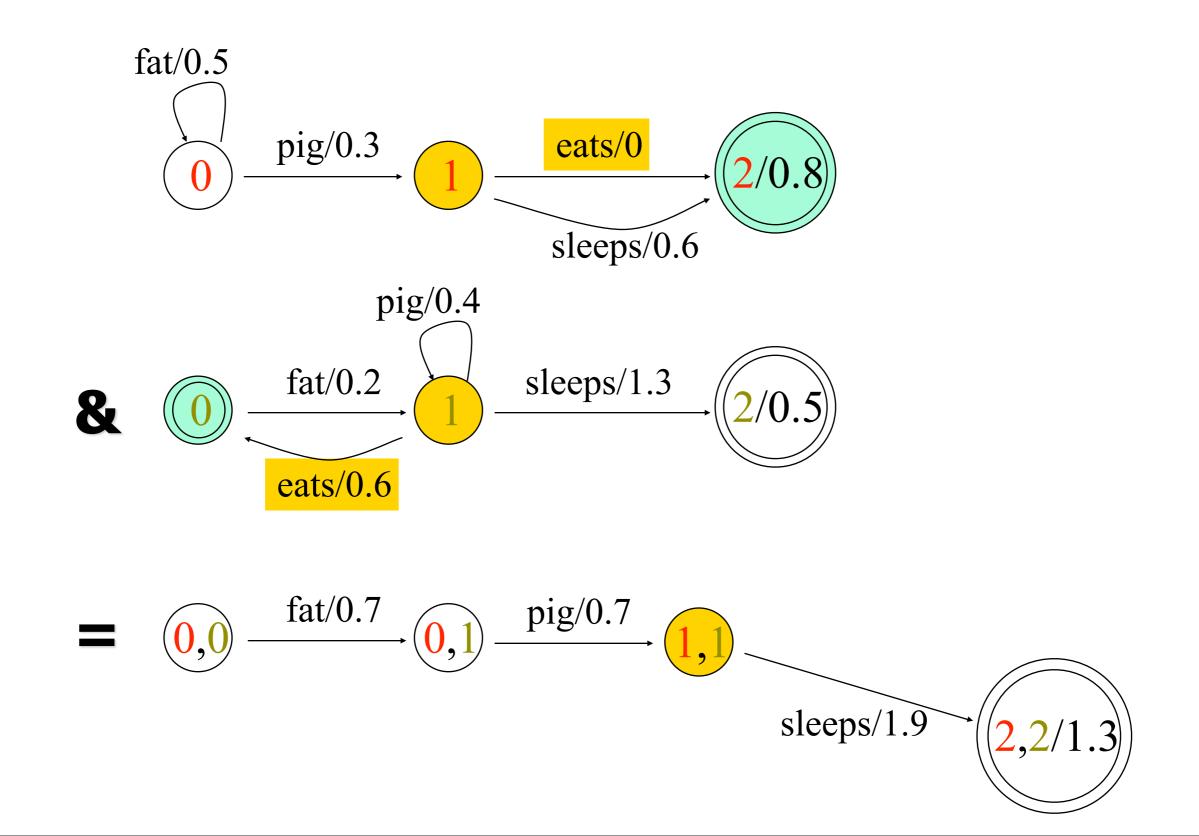


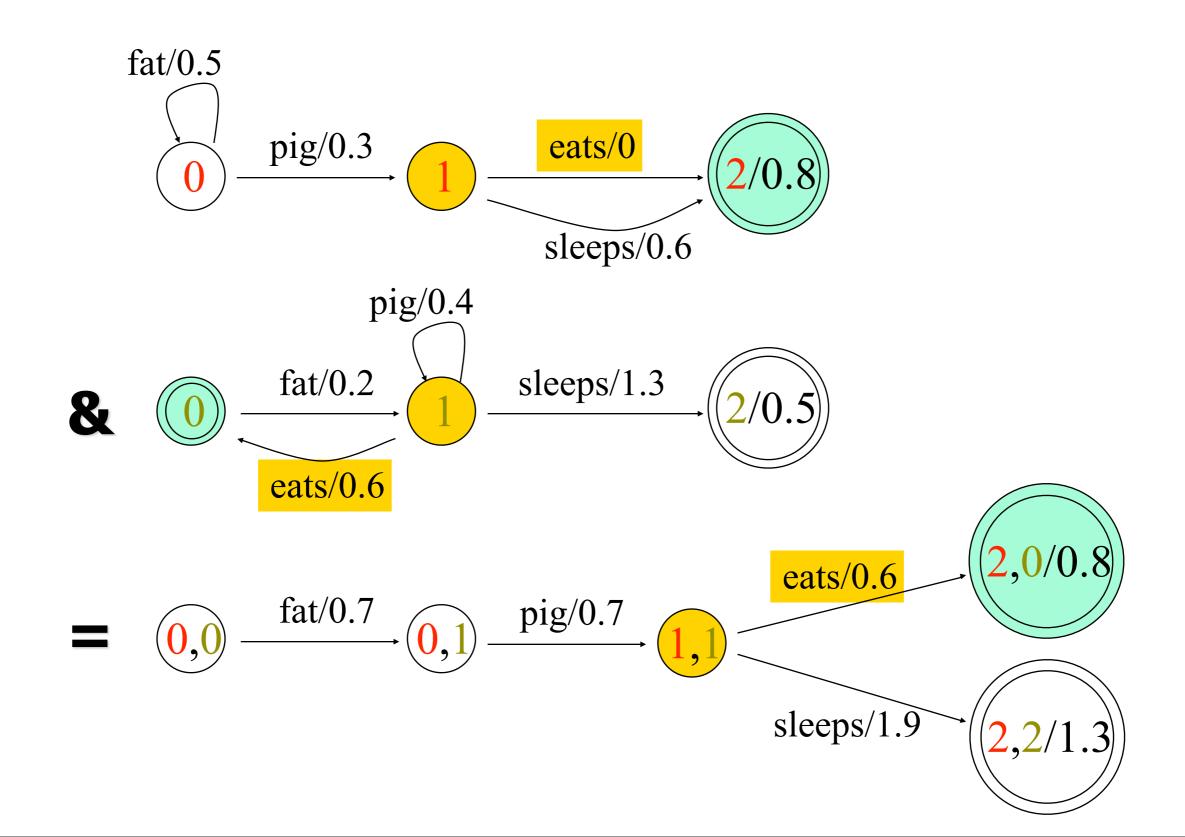
So must the new machine: along path 0,1 1,1

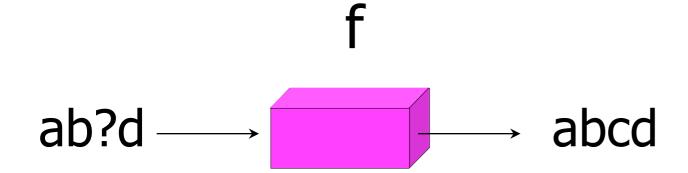


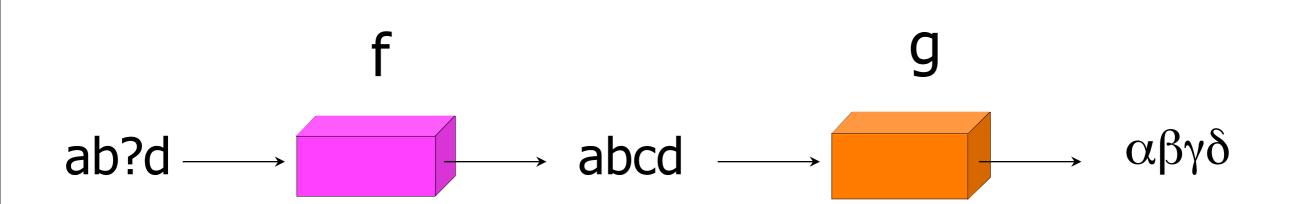
So must the new machine: along path 1,1 2,2

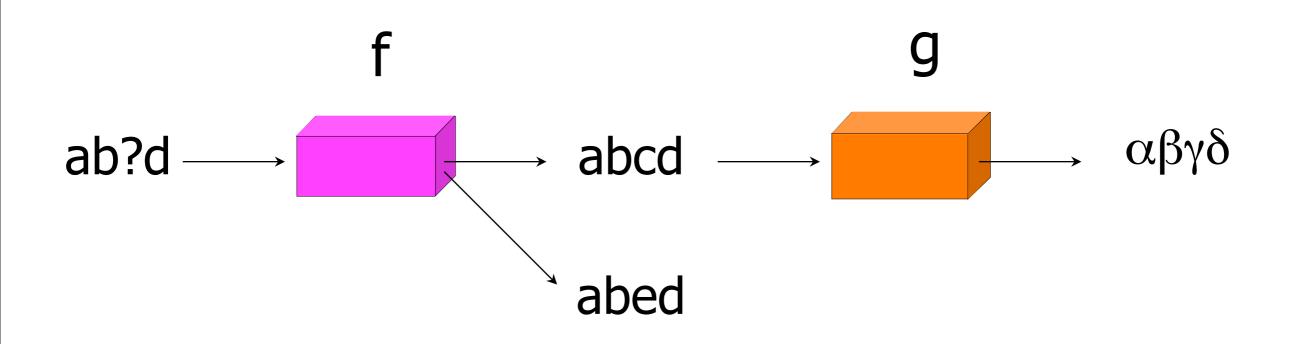


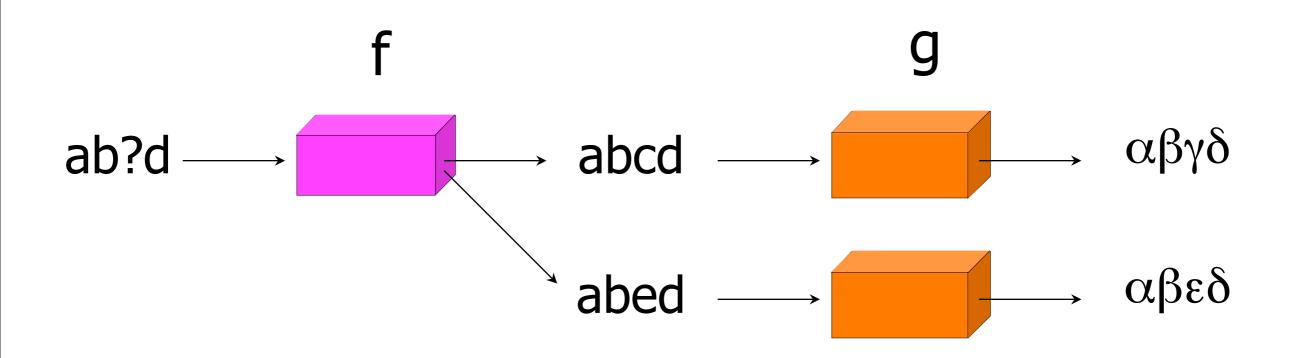


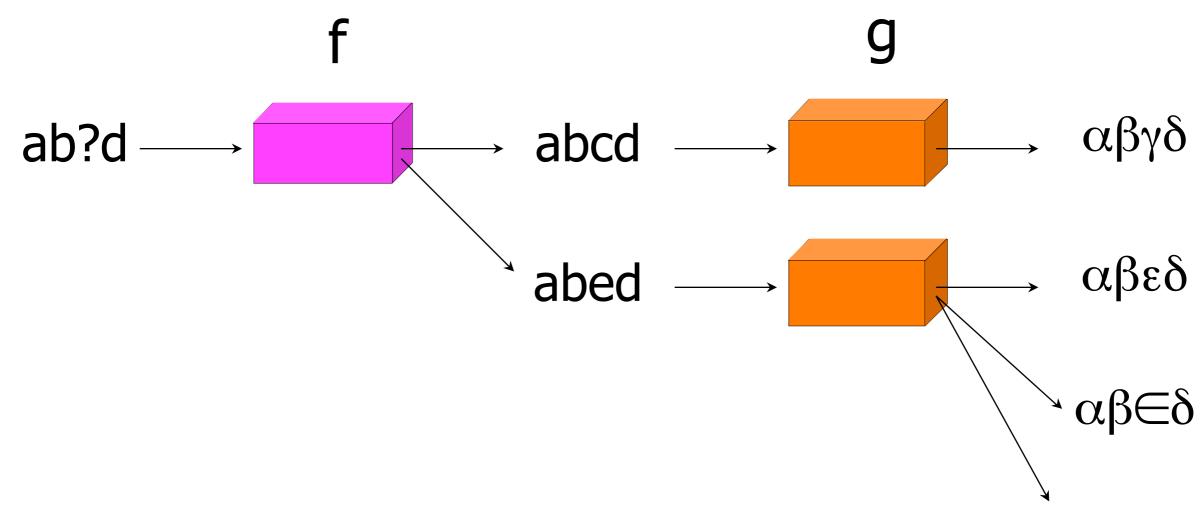


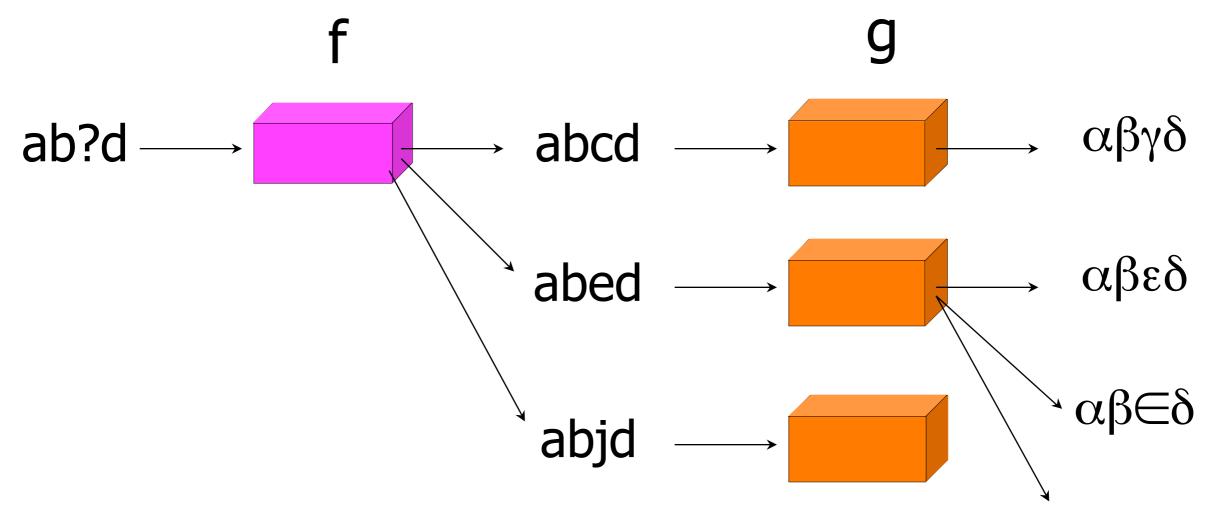


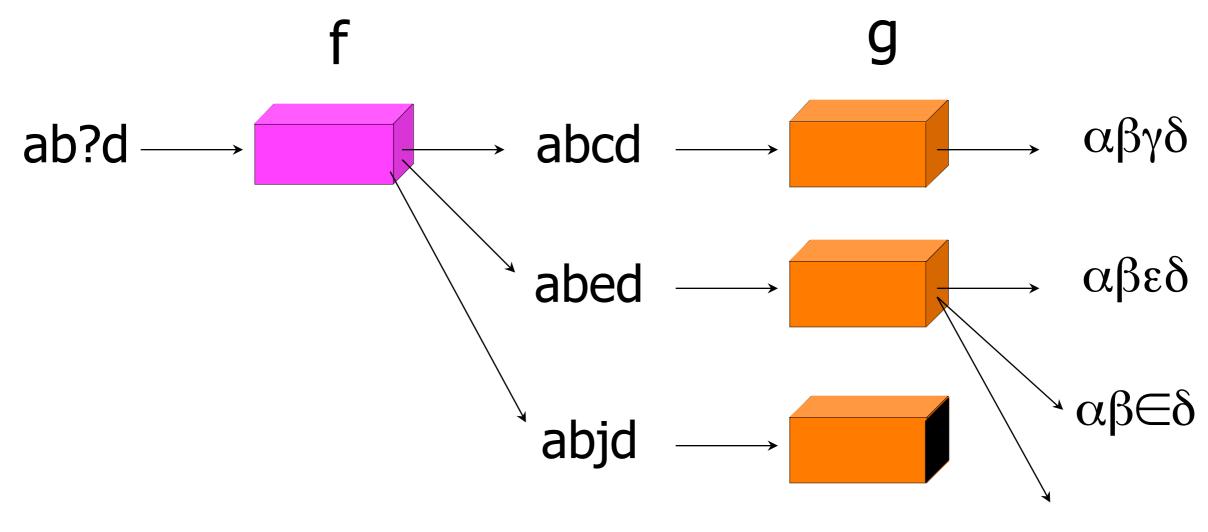


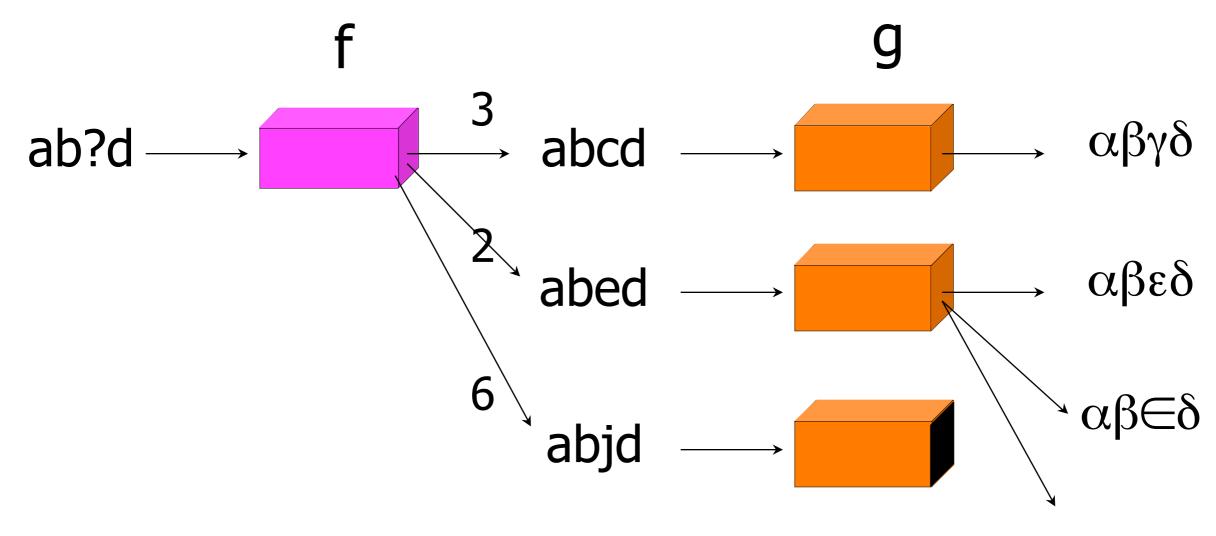


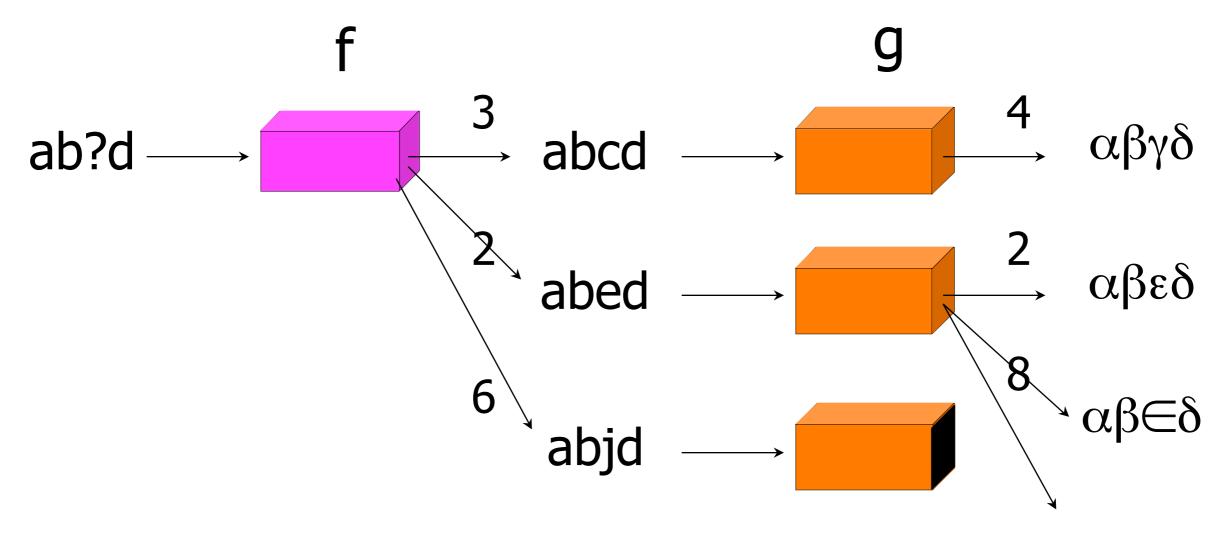


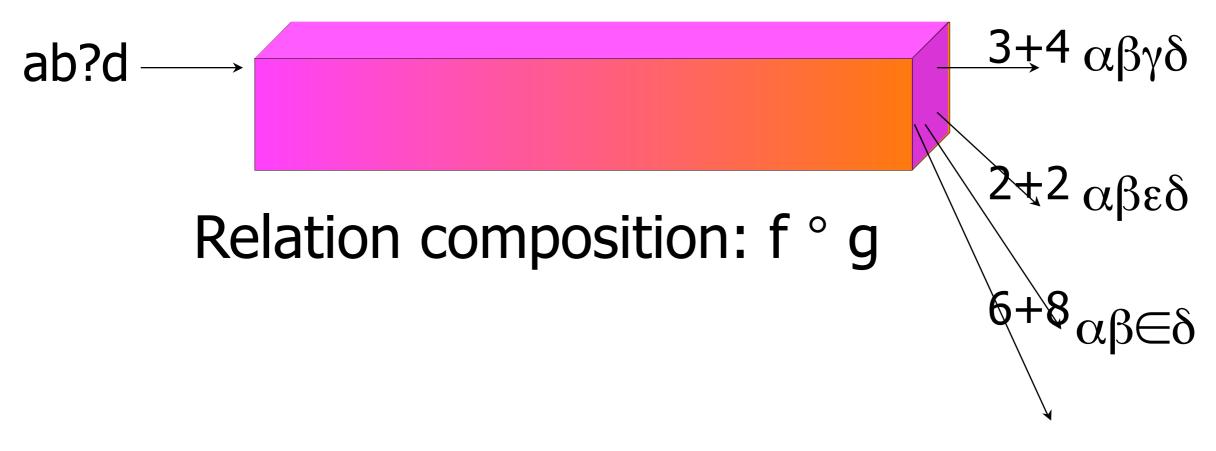




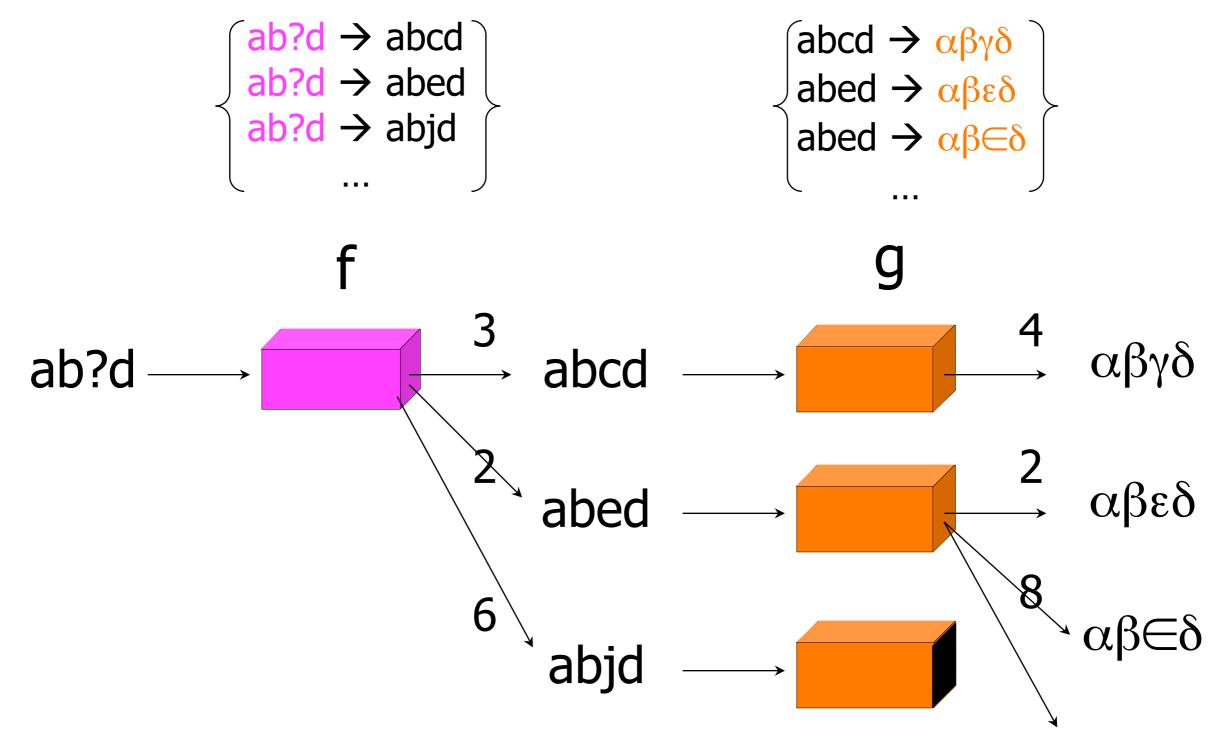


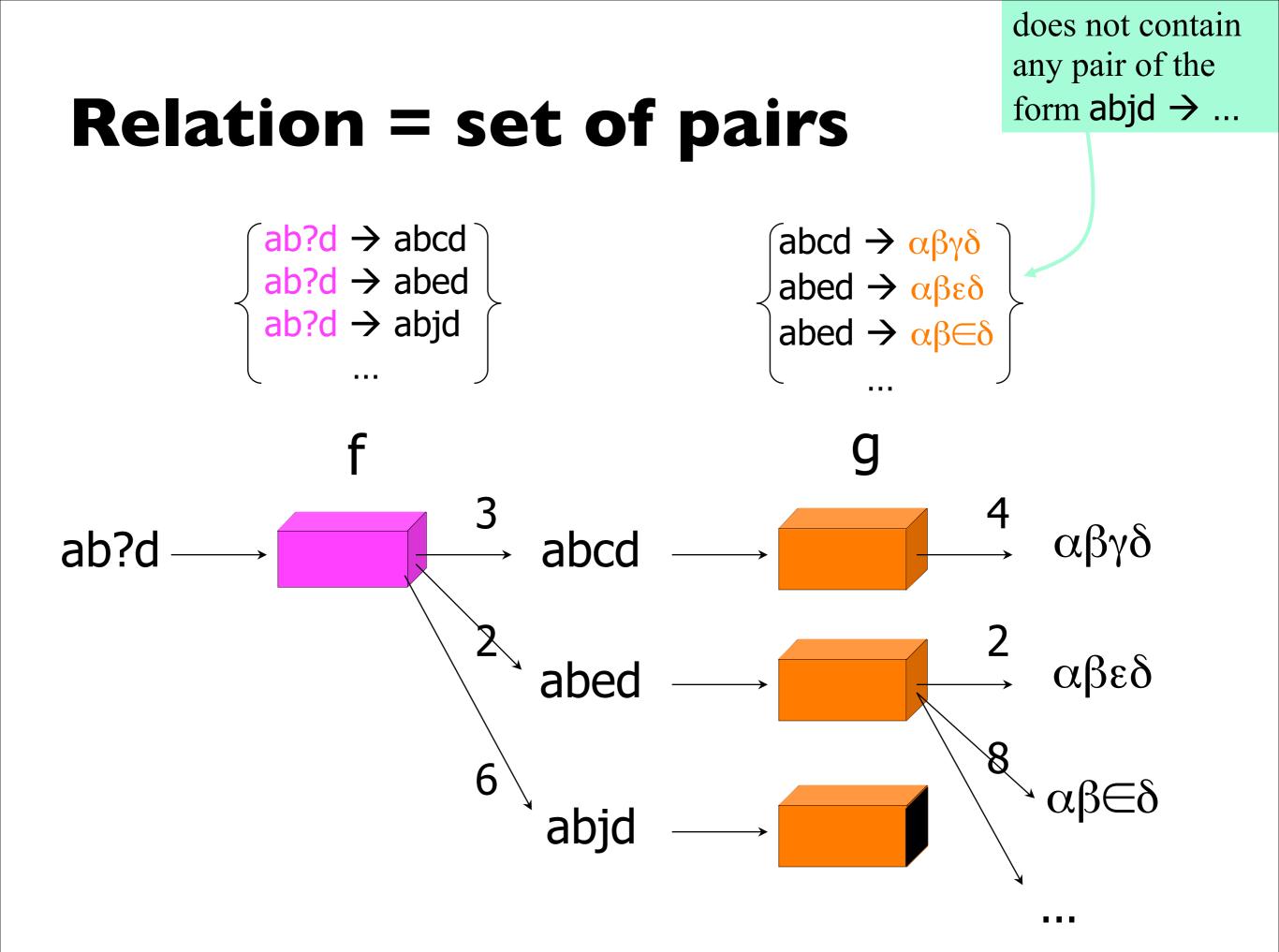




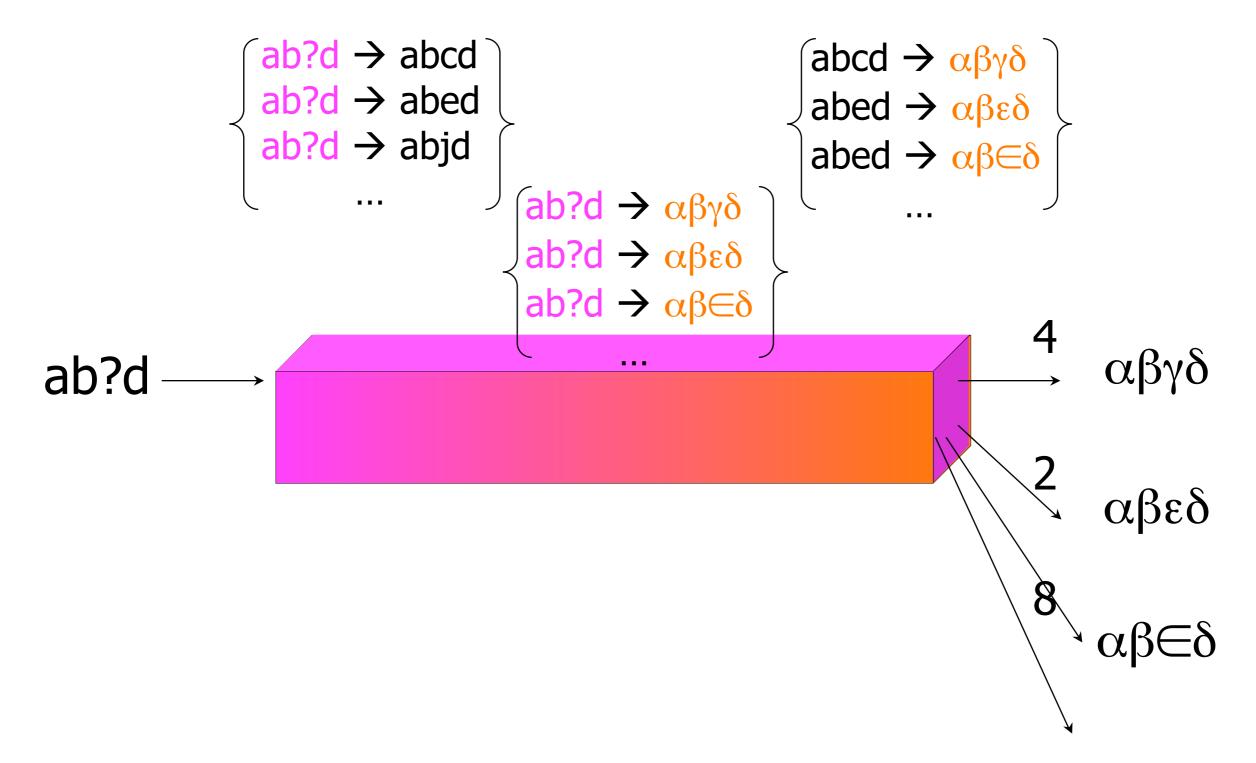


#### **Relation = set of pairs**

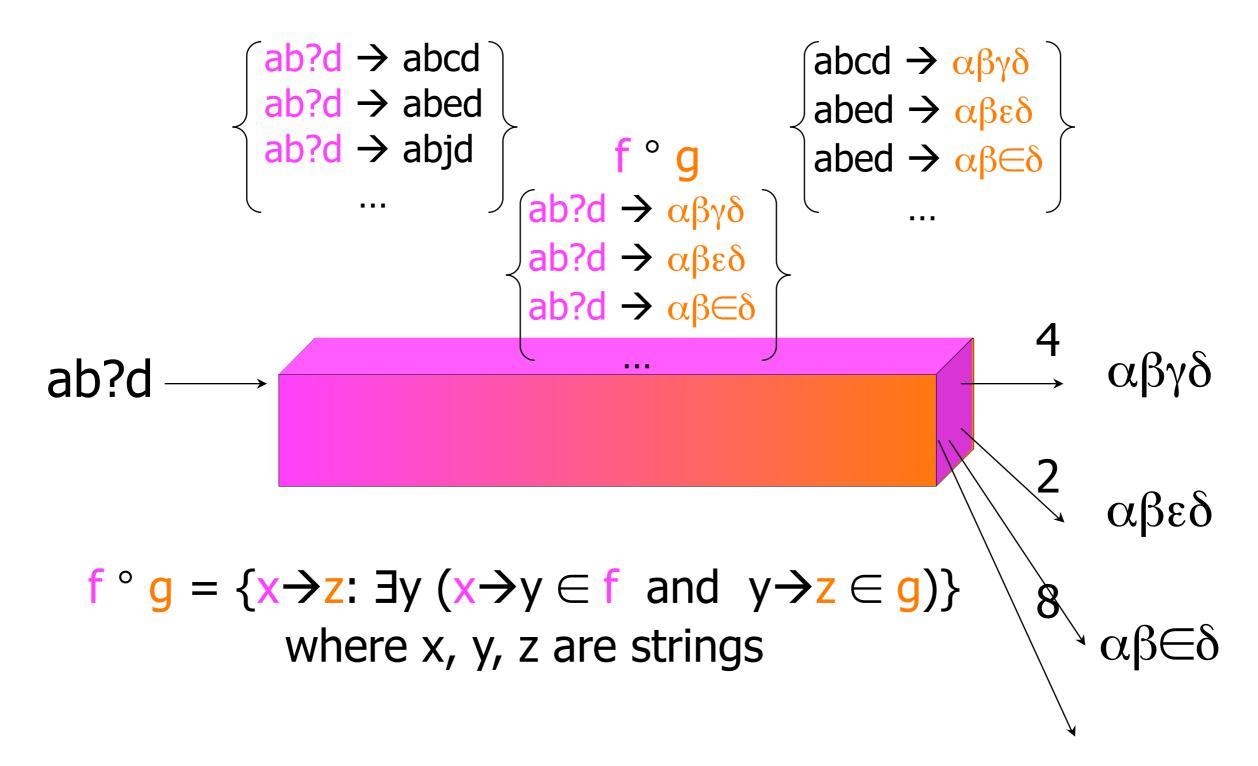




#### **Relation = set of pairs**

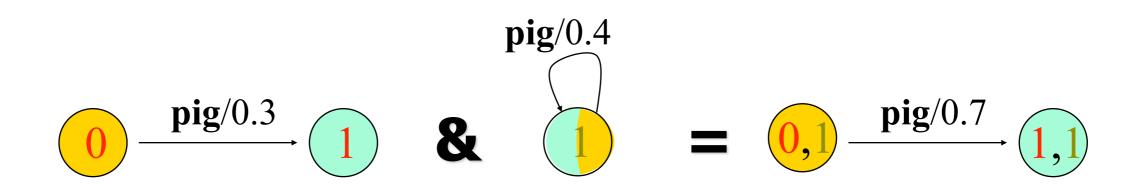


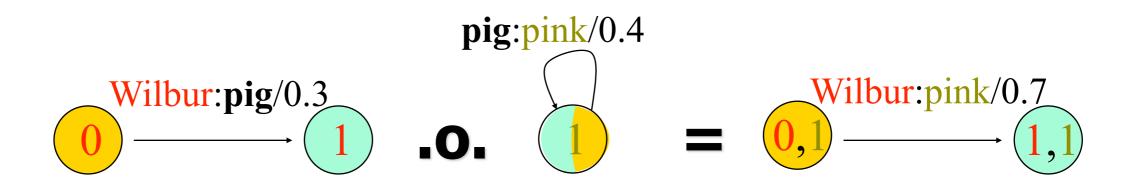
#### **Relation = set of pairs**



## Intersection vs. Composition

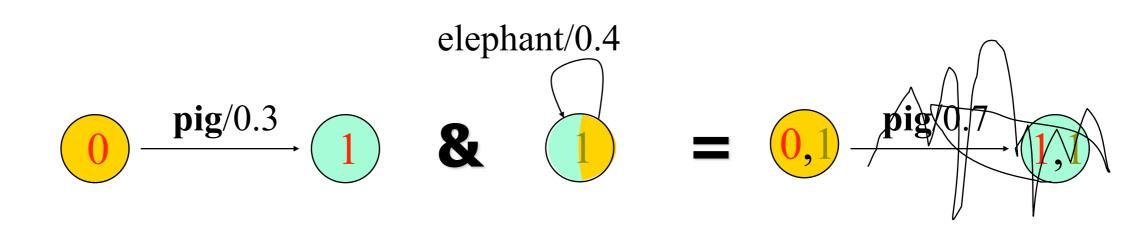
#### Intersection



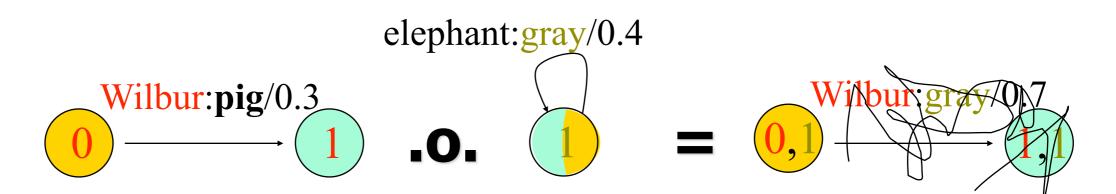


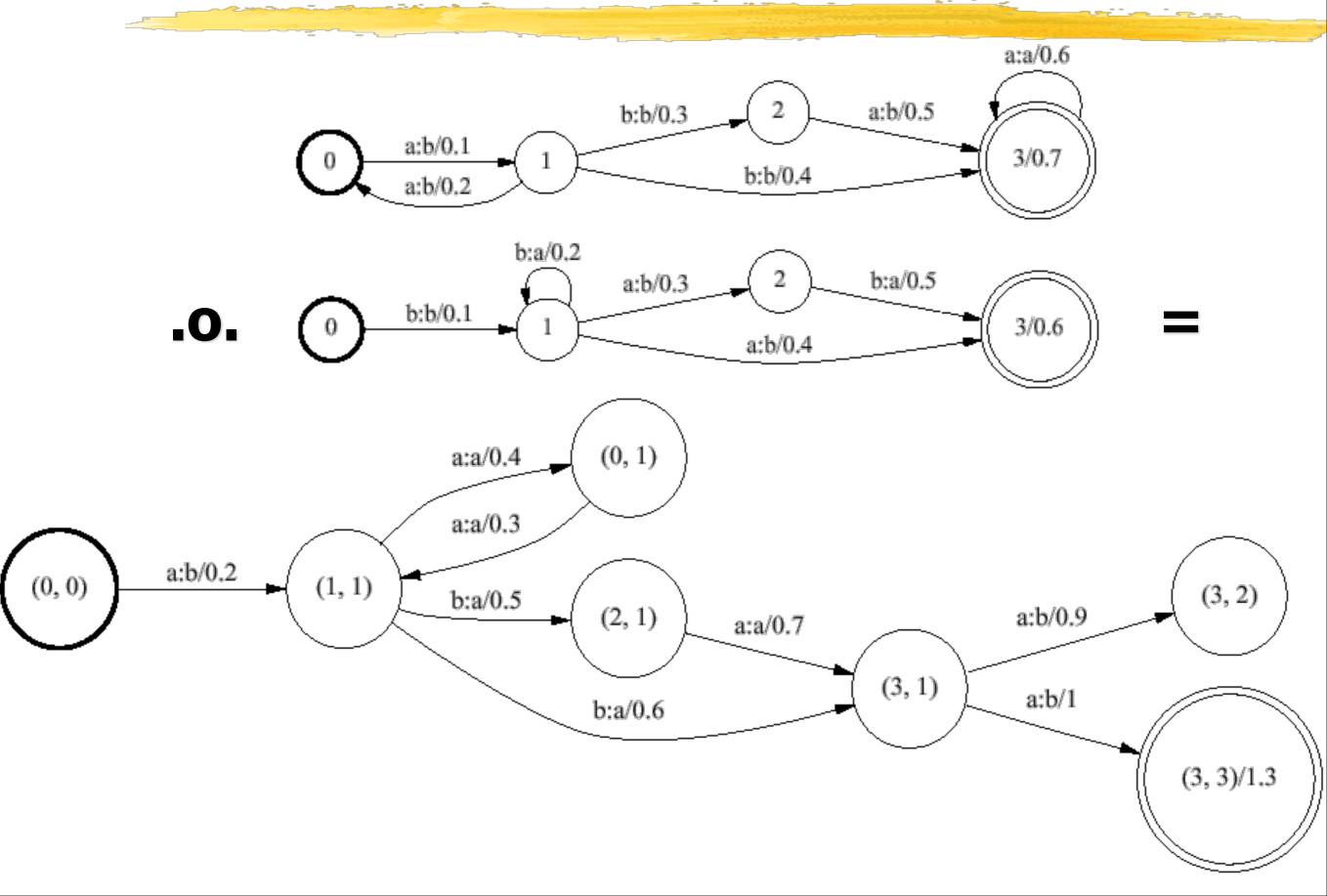
# Intersection vs. Composition

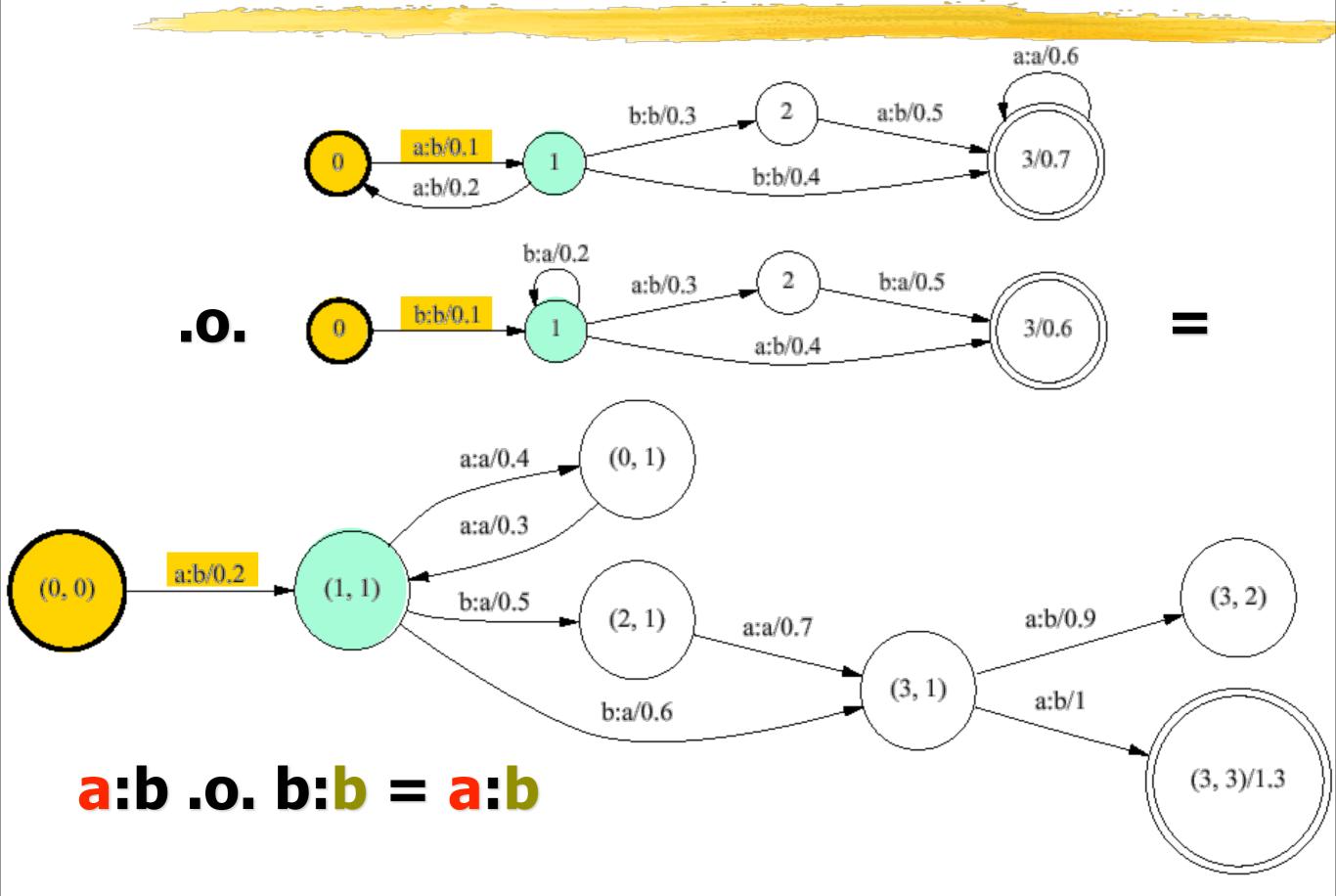
#### Intersection mismatch

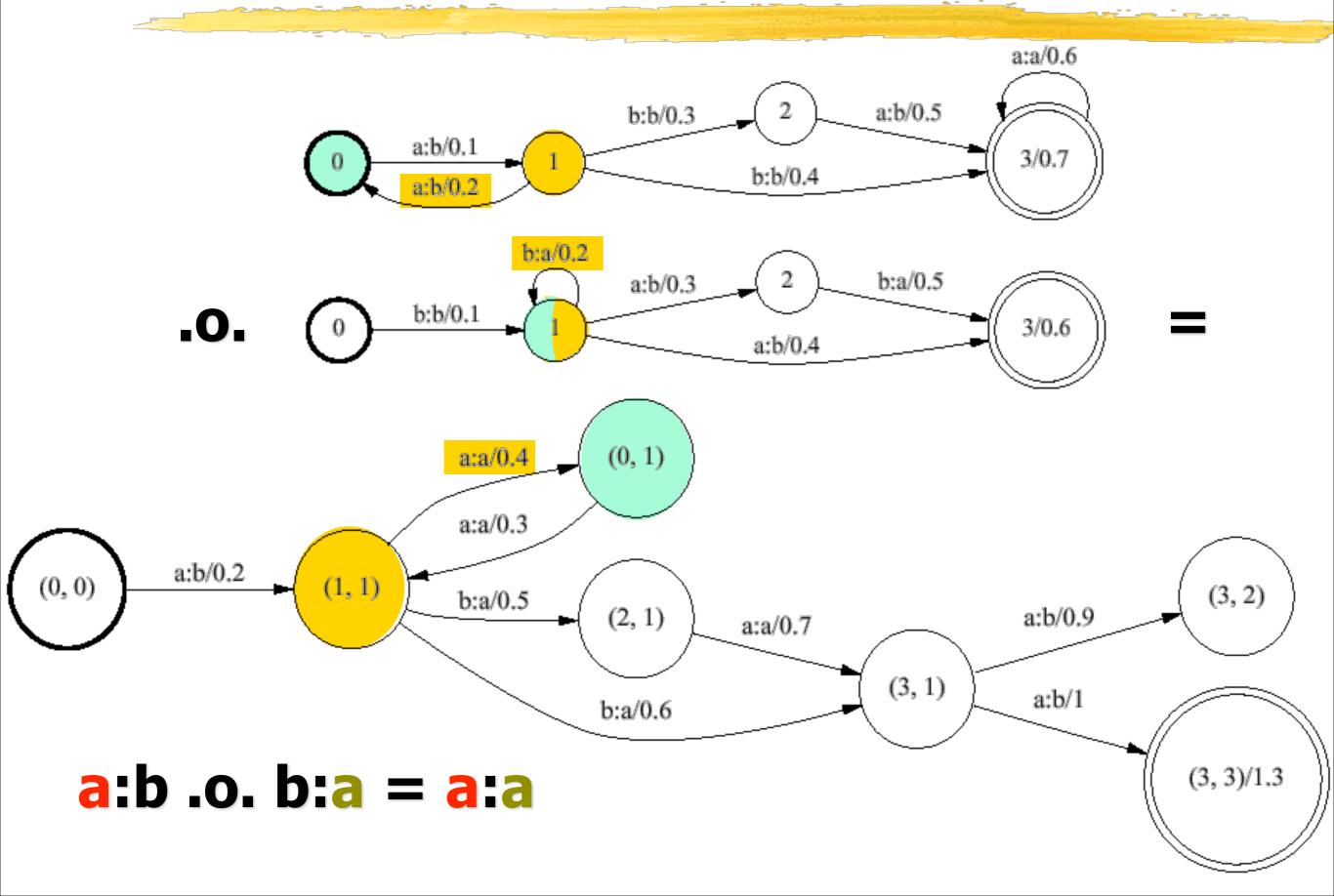


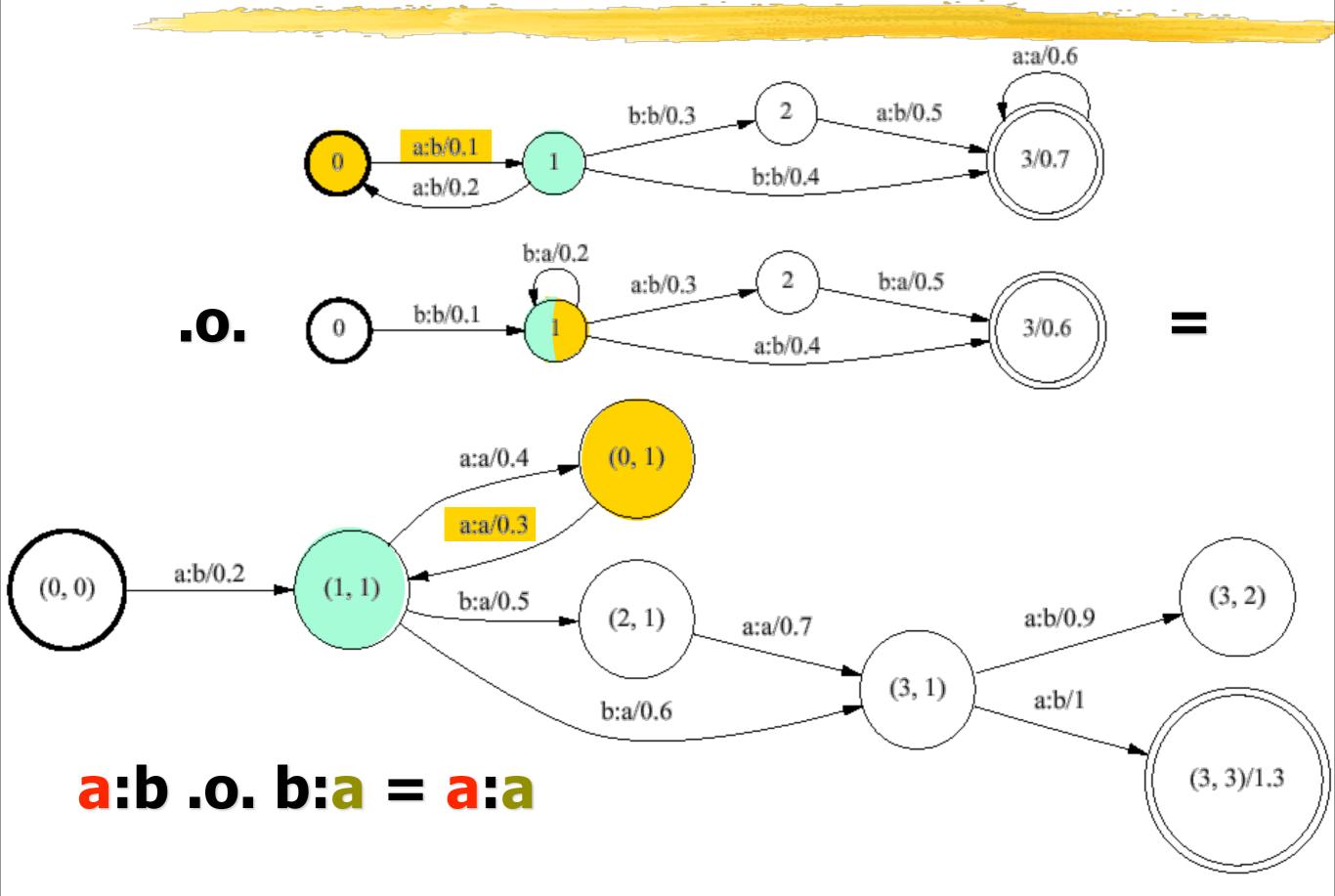
#### **Composition mismatch**

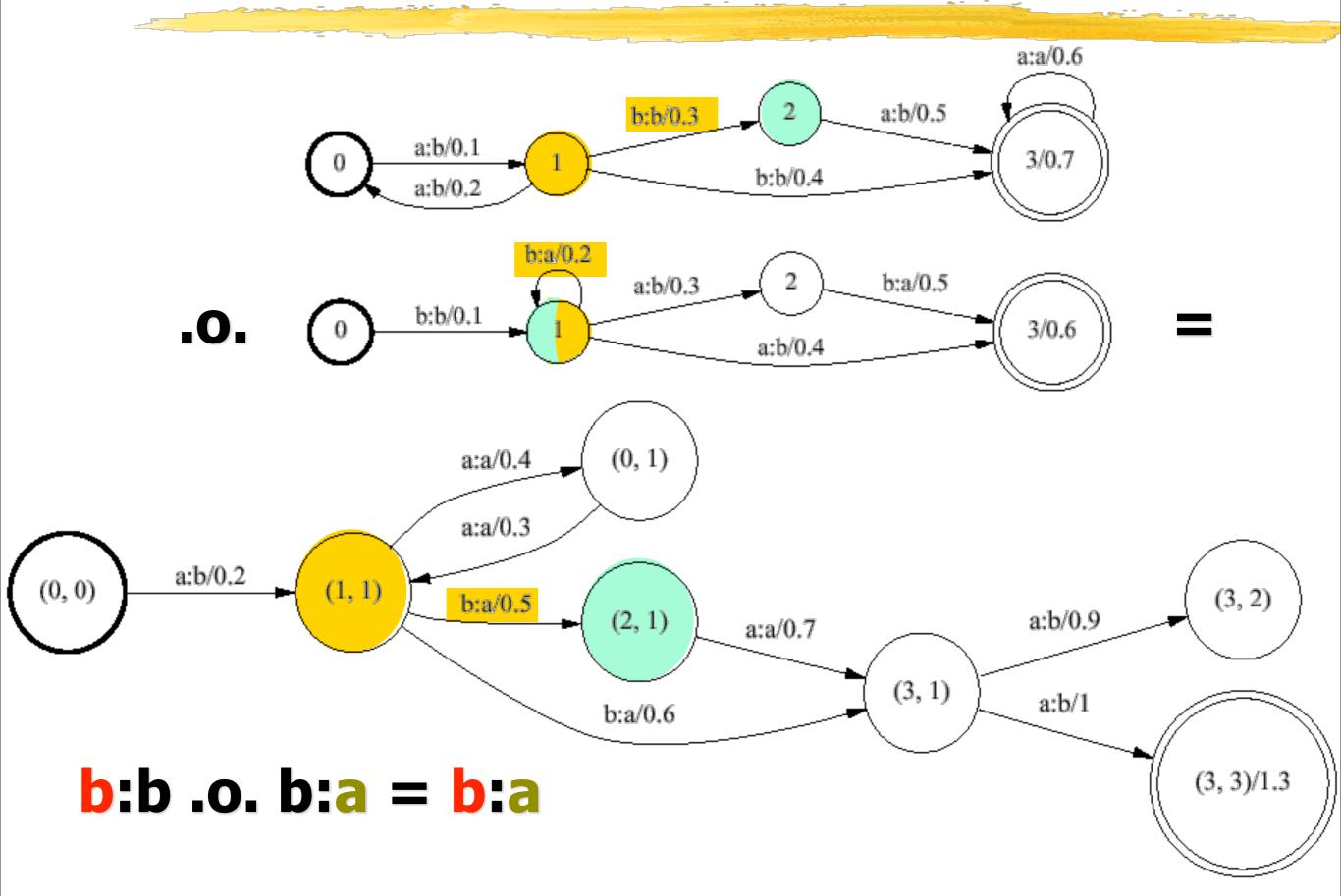


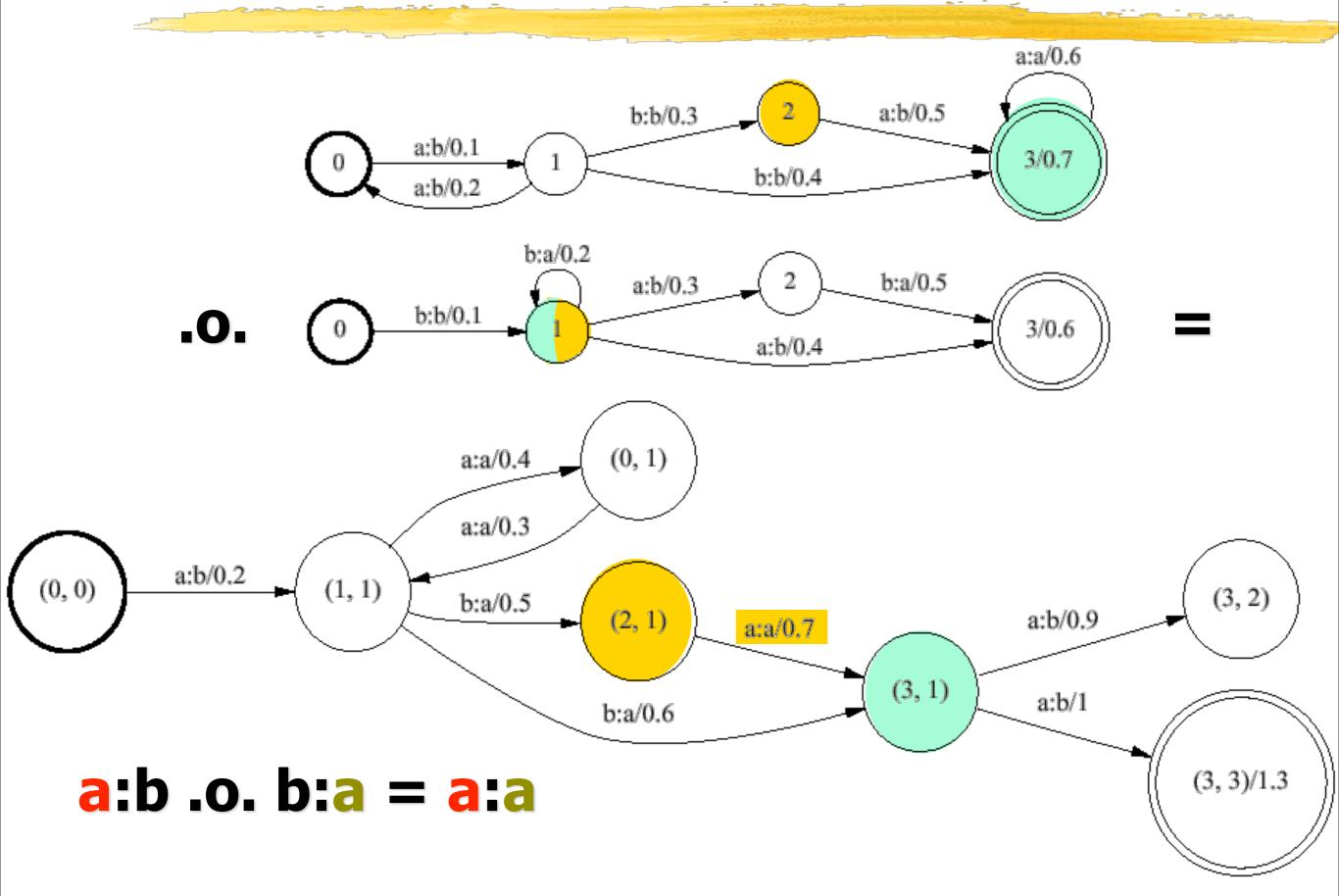


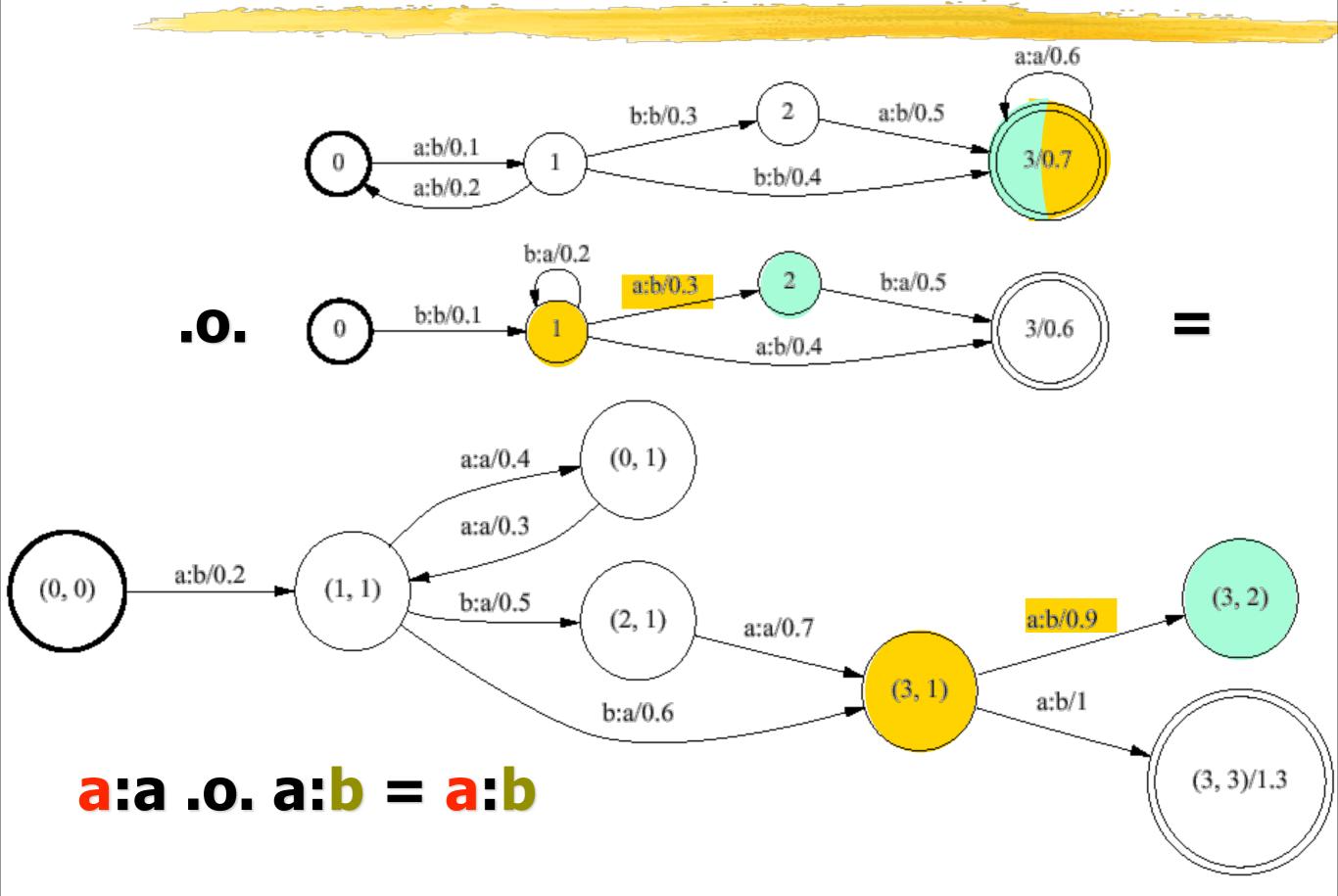


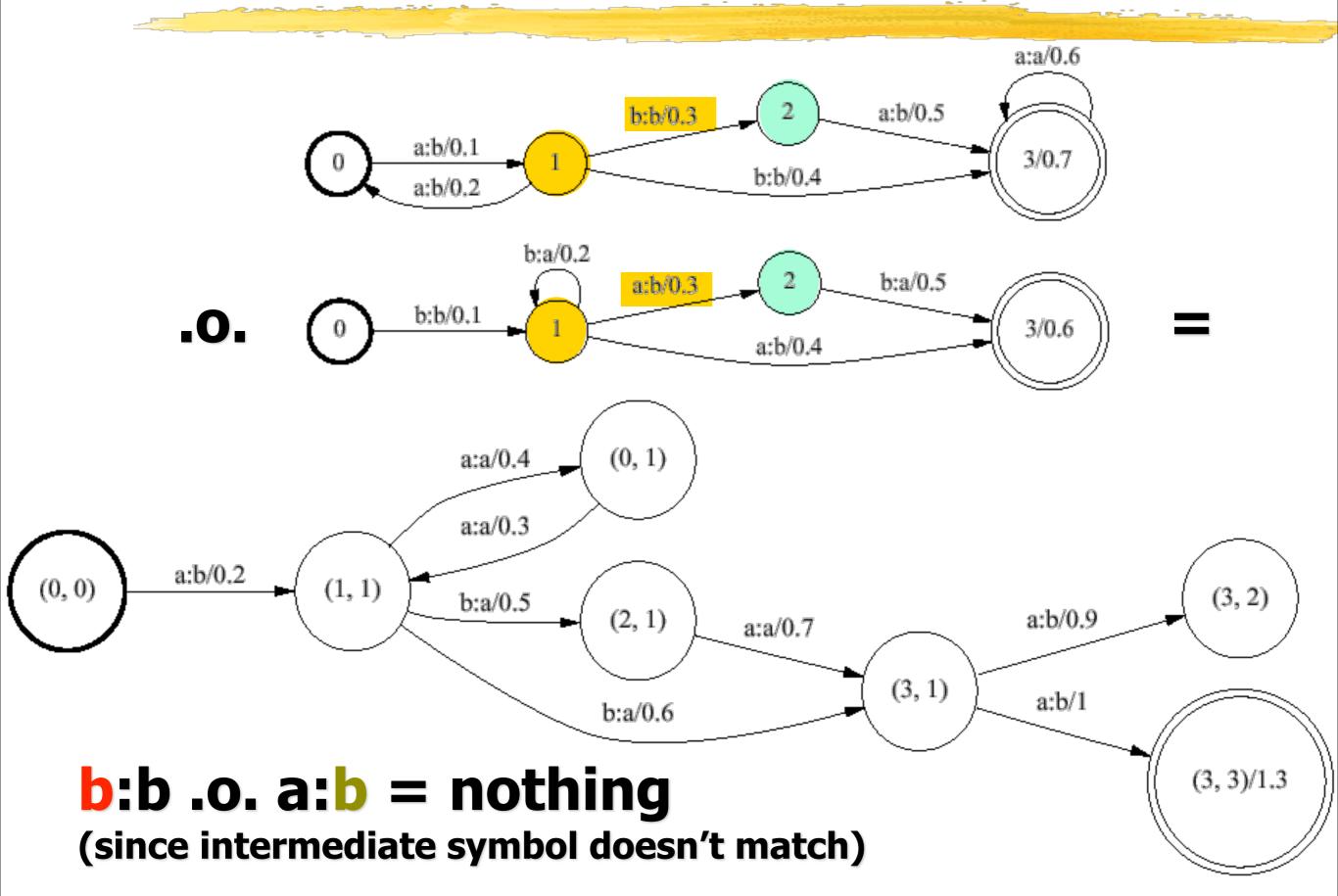


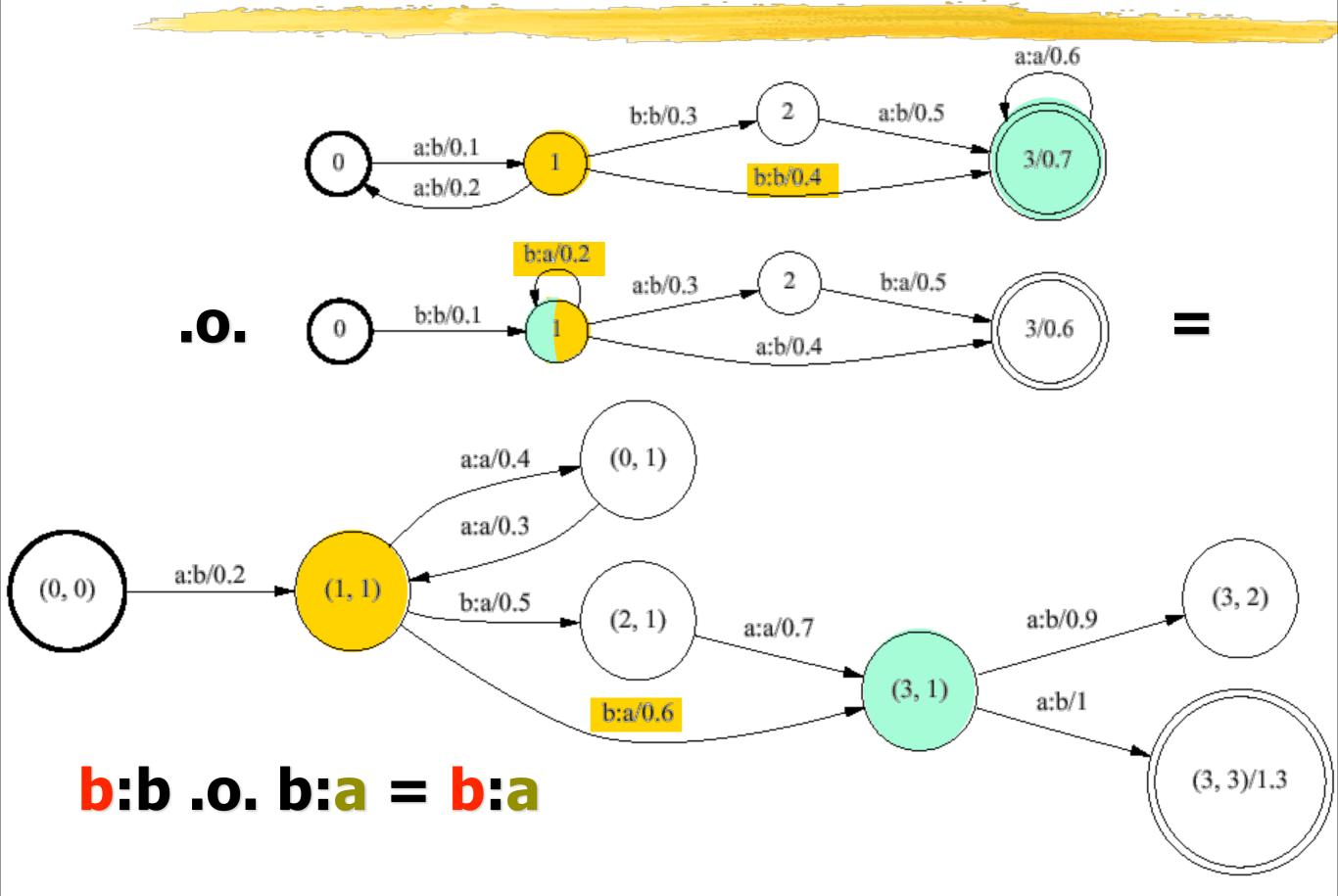


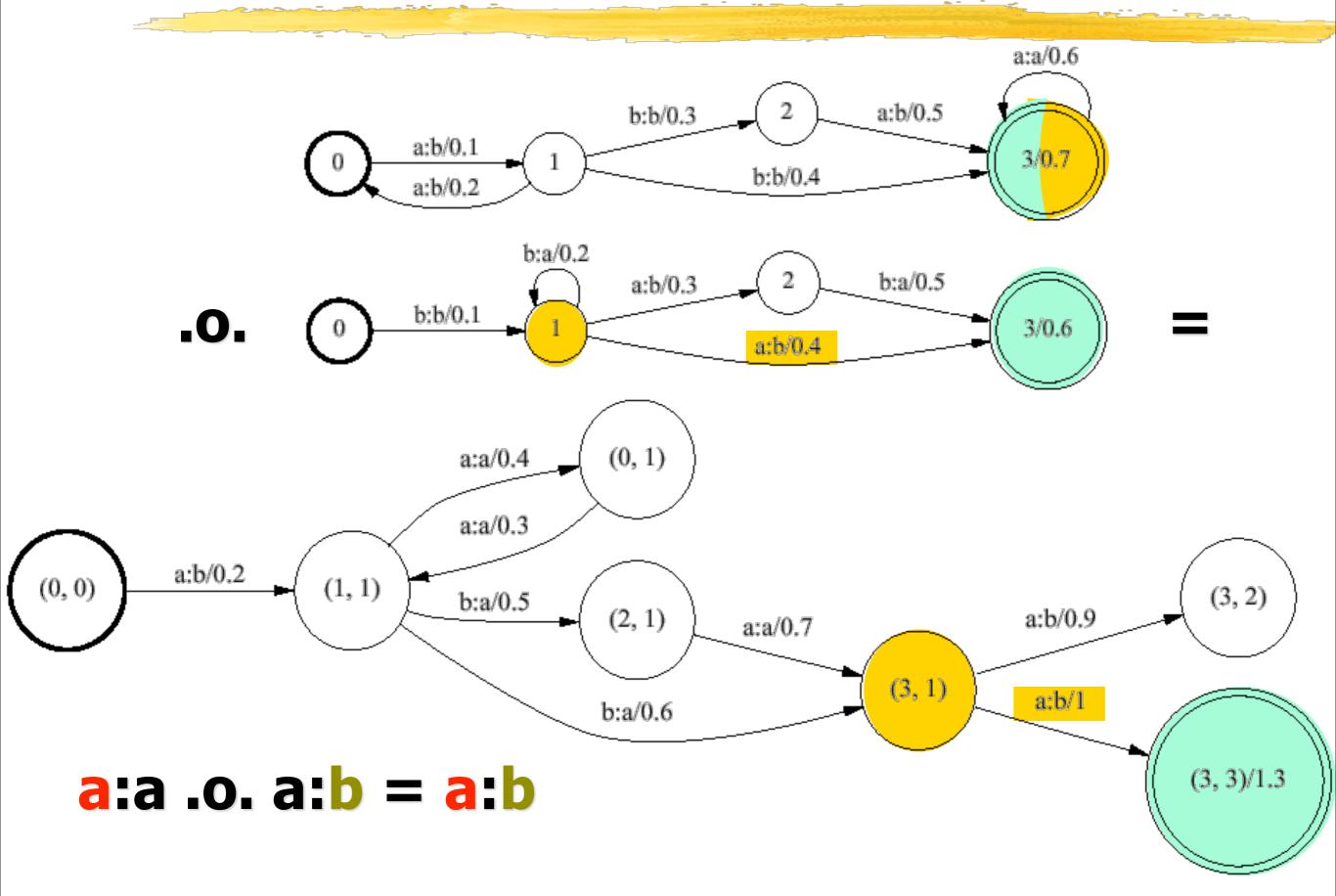




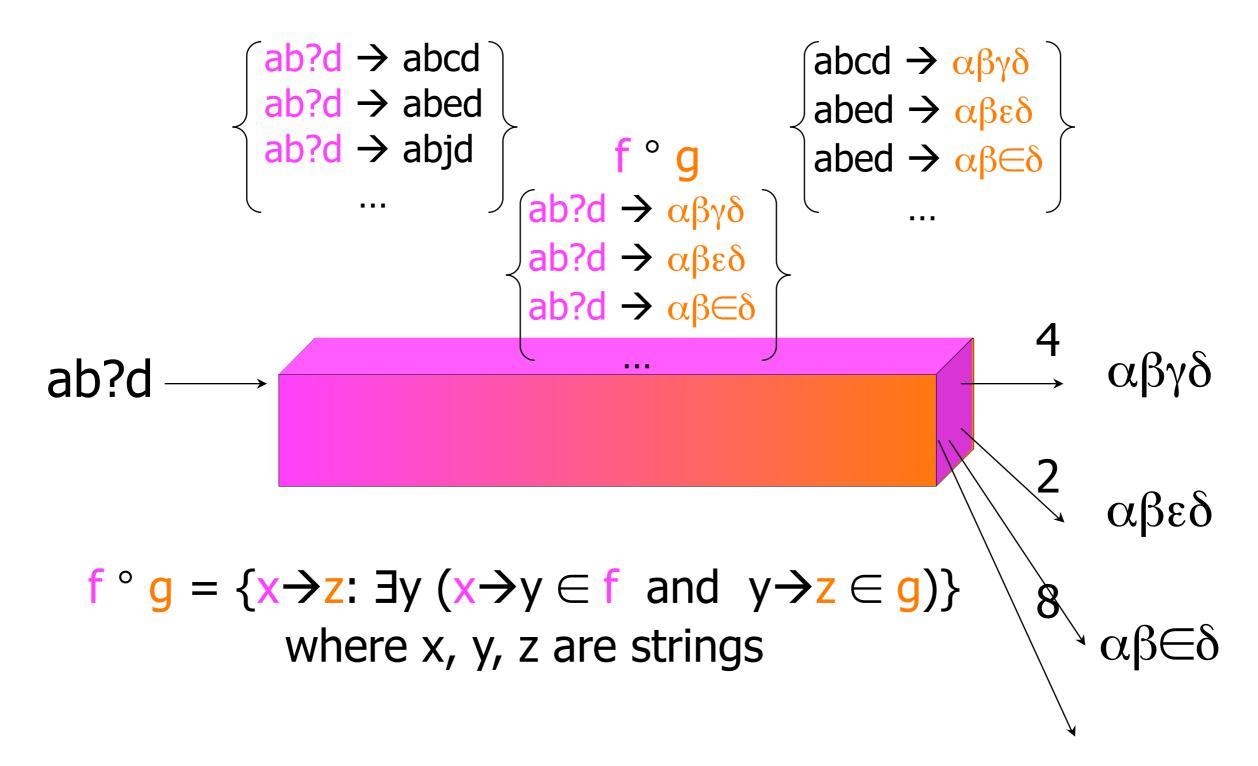








#### **Relation = set of pairs**



We've defined A .o. B where both are FSTs

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA
- Two relations (FSTs):

A ° B = { $x \rightarrow z$ :  $\exists y (x \rightarrow y \in A \text{ and } y \rightarrow z \in B)$ }

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA
- Two relations (FSTs):

A ° B = { $x \rightarrow z$ :  $\exists y (x \rightarrow y \in A \text{ and } y \rightarrow z \in B)$ }

Set and relation:

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA
- Two relations (FSTs): A  $\circ$  B = {x $\rightarrow$ z:  $\exists$ y (x $\rightarrow$ y  $\in$  A and y $\rightarrow$ z  $\in$  B)}
- Set and relation:
  - A ° B = { $x \rightarrow z$ :  $x \in A \text{ and } x \rightarrow z \in B$  }

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA
- Two relations (FSTs): A  $\circ$  B = {x $\rightarrow$ z:  $\exists$ y (x $\rightarrow$ y  $\in$  A and y $\rightarrow$ z  $\in$  B)}
- Set and relation:
  - A ° B = { $x \rightarrow z$ :  $x \in A \text{ and } x \rightarrow z \in B$  }
- Relation and set:

 $A \circ B = \{x \rightarrow z: x \rightarrow z \in A \text{ and } z \in B\}$ 

- We've defined A .o. B where both are FSTs
- Now extend definition to allow one to be a FSA
- Two relations (FSTs): A  $\circ$  B = {x $\rightarrow$ z:  $\exists$ y (x $\rightarrow$ y  $\in$  A and y $\rightarrow$ z  $\in$  B)}
- Set and relation:
  - $A \circ B = \{x \rightarrow z: x \in A \text{ and } x \rightarrow z \in B\}$
- Relation and set:

 $A \circ B = \{x \rightarrow z: x \rightarrow z \in A \text{ and } z \in B\}$ 

• Two sets (acceptors) – same as intersection:
A ° B = {x:
x  $\in$  A and
x  $\in$  B }

#### **Composition and Coercion**

- Really just treats a set as identity relation on set {abc, pqr, ...} = {abc→abc, pqr→pqr, ...}
  Two relations (FSTs): A ° B = {x→z: ∃y (x→y ∈ A and y→z ∈ B)}
- Set and relation is now special case (if ∃y then y=x):
  - $A \circ B = \{x \rightarrow z: x \rightarrow x \in A \text{ and } x \rightarrow z \in B\}$
- Relation and set is now special case (if ∃y then y=z):
- $A \circ B = \{x \rightarrow z: x \rightarrow z \in A \text{ and } z \rightarrow z \in B \}$

• Two sets (acceptors) is now special case:  $A \circ B = \{x \rightarrow x: x \rightarrow x \in A \text{ and } x \rightarrow x \in B \}$ 

Feed string into Greek transducer:

#### Feed string into Greek transducer:

• {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- {abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- {abed} .o. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- $[{abed} .0. Greek] .I = {\alpha\beta\epsilon\delta, \alpha\beta\in\delta}$

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- {abed} .o. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- $[{abed} .0. Greek] . I = {\alpha\beta\epsilon\delta, \alpha\beta\in\delta}$

#### Feed several strings in parallel:

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- {abed} .o. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- [{abed} .0. Greek].I = { $\alpha\beta\epsilon\delta$ ,  $\alpha\beta\in\delta$ }

#### Feed several strings in parallel:

{abcd, abed} .o. Greek

= {abcd  $\rightarrow \alpha\beta\gamma\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ }

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- {abed} .o. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- [{abed} .0. Greek].I = { $\alpha\beta\epsilon\delta$ ,  $\alpha\beta\in\delta$ }

#### Feed several strings in parallel:

{abcd, abed} .o. Greek

= {abcd  $\rightarrow \alpha\beta\gamma\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

[{abcd,abed} .0. Greek].| = {αβγδ, αβεδ, αβ∈δ}

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ }
- {abed} .o. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- [{abed} .0. Greek].I = { $\alpha\beta\epsilon\delta$ ,  $\alpha\beta\in\delta$ }
- Feed several strings in parallel:
  - {abcd, abed} .o. Greek

= {abcd  $\rightarrow \alpha\beta\gamma\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

- $[{abcd,abed} .0. Greek] | = {\alpha\beta\gamma\delta, \alpha\beta\epsilon\delta, \alpha\beta\in\delta}$
- Filter result via No $\varepsilon = \{\alpha\beta\gamma\delta, \alpha\beta\in\delta, ...\}$

#### Feed string into Greek transducer:

- {abed  $\rightarrow$  abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ }
- {abed} .0. Greek = {abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }
- [{abed} .0. Greek].I = { $\alpha\beta\epsilon\delta$ ,  $\alpha\beta\in\delta$ }
- Feed several strings in parallel:
  - {abcd, abed} .o. Greek

= {abcd  $\rightarrow \alpha\beta\gamma\delta$ , abed  $\rightarrow \alpha\beta\epsilon\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

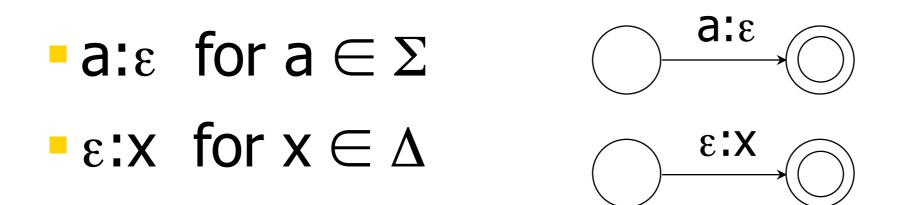
- $[{abcd,abed} .0. Greek] | = {\alpha\beta\gamma\delta, \alpha\beta\epsilon\delta, \alpha\beta\in\delta}$
- Filter result via No $\varepsilon = \{\alpha\beta\gamma\delta, \alpha\beta\in\delta, ...\}$

- {abcd,abed} .o. Greek .o. No $\epsilon$ 

= {abcd  $\rightarrow \alpha\beta\gamma\delta$ , abed  $\rightarrow \alpha\beta\in\delta$ }

# What are the "basic" transducers?

- The operations on the previous slides combine transducers into bigger ones
- But where do we start?



• Q: Do we also need a:x? How about  $\varepsilon:\varepsilon$ ?

## Reading

 Okan Kolak, William Byrne, and Philip Resnik. A Generative Probabilistic OCR Model for NLP Applications. In *HLT-NAACL*, 2003.

<u>http://aclweb.org/anthology-new/N/N03/</u> <u>N03-1018.pdf</u>