



- Solution: Log-sum-exp SR (backward algorithm) resp. arctic SR (Viterbi algorithm). Note: The log-sum-exp SR returns  $\log(Z)$  instead of  $Z$
- **Scoring functions** — **Hidden Markov Model**:  
 $\text{score}(t_{N-1}, t_N, w) =$   
 $\text{transition}(t_{N-1}, t_N) + \text{emission}(t_N, w_N) =$   
 $\text{transition probability (tag-tag pairs)} + \text{emission probability (word-tag pairs)}$  except for  $t_N = \text{EOS}$ , where emission probability = 0

## 18 Syntax — Syntactic Description

- **Constituent**: Coherent unit resp. nodes in tree
- **Terminals**: Words, *non-terminals*: Abstractions over words (e.g. NP, VP, etc.)
- **Grammar**: Production rules
- Can be represented as *parse tree*

### Formulation Probabilistic CFGs

- $(N, S, \Sigma, \mathcal{R}, P)$  —
- Non-terminals  $N = \{N_1, N_2, \dots\}$
  - Start non-terminal  $S$
  - Terminals  $\Sigma = \{a_1, a_2, \dots\}$
  - Production rules  $N \rightarrow N\alpha$ , where  $N$  is non-terminal and  $\alpha \in (N \cup \Sigma)^*$
  - Probabilities  $\mathcal{P}$  for each rule, locally normalized over each transition:  
 $\sum_k p(N \rightarrow \alpha_k) = 1$  where  $N \rightarrow \alpha_1, \dots, N \rightarrow \alpha_k$  are expansions of node  $N$
  - String  $s$  of length  $M$

- Probability of a tree:  $p(t) = \prod_{r \in t} p(r) = p(S \rightarrow S_1 S_2) M^{-1} \times p(S \rightarrow X) M \times p(X \rightarrow \sigma) M$  where:
- $p(S \rightarrow S_1 S_2) M^{-1}$ : repeats the rule  $S \rightarrow SS$   $M-1$  times to obtain  $M$  starting nodes
  - $X$ : non-terminal,  $\sigma$ : terminal

Probability of a string:  $p(s) = \sum_{t \in T(s)} p(t)$

- Chomsky Normal Form (CNF)** —
- Production rules:  $N_1 \rightarrow N_2 N_3$  are non-terminal productions,  $N \rightarrow a$  are terminal productions,  $S \rightarrow \epsilon$
  - Prohibits cyclic rules
  - Transform CFG to CNF:
    - Remove  $A \rightarrow \epsilon$ , e.g. if CFG contains  $S \rightarrow AB|E$ ,  $B \rightarrow b|e$ , change to  $S \rightarrow AB|A|e$ ,  $B \rightarrow b$
    - Remove  $A \rightarrow B$ , e.g. if CFG contains  $S \rightarrow B$ ,  $B \rightarrow b$ , change to  $S \rightarrow b$
    - Convert long productions  $A \rightarrow B_1 \dots B_k$  with  $k > 2$  and mixed rules  $A \rightarrow aB$  with intermediate non-terminals

- Weighted CFGs** —
- A more general formulation of PCFGs, globally normalized

- $p(t) = \frac{\prod_{r \in t} \exp(\text{score}(r))}{\sum_{t' \in T} \prod_{r' \in t'} \exp(\text{score}(r'))}$
- Challenge:  $Z$  is infinitely large ( $\Sigma^*$ )
- Solution:  $p(t|s) = \frac{\prod_{r \in t} \exp(\text{score}(r))}{\sum_{t' \in T(s)} \prod_{r' \in t'} \exp(\text{score}(r'))}$
- Challenge:  $Z$  is still potentially infinitely large (cycles)
- Solution: Revert to CNF. Then,  $|Z|$  is the number of rooted binary trees, i.e. Catalan number  $C_{M-1}$
- Then:  $p(t|s) = \frac{\prod_{N_i \rightarrow N_j N_k, e \in t} \exp(\text{score}(N_i \rightarrow N_j N_k))}{\sum_{t' \in T(s)} \prod_{N_i \rightarrow N_j N_k, e \in t'} \exp(\text{score}(N_i \rightarrow N_j N_k))}$
- $\prod_{N_i \rightarrow a, e \in t} \exp(\text{score}(N_i \rightarrow a)) \times \prod_{N_i \rightarrow a, e \in t} \exp(\text{score}(N_i \rightarrow a))$
- consisting of non-terminal productions  $\times$  terminal productions

- Span** —
- **Admissible** if  $X \rightarrow w[i:j]$
  - For tree where nodes expand along right diagonal of tree only, only 1 admissible span per span size:  $[M - \text{span size} + 1, M + 1]$

- Optimization CKY Algorithm** —
- Can: Compute  $Z$  (inside SR or log-sum-exp SR), find best parse and its probability (Viterbi or arctic SR), determine if a given string is admissible by the grammar (Boolean SR)
  - Requires grammar in CNF
  - Chart  $[i, j, X]$  is probability that non-terminal  $X$  generates the subtree resp. substring  $s_i, \dots, s_j$

- For  $m=1, \dots, M$ : Terminal productions: For  $X \rightarrow s_m$ :  
 Chart  $[m, m+1, X] \leftarrow \exp(\text{score}(X \rightarrow s_m))$

- For span  $= 2, \dots, M$ : For  $i=1, \dots, M - \text{span} + 1$ :  
 $k \leftarrow i + \text{span} - 1$  For  $j=i+1, \dots, k-1$ :  
 Non-terminal productions: For  $X \rightarrow YZ$ :  
 Chart  $[i, k, X] \oplus \exp(\text{score}(X \rightarrow YZ)) \otimes \text{Chart}[i, j, Y] \otimes \text{Chart}[j, k, Z]$
- Return Chart  $[1, M+1, S]$

Runtime:  $O(M^3 |R|)$  (for tree where nodes expand along right diagonal of tree only),  
 $i, m = \text{span size} + 1, j = i + 1$ , then  $O(M |R|)$

**Pumping lemmas** Shows that languages that require strict equality between counts of symbols cannot be generated by a CFG:

## 19 Syntax — Dependency Formulation

String  $w$

- Spanning tree**:
- $N$  nodes + 1 root node (words in sentence  $|w|$  or, if there are dependency parsing rules available, items in rules)
  - $N-1$  edges, of which 1 edge is fixed (root)
  - Directed spanning trees:  $N(N-1)$
  - Directed, labeled spanning trees:  $N(N-1) \times n$  where  $x$  is given by the number of the destination nodes (e.g.  $N$ ) raised to the power of the number of edges that can vary (e.g.  $N-2$ )

- Probability of spanning tree** —  $\frac{\exp(\text{score}(t, w))}{\sum_{t'} \exp(\text{score}(t', w))}$
- Challenge: Naively:  $(N-1) N^{N-2}$  spanning trees with one root
  - Solution: **Edge factored assumption**  $p(t|w) = \frac{\prod_{(i,j) \in t} \exp(\text{score}(i, j, w)) \exp(\text{score}(r, i))}{\sum_{t'} \prod_{(i,j) \in t'} \exp(\text{score}(i, j, w)) \exp(\text{score}(r, i))}$  where  $(i \rightarrow j)$  is an edge,  $r$  is the root
  - **Matrix-Tree Theorem** — Counts number of directed spanning trees in  $O(N^3)$  time:

- **Adjacency matrix**: One entry for each node  $i$  to node  $j$ , if they were connected via an edge  $i \rightarrow j$   $A_{ij} = \exp(\text{score}(i, j, w))$  with 0 on diagonal and non-null on off-diagonal
- **Root vector**: One entry for each node  $j$ , if it were the root node  $\rho_j = \exp(\text{score}(j, w))$
- Construct **Laplacian matrix**, accounting for constraint that there is only one root node:

$$L_{ij} = \begin{cases} \rho_i & \text{if } i=1 \\ \sum_{i'=1, i' \neq j} A_{i'j} & \text{if } i=j \\ -A_{ij} & \text{otherwise} \end{cases} \quad \text{i.e.}$$

- $(-1)^{\text{st}}$  row of  $L$  contains root scores
- $(-1)$ -diagonal of  $L$  (except  $L_{1,1}$ ) contains sum within each column of  $A$  (except  $A_{i,i}$ )
- off-diagonal of  $L$  (except  $1^{\text{st}}$  row) contains elements of  $A$  with  $\times$  with  $-1$
- According to matrix tree theorem:  $|L| = \det(L) = Z = \text{number of trees in graph}$

- Optimization CLE algorithm** in  $O(N^2)$  time:
- Greedy algorithm selects best incoming edge for each node, except the root
  - This can cause cycles, which we contract
  - Break cycle: For each enter edge, break cycle by removing edges that are also incoming at the node where the enter edge is incoming
  - Re-weight: To enter edge, add weights of remaining edges that are strictly on (not in) the cycle

- If there are multiple edges emanating from the root: For each root edge, calculate cost of deleting this edge: Cost = Weight of root edge - weight of next-best incoming edge to target node (treating cycles as single nodes)
- Preliminarily remove edge with lowest cost, but keep target node intact
- Preliminarily repeat steps 3,6 as needed
- Re-run greedy algorithm in contracted form
- If this leads to a cycle: Undo removal, contract (treating cycles as single nodes), then re-expand, Otherwise: Re-expand (SR)

## 20 Semantic Parsing

- Lambda Calculus** Basic components:
- Logical constants: Objects vs. relations
  - Variables: Undetermined logical constants, free vs. bound
  - Objects  $x, y, z, \dots$  in  $\lambda x.f(x)$
  - Relations  $P, Q, R, \dots$  in  $\lambda P.P(\dots)$

- Relation has *arity* which is number of objects it relates, e.g.  $P(x, y)$  has arity 2
- **Literals**: Formed by applying relations to objects, e.g. likes(ALEX, y),  $P(x, y)$
- Abstraction:
  - Let  $M, N$  be terms and  $x$  be a variable
  - $\lambda x$  is an *abstraction*
  - $\lambda x.M$  is a function with input  $x$  and *scope* of  $\lambda$  is  $M$ . It replaces every free occurrence of  $x$  in  $M$  with whatever the function is applied to
  - $\lambda x.MN$  is a function applied to  $N$
  - Output of  $\lambda x.MN$ :  $M[x := N]$
  - Note:
    - In the output of  $\lambda x.MN$ , only  $M[x := N]$  remains (for relations  $P(x)$ , the relation  $P$  is exchanged, but the arguments  $x$  remain), and  $\lambda x$  and  $N$  disappear
    - If  $M$  doesn't contain variable, it is returned as is, e.g.  $\lambda x.y.xN$  returns  $y$
    - Work from outside in
    - Determine what is applied to what based on what is the argument in the CCG expression  $X|Y$

Application:  $M, N \rightarrow (MN)$  Parentheses are left-associative, e.g.  $(\lambda xy.x \lambda xy.y \lambda xy.y) = (((\lambda xy.x) \lambda xy.y) \lambda xy.y)$

- **Alpha conversion** — Renaming
- We can rename a variable in an abstraction together with all its occurrences in the scope of the abstraction, which are bound to the same abstraction
- Note: Only variables are renamed, not functions  $\lambda x.M$
- If the renamed variable remains bound to the same abstraction and the remaining variables remain free resp. bound as before, the renaming is valid
- **Beta reduction** — Application We can apply one **lambda** term to another if the free variables in  $N$  remain free in  $M[x := N]$

**LIG resp. CCG LIG**  $(N, S, I, \Sigma, \mathcal{R})$  —

- $N$ : Non-terminals  $N_1, N_2, N_3, \dots$
- $S$ : Start non-terminal
- $I$ : Indices  $f, g, h, \dots$
- $\Sigma$ : Alphabet of terminals  $a_1, a_2, a_3, \dots$
- $\mathcal{R}$ : Production rules:  $N[i\sigma] \rightarrow \alpha M[\sigma]\beta$ ,  $N[i\sigma] \rightarrow \alpha M[f\sigma]\beta$ ,  $N[f\sigma] \rightarrow \alpha M[\sigma]\beta$  where  $N$  and  $M$  are non-terminals,  $\alpha$  and  $\beta$  are sequences of terminals and non-terminals,  $\sigma$  is associated with the non-terminal on the LHS and passed to exactly one non-terminal on RHS, base symbol  $\sigma$  is usually combined with a marker  $f$ , where markers encode a count and can be popped from (rule 3) or pushed to (rule 2) the stack

- CCG**  $(V, \mathcal{V}, \Sigma, \mathcal{R}, f, R)$
- $N$ : Non-terminals
  - $\mathcal{V}$ : Non-terminals
  - $\mathcal{V}_T$ : Terminals
  - $\mathcal{V}_P$ : Production rules
  - $f$ :  $\mathcal{V}_T \rightarrow \mathcal{C}(\mathcal{V}_N)$

- Categories:
- **Atomic categories**: Terminals
  - **Complex categories**:
    - Built from atomic categories via operators
    - Function with pattern:  $X|Y \quad Y \rightarrow X$ , where  $X|Y$  is function, with  $|$  being an operator,  $Y$  being argument, and  $X$  being output, and applying function to  $Y$  yields  $X$
    - Have *arity* which is number of arguments it relates, e.g.  $X|Y|Z$  has arity 2
  - $X = A_1 m X_m \dots | \cdot | X_1$  where
    - $A_i$ : Atomic category, output of  $X$
    - $X_1, \dots, X_m$ : Arbitrary categories, arguments of  $X$
    - $m$  = Arity of  $X$

- Operators: Backward and forward slash
- Note: Operators are read from outside in

- Rules: Specify how categories can be combined into other categories
- **Function application**:
    - Forward  $(>)$ :  $X: X|Y \quad Y \rightarrow X$
    - Backward  $(<)$ :  $Y \quad X|Y \rightarrow X$
    - CCGs that only have application rules, have power of CFG
  - **Function composition**: Forward  $((B >)$ :  $(X|Y) \quad (Y|Z) \rightarrow (X|Z)$ , etc.
  - **Higher-order rules**: Forward  $(>^n)$ :  $X|Y \quad Y|n Y_n \dots | Y_1 \rightarrow X|n Y_n \dots | Y_1$ , etc.
  - **Type raising**: Forward  $(T >)$ :  $X \rightarrow T/(T|X)$ , etc.

- $2^{|V|N}$  by forward and  $2^{|V|N}$  backward rules, but infinitely many rule instances
- Lexicon: Associates terminals with categories
  - E.g. Atomic: Harry =  $NP$
  - E.g. Complex: walks =  $(S|NP)$
- CCG parsing:  $\frac{A_1}{B} \dots \frac{A_k}{B}$  where  $B$  is a consequence of  $A_1, \dots, A_k$  CKY-style parsing algorithm:
  - **Axioms**:  $[X, i, i+1]$  where  $w_{i+1} = X$  is a lexicon entry
  - **Inference rules**:  

$$\frac{[X|Y, i, j] \quad [Y, j, k]}{[X, Y, i, k]} \quad X|Y \quad Y\beta \rightarrow X\beta$$
 resp.  

$$\frac{[Y, j, i, j] \quad [X|Y, j, k]}{[Y, j, i, k]} \quad Y\beta \quad X|Y \rightarrow X\beta$$
  - Polynomial time algorithm:
    - Arity bounded by *grammar constant*  $C_g$ :  $\text{ar}(X) \leq C_g$
    - Challenge: Categories with arity  $> C_g$  can no longer be derived
    - Solution: Decompose longer derivations into smaller pieces: **Derivation contexts**  $c$
    - New inference rules:
      - 1) If the composition of two categories  $(X|Y \quad Y\beta)$  would result in a context item with arity  $\text{ar}(X\beta) > C_g$ : Context item:  $\frac{[X|Y, i, j] \quad [Y, j, k]}{[Y, \beta, i, j, k]}$  where  $\beta$  are all non-terminals and edge operators, except for  $Y$
      - 2.1) If the further composition of a context item with another category  $(X|Y \quad Y\beta)$  would result in a category with arity  $\text{ar}(X\beta) \leq C_g$ : Context item can be recombinced with the original derivation:  $\frac{[X|Y, i', j'] \quad [Y, \beta, i, j', j']}{[X, Y, i', j', k]}$
      - 2.2.1) If the further composition of a context item with another category  $(\beta|Z \quad Z\gamma)$  would result in a category with arity  $\text{ar}(Y\beta\gamma) \leq C_g$ : Context item can be extended normally:  $\frac{[Y, \beta, Z, i, i', j', j'] \quad [Z, \gamma, j, k]}{[Y, \beta, \gamma, i, i', j', k]}$  where  $\beta\gamma$  refers to  $\beta$  without  $|Z|$  (can be  $\epsilon$ )
      - 2.2.2) Otherwise: A new context is derived from the previous context:  $\frac{[Y, \beta, Z, i, i', j', j'] \quad [Z, \gamma, j, k]}{[Z, \gamma, i, i, j, k]}$  item
    - 3.1) When the arity is small enough, the new context item can be recombinced with the previous context:  $\frac{[Y, \beta, Z, i'', i', j', j''] \quad [Z, \epsilon, i, i, i', j'']}{[Y, \beta, Z, i, i', j', j']}$

- $C_g \geq \max\{f, a, n\}$ : At least as large as maximal  $|Y|$  (determined by largest arity  $a$  in lexicon) and  $\beta$  (determined by the maximum degree  $n$  of composition rules)
- CCG can be paired with Lambda calculus
- **SK-Calculus**
  - Alternative to Lambda calculus
  - Variables:  $x, y, z, \dots$
  - **Primitive functions resp. combinators**:
    - $S$ :  $Sxyz = (xz(yz)) = ((xz)(yz)) = yx$
    - $K$ : Constant function:  $Kxy = ((Kx)y) = x$
    - $I$ : Identity function:  $Ixx = x$
    - $SKK$  and  $I$  are equivalent:  $S(KK)x = SKKx = Kx(Kx) = x$
  - Parentheses are left-associative, e.g.  $(Kxy)z = (((Kx)y)z)$

- Initial state  $(q_1^i, q_2^i) \in I_1 \times I_2$  and final state  $(q_1^f, q_2^f) \in F_1 \times F_2$
- For  $q_1, q_2 \in Q_1 \times Q_2$ : For  $a, b/w_1 \rightarrow q_1^i, q_2 \xrightarrow{c/d/w_2} q_2^j$   $E_{T_1}(q_1) \times E_{T_2}(q_2)$ : If  $b=c$ ,  $(q_1, q_2)$  and  $q_1^i, q_2^j$  are accessible and co-accessible (reachable from  $(q_1^i, q_2^i)$ , can reach  $(q_1^f, q_2^f)$ ) in  $\mathcal{T}$ : Add new  $(q_1, q_2)$ ,  $(q_1^i, q_2^j)$  and  $(q_1, q_2) \xrightarrow{a.d/w_1 w_2} (q_1^i, q_2^j)$  to  $\mathcal{T}$

- For  $(q_1^i, q_2^j) : \lambda T = \lambda_1(q_1^i) \otimes \lambda_2(q_2^j)$
- Return  $\mathcal{T}$
- Challenge: Runtime  $O(|Q_1| |Q_2|)$
- Solution: **Accessible algorithm**: Construct all possible pairs of initial states and then expand outward, adding accessible states

**Pathsum and Algorithms** Pathsum:

- $Z(\mathcal{A}) = \bigoplus_{\pi \in \Pi(\mathcal{A})} w(\pi)$  Backward algorithm: Computes finite  $Z(\mathcal{A})$  in acyclic WFSA in  $O(|\delta|)$  time:
- For  $q \in \text{Rev-Top}(\mathcal{A})$  (starts with final state):
  - If  $q \in F$ :  $\beta(q) = \rho(q)$   
 Else:  $\beta(q) = \bigoplus_{a/w} a/w \quad w \otimes \beta(q')$   

$$q \xrightarrow{a/w} q'$$
  - Return:  $\bigoplus_{q \in I} \lambda(q') \otimes \beta(q')$

- 21 WFSA**
- FSA**  $\mathcal{A} = (\Sigma, Q, I, F, \delta)$
- Alphabet  $\Sigma$  with  $a, b, c, \dots$
  - States  $Q$ , initial states  $I$ , final states  $F$

- Transitions  $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$  resp.  $a \in \Sigma \cup \{\epsilon\}$   
 $q \xrightarrow{a} q'$
- Sequentially reads individual symbols of an input string  $s$  and transitions from state  $q$  to state  $q'$  upon reading a symbol  $a$  iff  $(q, a, q') \in \delta$
- If, after reading the last symbol, ends up in a state  $q_f \in F$ , automaton *accepts* the string

- WFSA**  $\mathcal{A} = (\Sigma, Q, I, F, \delta, \lambda, \rho)$
- Transitions weighted with SR
  - $I = \{q \in Q | \lambda(q) \neq 0\} \subseteq Q$ ,  $F = \{q \in Q | \rho(q) \neq 0\} \subseteq Q$
  - $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times \mathbb{K} \times Q$  resp.  $a \in \Sigma \cup \{\epsilon\} / w \rightarrow q'$
  - $\lambda: Q \rightarrow \mathbb{K}$  resp.  $\rho: Q \rightarrow \mathbb{K}$  are initial resp. final weighting functions,  $\bar{0}$  for non-initial resp. non-final  $q$

- (W)FSA Terminology** Paths:  $\pi$
- Element of  $\delta^*$ :  $q_1/w_1 \rightarrow q_2, q_2 \dots q_N$
  - $p(\pi) = q_1$  is the beginning and  $q(\pi) = q_N$  the ending state of the path
  - **Length** of is the number of transitions, *yield* the concatenation of symbols
  - **Inner path weight**:  $w_I(\pi) = \bigotimes_{N=1}^n w_n$
  - **Path weight**:  $w(\pi) = \lambda(p(\pi)) \otimes w_I(\pi) \otimes \rho(q(\pi))$
  - A path is *accepting* iff  $w(\pi) \neq \bar{0}$

- Transitions:
- Outgoing from  $q$ :  $E_{\mathcal{A}}(q) = \{a, t, w | (q, a, w, t) \in \delta\}$
  - Incoming to  $q$ :  $E^{-1}(q) = \{a, s, w | (s, w, a, q) \in \delta\}$
  - $\mathcal{A}$  is **unambiguous** iff for string there is maximum one accepting path

- WFST**  $\mathcal{T} = (\Sigma, \Omega, Q, I, F, \delta, \lambda, \rho)$
- $\Sigma$  is the input alphabet,  $\Omega$  is the output alphabet
  - $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Omega \cup \{\epsilon\}) \times \mathbb{K} \times Q$  resp.  $a \in \Sigma \cup \{\epsilon\}; b \in \Omega \cup \{\epsilon\} / w \rightarrow q'$

### WFST Compositions

- $\tau_1 \circ \tau_2$  from  $\mathcal{T}_1 = (\Sigma, \Omega, Q_1, I_1, F_1, \delta_1, \lambda_1, \rho_1)$  and  $\mathcal{T}_2 = (\Omega, \Gamma, Q_2, I_2, F_2, \delta_2, \lambda_2, \rho_2)$ :  $\mathcal{T} = (\Sigma, \Gamma, Q, I, F, \delta, \lambda, \rho)$  such that:  $\mathcal{T}(x, y) = \bigoplus_{z \in \Omega} \mathcal{T}_1(x, z) \otimes \mathcal{T}_2(z, y)$  where  $\mathcal{T}(i, j)$  is weight assigned to mapping from input  $i$  to output  $j$

- Naive algorithm:
- Initial state  $(q_1^i, q_2^i) \in I_1 \times I_2$  and final state  $(q_1^f, q_2^f) \in F_1 \times F_2$
  - For  $q_1, q_2 \in Q_1 \times Q_2$ : For  $a, b/w_1 \rightarrow q_1^i, q_2 \xrightarrow{c/d/w_2} q_2^j$   $E_{T_1}(q_1) \times E_{T_2}(q_2)$ : If  $b=c$ ,  $(q_1, q_2)$  and  $q_1^i, q_2^j$  are accessible and co-accessible (reachable from  $(q_1^i, q_2^i)$ , can reach  $(q_1^f, q_2^f)$ ) in  $\mathcal{T}$ : Add new  $(q_1, q_2)$ ,  $(q_1^i, q_2^j)$  and  $(q_1, q_2) \xrightarrow{a.d/w_1 w_2} (q_1^i, q_2^j)$  to  $\mathcal{T}$
  - For  $(q_1^i, q_2^j) : \lambda T = \lambda_1(q_1^i) \otimes \lambda_2(q_2^j)$
  - Return  $\mathcal{T}$
  - Challenge: Runtime  $O(|Q_1| |Q_2|)$
  - Solution: **Accessible algorithm**: Construct all possible pairs of initial states and then expand outward, adding accessible states

- Pathsum and Algorithms** Pathsum:
- $Z(\mathcal{A}) = \bigoplus_{\pi \in \Pi(\mathcal{A})} w(\pi)$  Backward algorithm: Computes finite  $Z(\mathcal{A})$  in acyclic WFSA in  $O(|\delta|)$  time:
- For  $q \in \text{Rev-Top}(\mathcal{A})$  (starts with final state):
  - If  $q \in F$ :  $\beta(q) = \rho(q)$   
 Else:  $\beta(q) = \bigoplus_{a/w} a/w \quad w \otimes \beta(q')$   

$$q \xrightarrow{a/w} q'$$
  - Return:  $\bigoplus_{q \in I} \lambda(q') \otimes \beta(q')$

Forward: Exchange  $I$  and  $F$  and use  $q' \xrightarrow{a/w} q$

- Lehmann's**: Computes infinite  $Z(\mathcal{A})$  in cyclic WFSAs:
- First WFSAs as a matrix:
- **Adjacency matrix** for all  $\sigma \in \Sigma \cup \{\epsilon\}$ :  $\mathbf{W}(a) \in \mathbb{R}^{|Q| \times |Q|}$ : From-states  $q_n$  in rows, entries as weights  $w$  for  $q_n \xrightarrow{a/w} q_m$  if it exists, else  $\bar{0}$
  - Can be collapsed:  $\mathbf{W} = \bigoplus_{a \in \Sigma \cup \{\epsilon\}} \mathbf{W}(a)$
  - Pathsum for paths of length exactly  $l$ :  $\mathbf{W}^l = \mathbf{W} \otimes \mathbf{W} \otimes \dots$
  - Compute  $Z(\mathcal{A})$  via Kleene closure of matrix  $\mathbf{R}$  over closed SR in  $O(|Q|^3)$  time:

- $\mathbf{R} \in \mathbb{R}^{|Q| \times |Q|}$  with entries as  $\ominus$ -sum of the weights for  $q_i \rightarrow q_k$ :  $\mathbf{R}_{ik} = \bigoplus_{\pi \in \Pi(q_i, q_k)} w_I(\pi)$  where cyclical terms (either single nodes or node combinations) are denoted with  $*$
- $Z(\mathcal{A}) = \bigoplus_{i, k=1}^{|Q|} \lambda(q_i) \otimes \mathbf{R}_{ik} \otimes \rho(q_k)$  where SR retains only pairs  $(i, k)$  where  $i$  is initial and  $k$  is final state  $\rho(q_f)$

- Compute  $\mathbf{R}$ :
- $\mathbf{R}(\bar{0}) = \mathbf{W}$
  - For  $j=1$  up to  $|Q|$ : For  $i=1$  up to  $|Q|$ : For  $k=1$  up to  $|Q|$ :  $\mathbf{R}_{ik}^{j+1} \leftarrow \mathbf{R}_{ik}^{j-1} \otimes (\mathbf{R}_{ij}^{j-1} \otimes \mathbf{R}_{jk}^{j-1})^* \otimes \mathbf{R}_{jk}^{j-1}$

- Return:  $\mathbf{I} \oplus \mathbf{R}(|Q|)$
- Transliteration Aim and approach** —

- Aim: Compute  $p(y|x) = \frac{\exp(\text{score}(x, y))}{Z}$  where  $\text{score}(x, y)$  is the pathsum  $\log \sum_{\pi \in \Pi(y)} w(\pi)$
- Approach:
  - $x \in \Sigma^*$  as  $\mathcal{T}_X$
  - $y$