



UPPSALA
UNIVERSITET

Collins' and Eisner's algorithms

Syntactic analysis (5LN455)

2012-12-16

Joakim Nivre

Department of Linguistics and Philology

Based on slides by Marco Kuhlmann



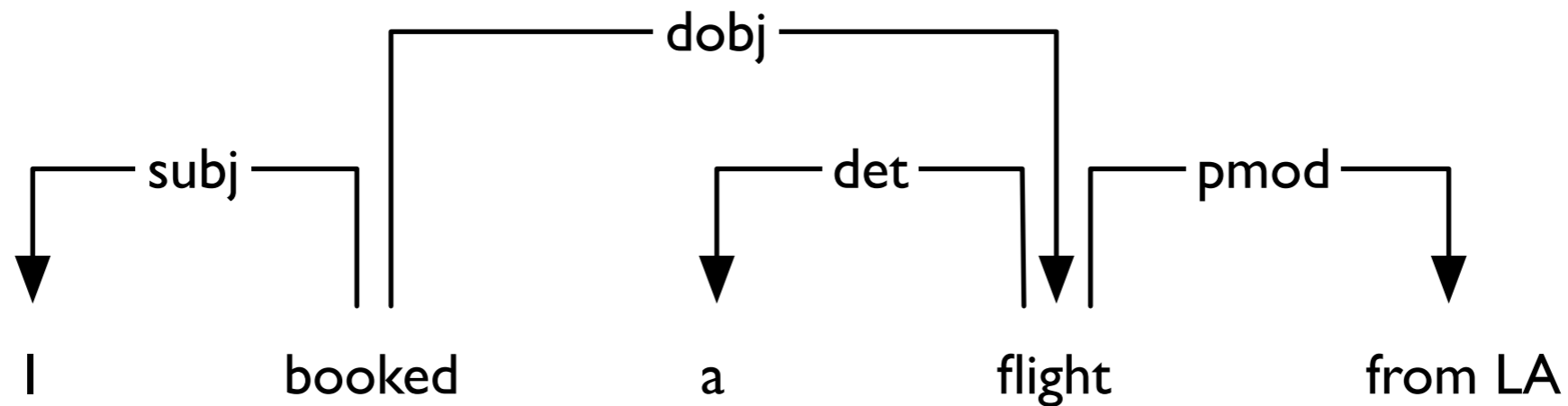


UPPSALA
UNIVERSITET

Collins' algorithm



Recap: Dependency trees



- In an arc $h \rightarrow d$, the word h is called the **head**, and the word d is called the **dependent**.
- The arcs form a rooted tree.



Recap: The arc-factored model

- To score a dependency tree, score the individual arcs, and combine the score into a simple sum.

$$\text{score}(t) = \text{score}(a_1) + \dots + \text{score}(a_n)$$

- Define the **score** of an arc $h \rightarrow d$ as the weighted sum of all features of that arc:

$$\text{score}(h \rightarrow d) = f_1 w_1 + \dots + f_n w_n$$



Collin's algorithm

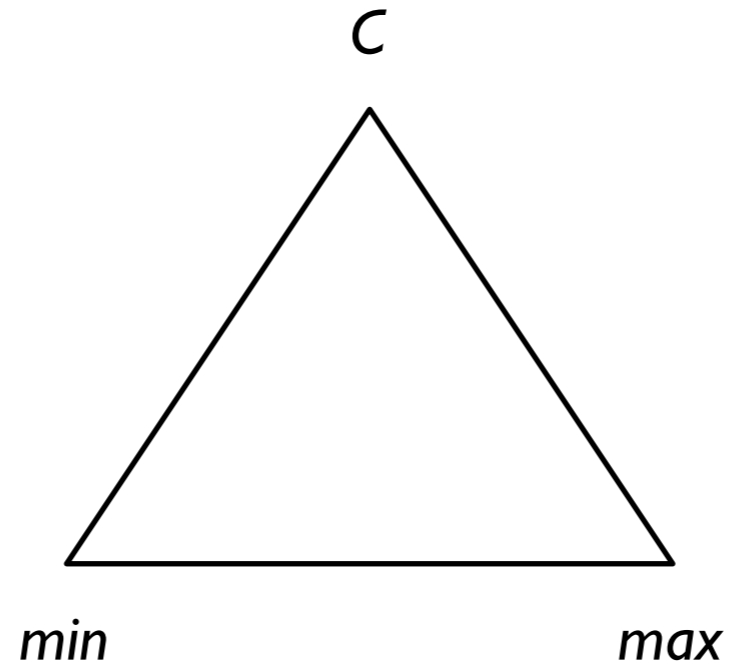
- Collin's algorithm is a simple algorithm for computing the highest-scoring dependency tree under an arc-factored scoring model.
- It can be understood as an extension of the CKY algorithm to dependency parsing.
- Like the CKY algorithm, it can be characterized as a bottom-up algorithm based on dynamic programming.



UPPSALA
UNIVERSITET

Collins' algorithm

Signatures



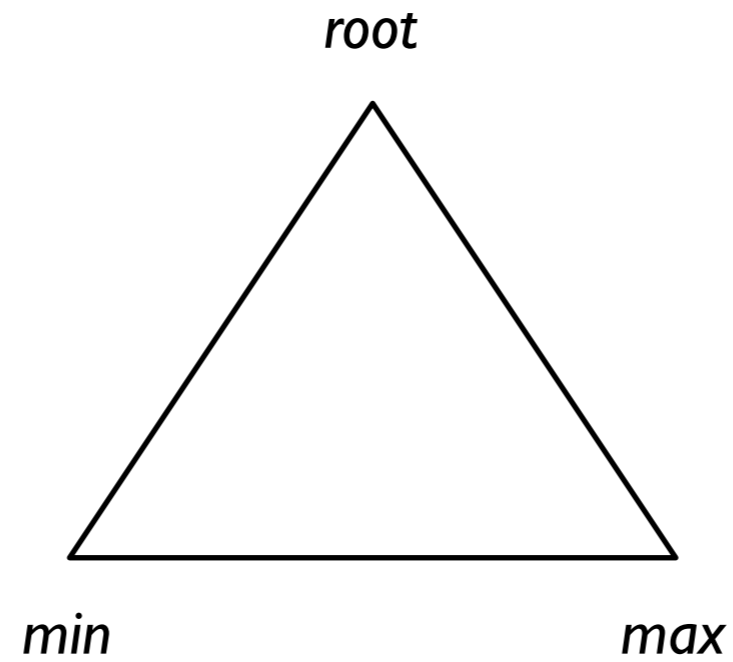
$[min, max, c]$



UPPSALA
UNIVERSITET

Collins' algorithm

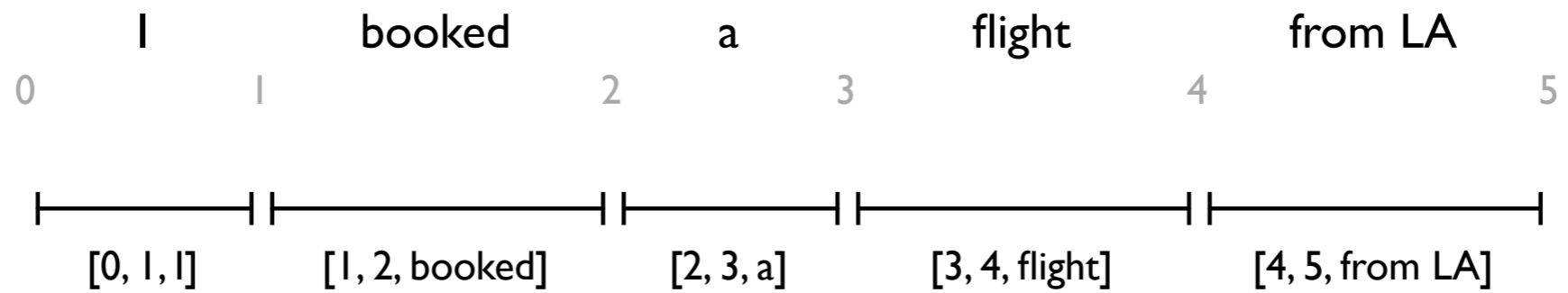
Signatures



[*min, max, root*]

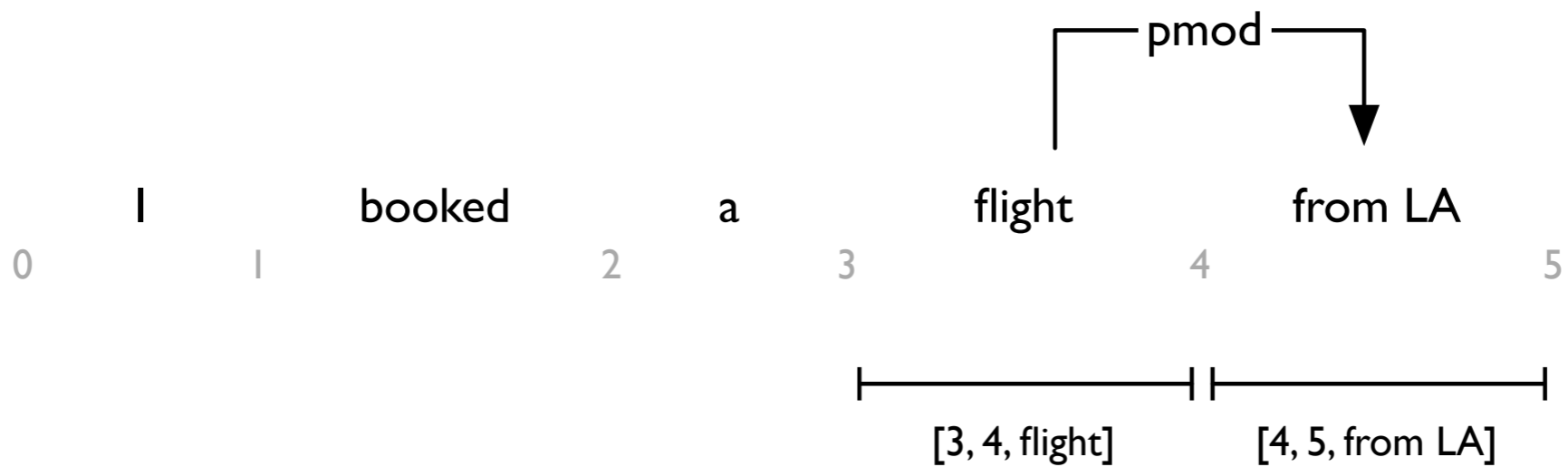


Initialization



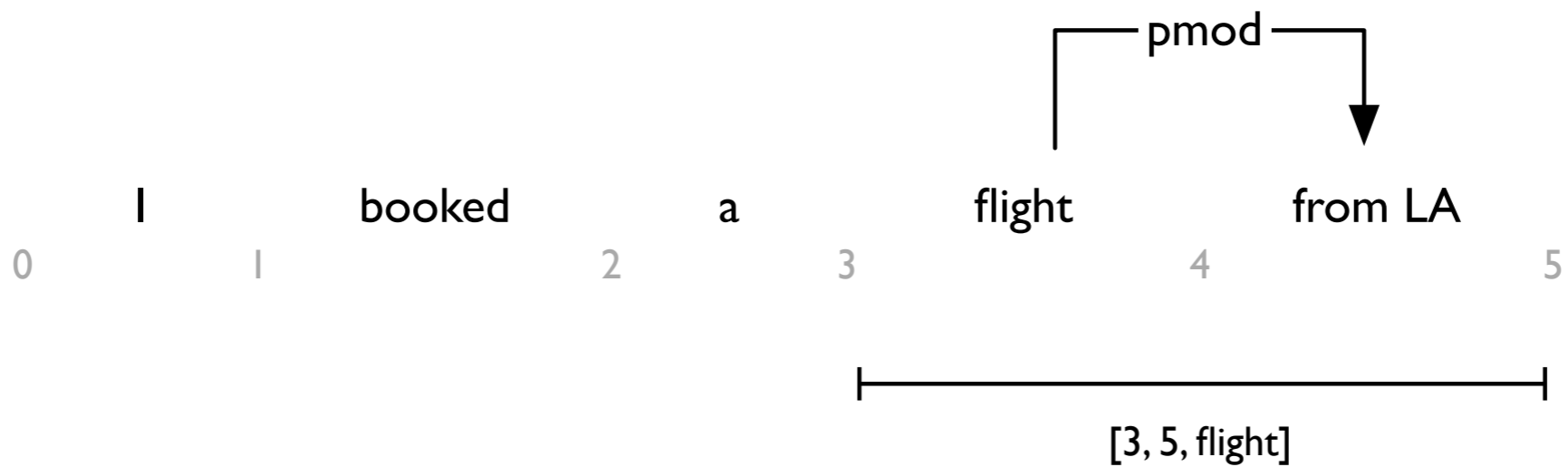


Adding a left-to-right arc





Adding a left-to-right arc

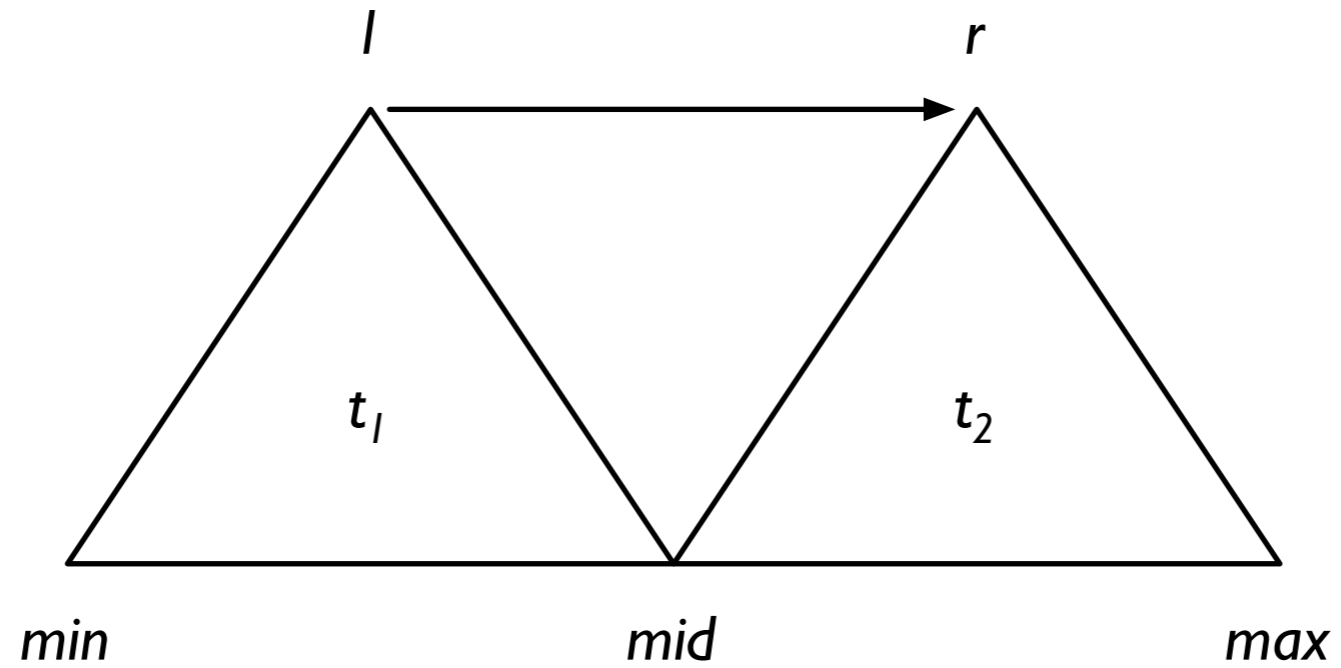




UPPSALA
UNIVERSITET

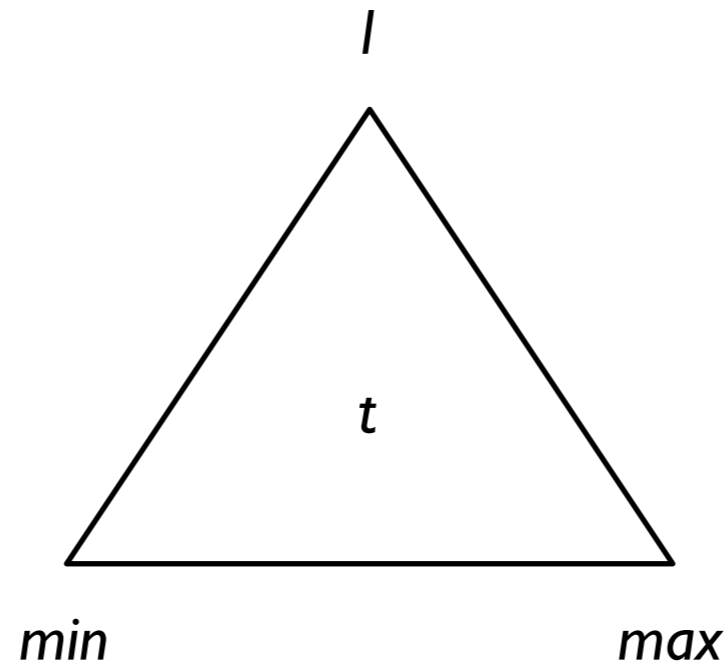
Collins' algorithm

Adding a left-to-right arc





Adding a left-to-right arc



$$\text{score}(t) = \text{score}(t_1) + \text{score}(t_2) + \text{score}(l \rightarrow r)$$

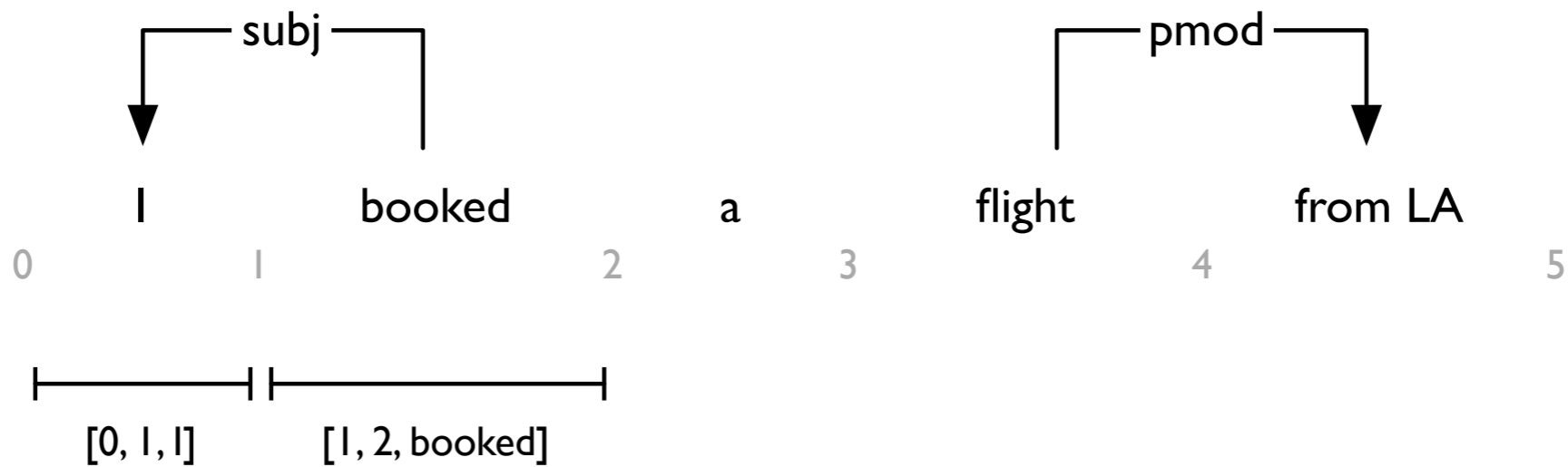


Adding a left-to-right arc

```
for each [min, max] with max - min > 1 do
    for each l from min to max - 2 do
        double best = score[min][max][l]
        for each r from l + 1 to max - 1 do
            for each mid from l + 1 to r do
                t1 = score[min][mid][l]
                t2 = score[mid][max][r]
                double current = t1 + t2 + score(l → r)
                if current > best then
                    best = current
            score[min][max][l] = best
```

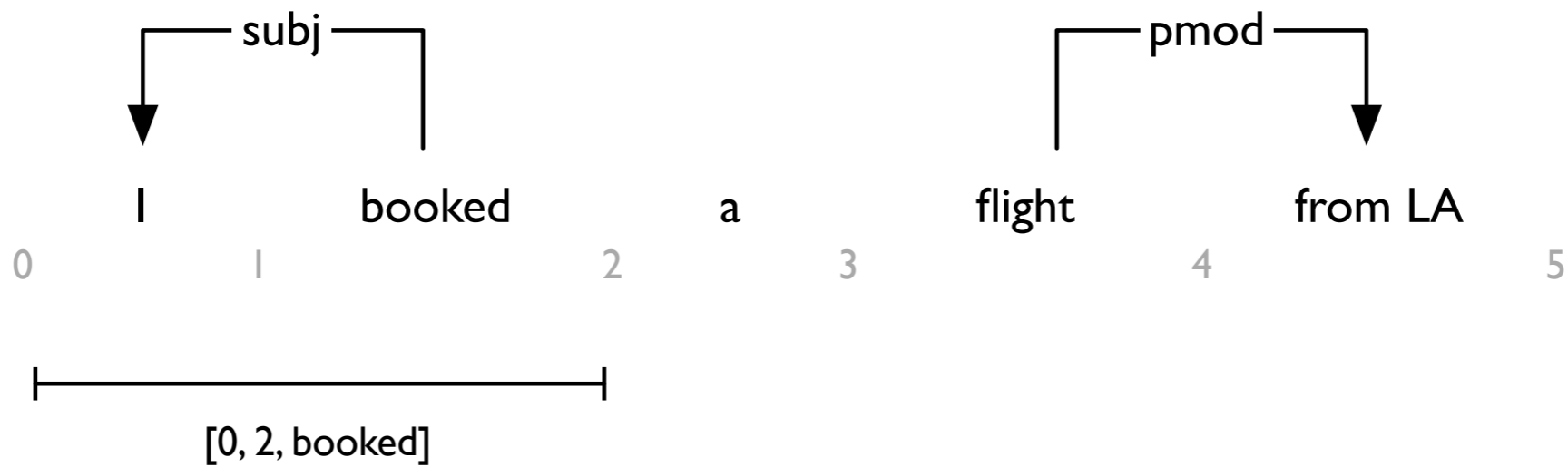


Adding a right-to-left arc





Adding a right-to-left arc

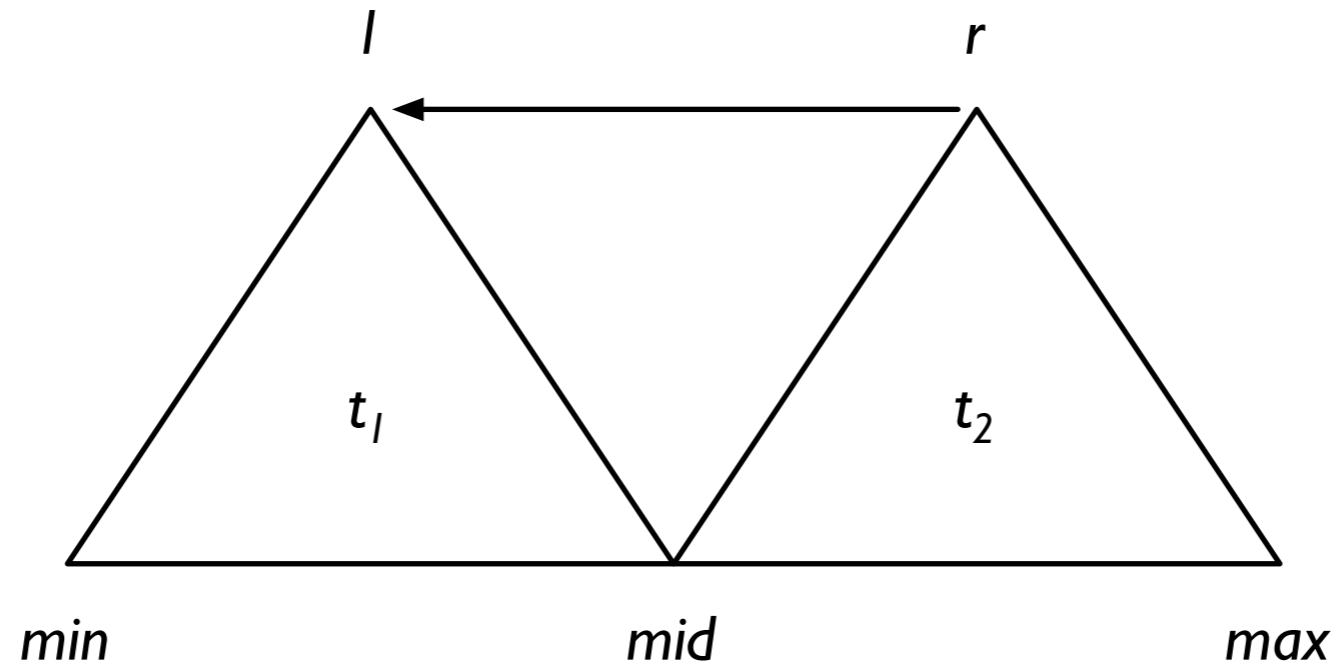




UPPSALA
UNIVERSITET

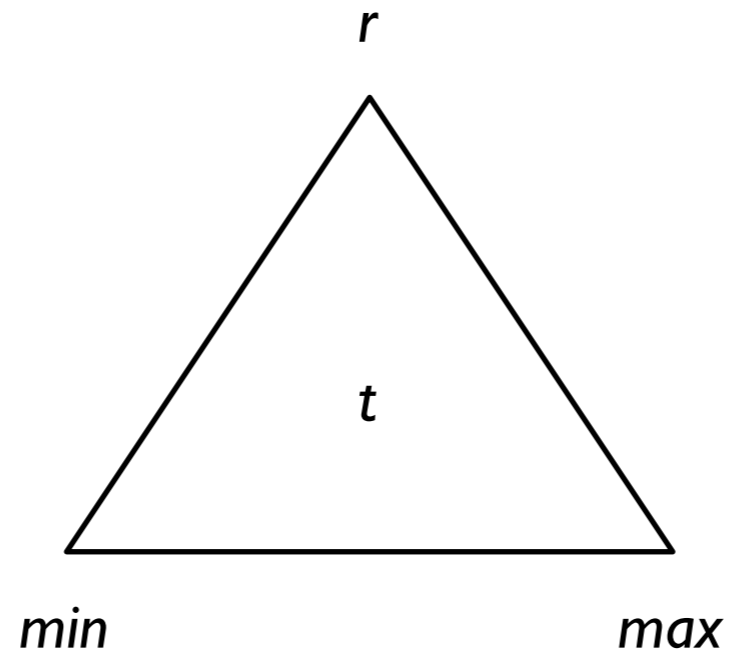
Collins' algorithm

Adding a right-to-left arc





Adding a right-to-left arc



$$\text{score}(t) = \text{score}(t_1) + \text{score}(t_2) + \text{score}(r \rightarrow l)$$

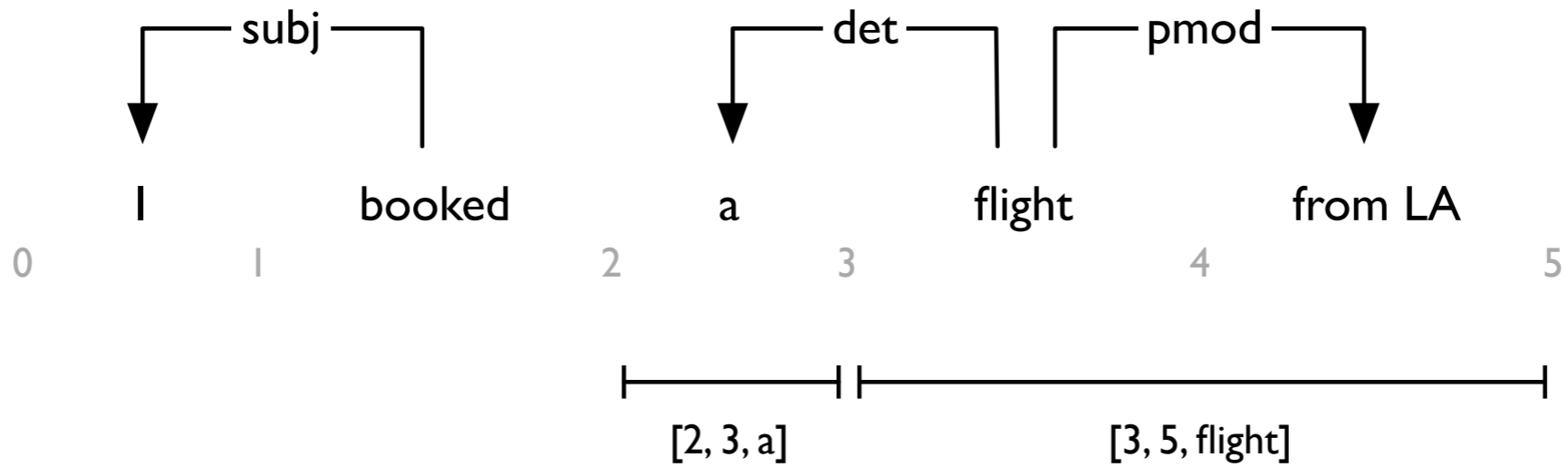


Adding a right-to-left arc

```
for each [min, max] with max - min > 1 do
    for each r from min + 1 to max - 1 do
        double best = score[min][max][r]
        for each l from min to r - 1 do
            for each mid from l + 1 to r do
                t1 = score[min][mid][l]
                t2 = score[mid][max][r]
                double current = t1 + t2 + score(r → l)
                if current > best then
                    best = current
            score[min][max][r] = best
```

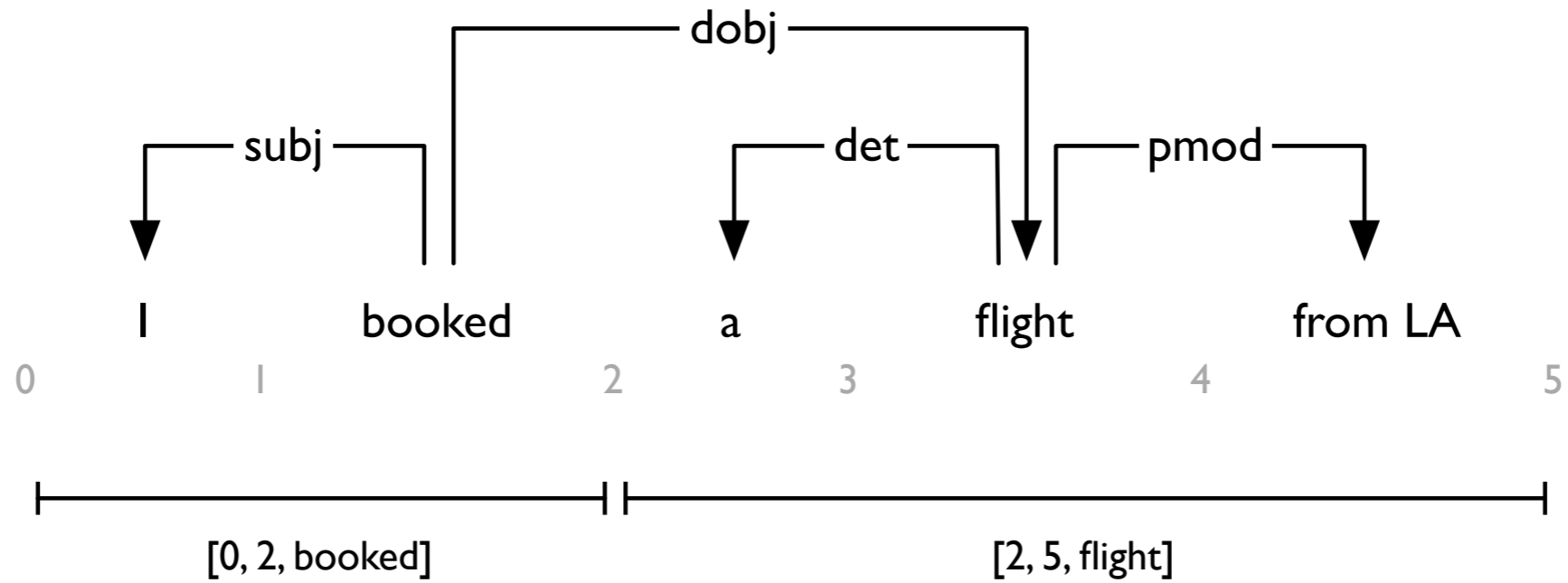


Finishing up



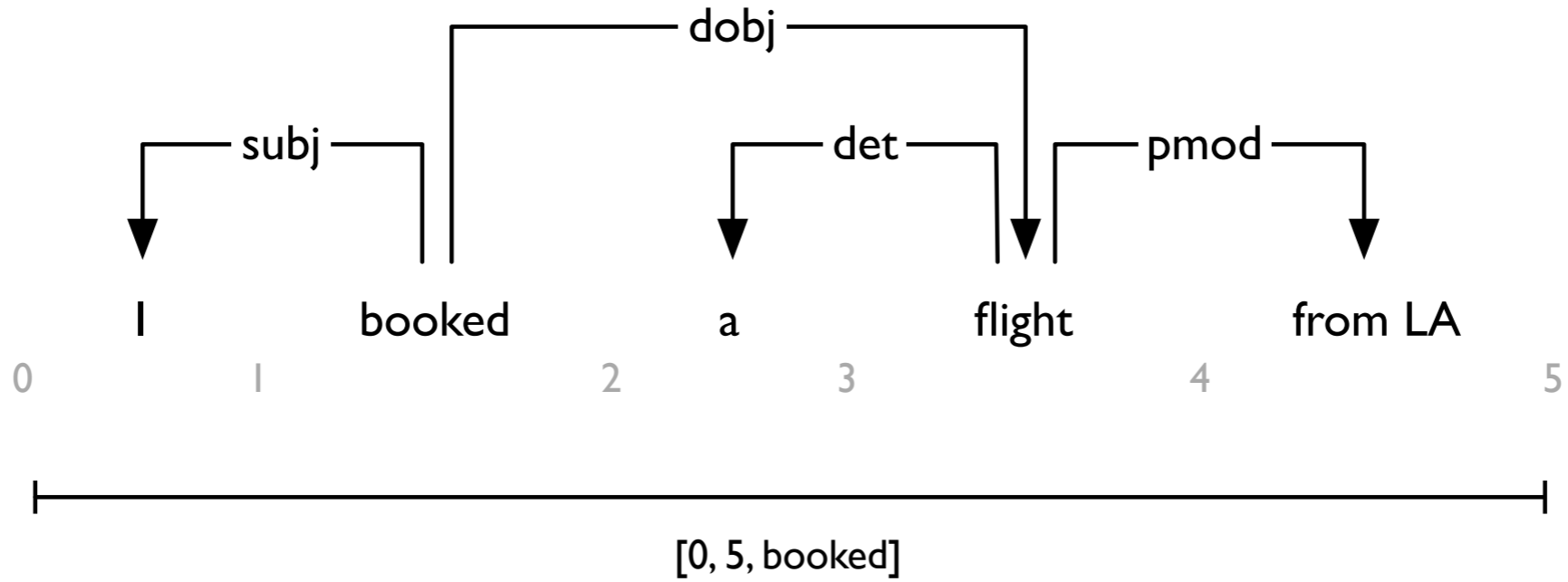


Finishing up





Finishing up





Complexity analysis

- **Space requirement:**
 $O(|w|^3)$
- **Runtime requirement:**
 $O(|w|^5)$



Extension to the labeled case

- It is important to distinguish dependencies of different types between the same two words.

Example: subj, dobj

- For this reason, practical systems typically deal with **labeled arcs**.
- The question then arises how to extend Collins' algorithm to the labeled case.



Naive approach

- Add an innermost loop that iterates over all edge labels in order to find the combination that maximizes the overall score.
- For each step of the original algorithm, we now need to make $|L|$ steps, where L is the set of all labels.



Smart approach

- Before parsing, compute a table that lists, for each head–dependent pair (h, d) , the label that maximizes the score of arcs $h \rightarrow d$.
- During parsing, simply look up the best label in the precomputed table.
- This adds (not multiplies!) a factor of $|L||w|^2$ to the overall runtime of the algorithm.



Summary

- Collins' algorithm is a CKY-style algorithm for computing the highest-scoring dependency tree under an arc-factored scoring model.
- It runs in time $O(|w|^5)$.
This may not be practical for long sentences.



UPPSALA
UNIVERSITET

Eisner's algorithm

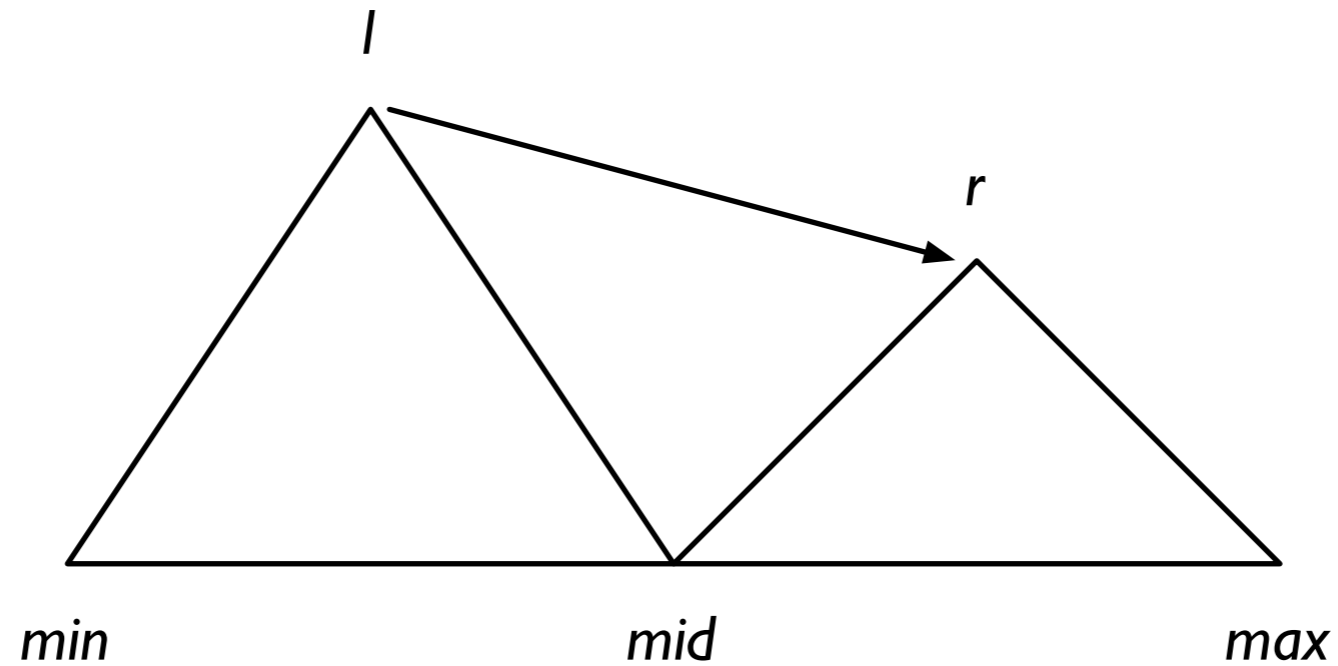


Eisner's algorithm

- With its runtime of $O(|w|^5)$, Collins' algorithm may not be of much use in practice.
- With Eisner's algorithm we will be able to solve the same problem in $O(|w|^3)$.



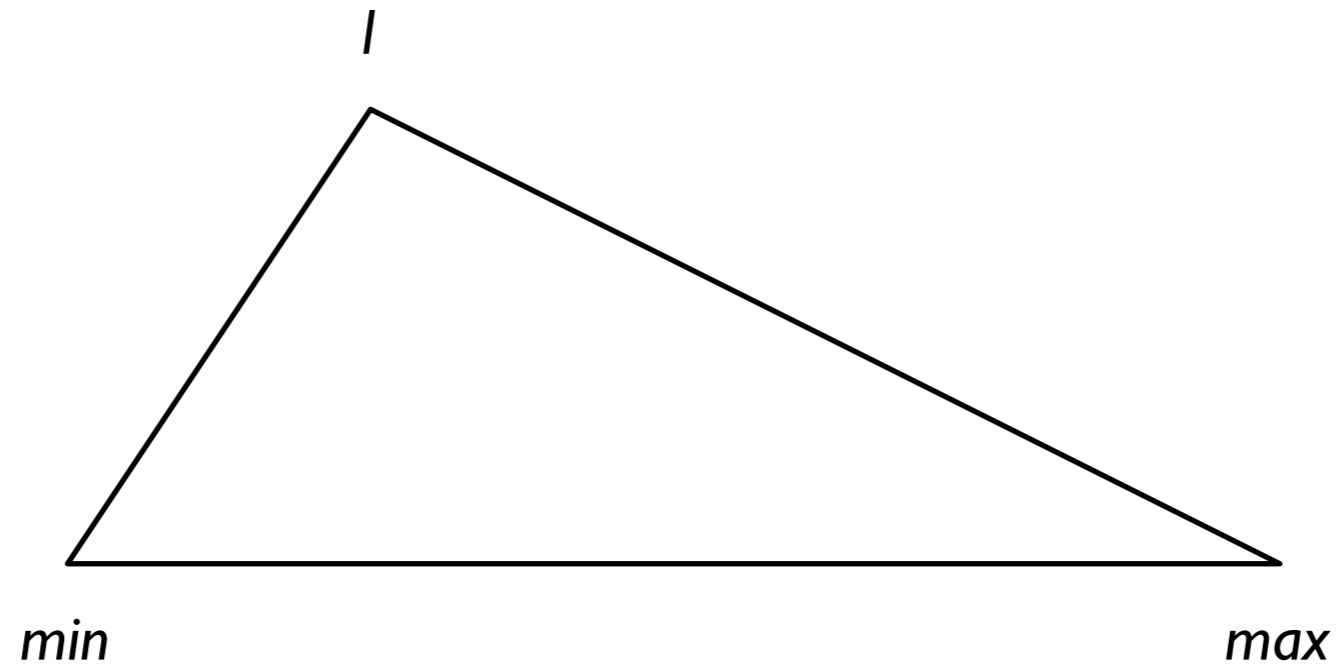
Basic idea



In Collins' algorithm, adding a left-to-right arc is done in one single step, specified by 5 positions.



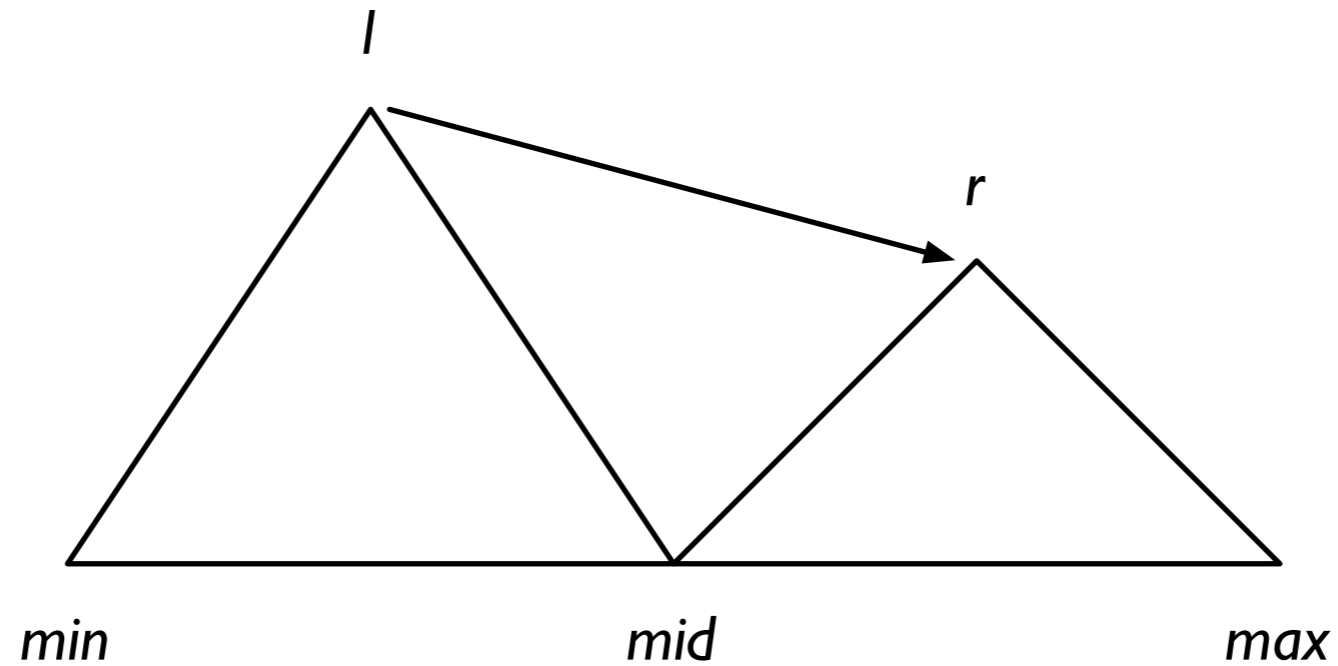
Basic idea



In Collins' algorithm, adding a left-to-right arc is done in one single step, specified by 5 positions.



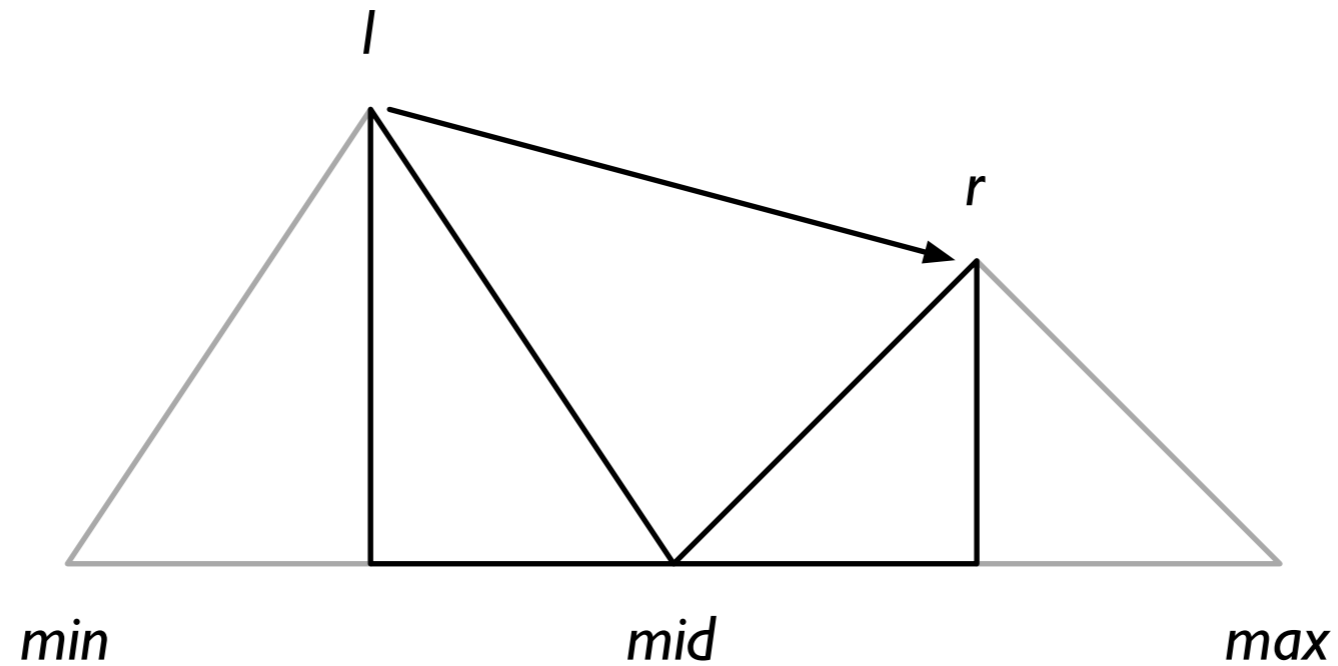
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



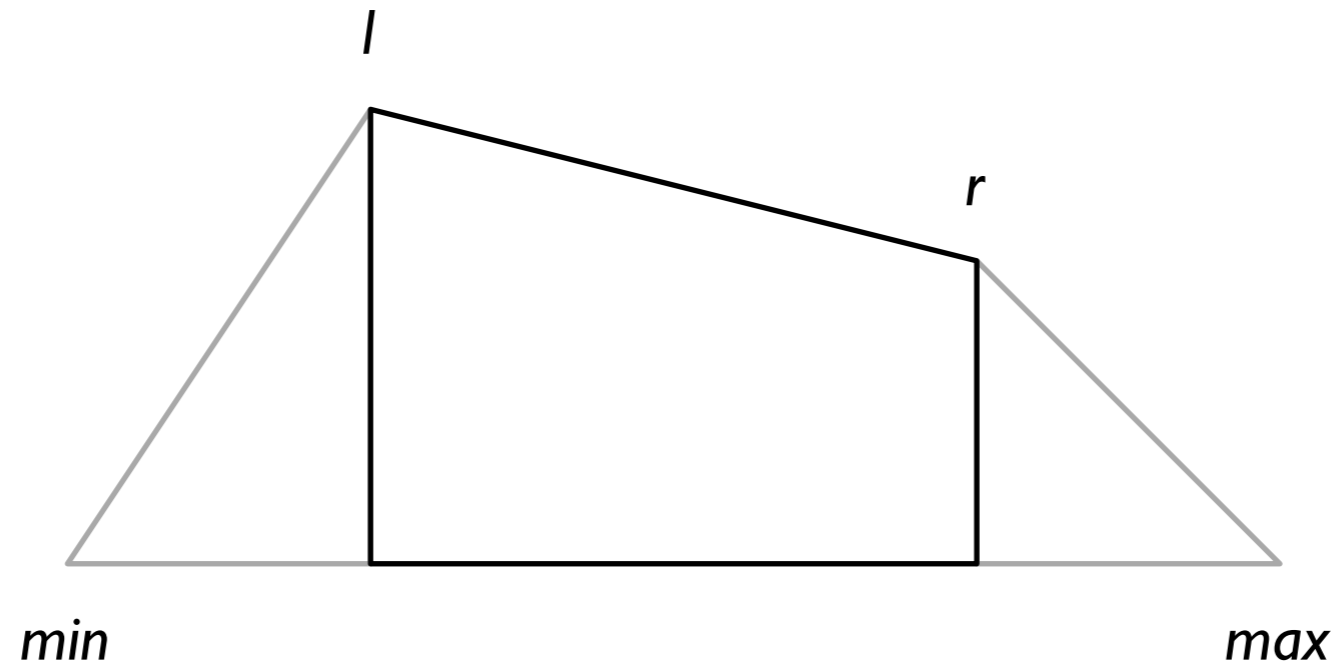
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



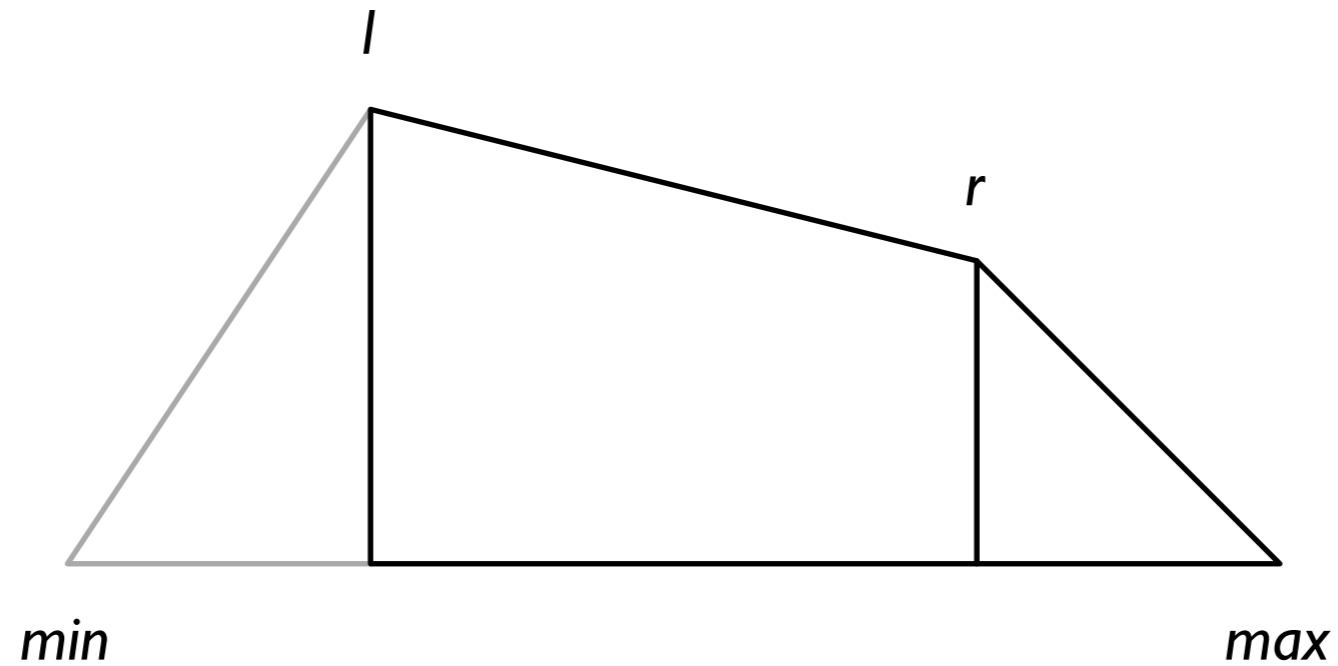
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



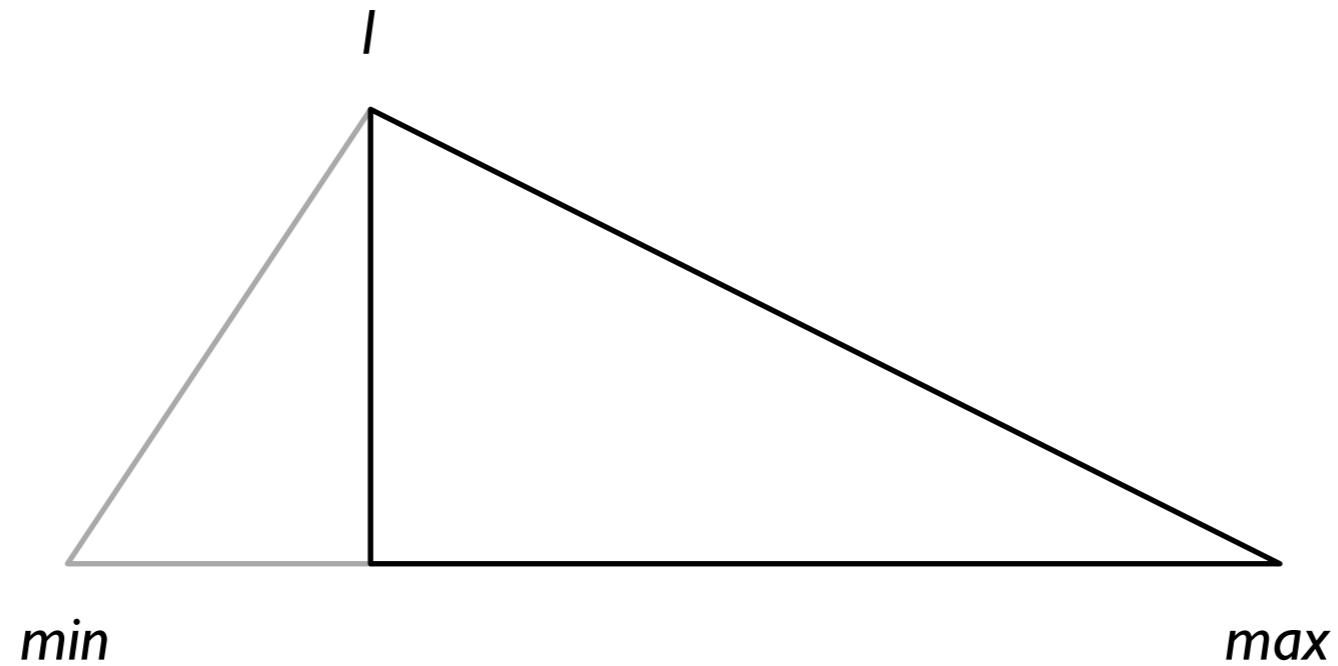
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



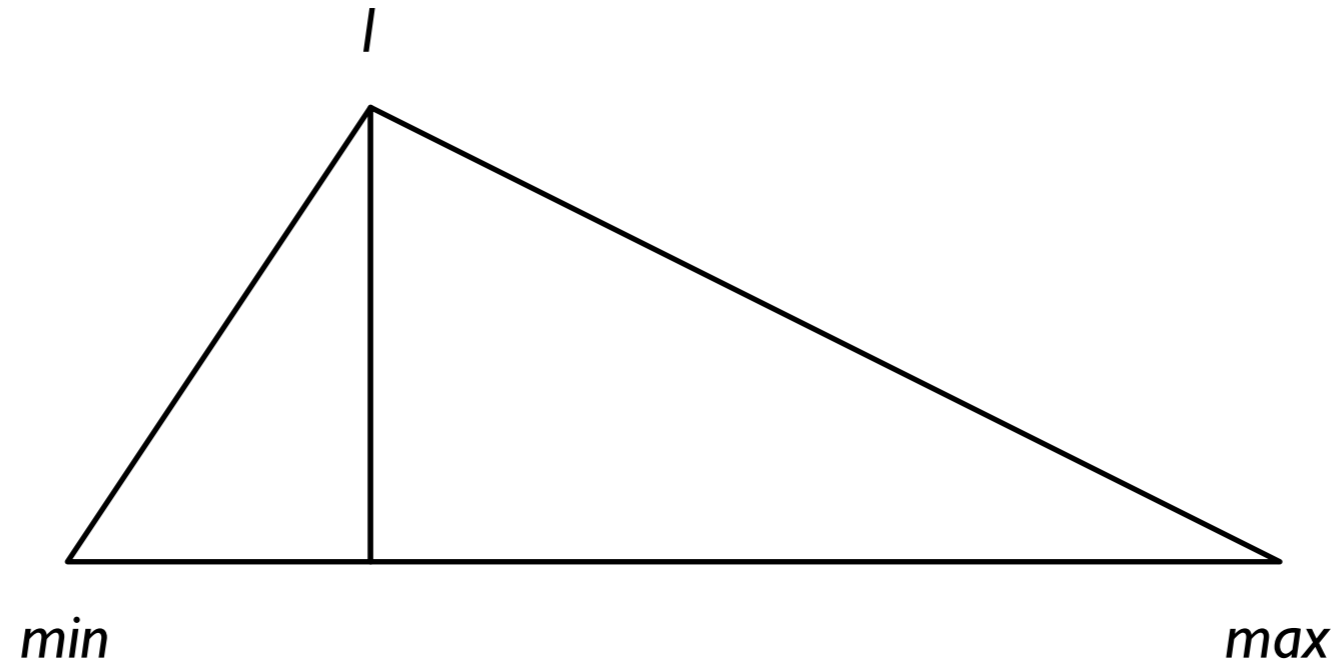
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



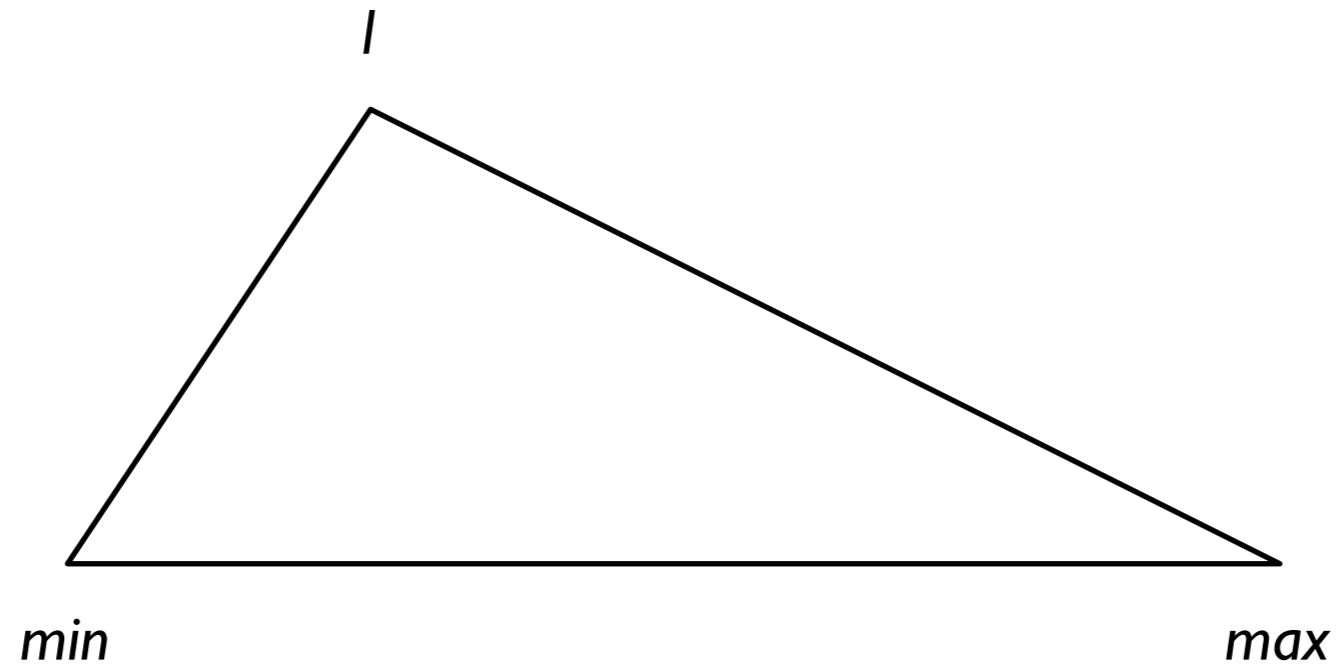
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



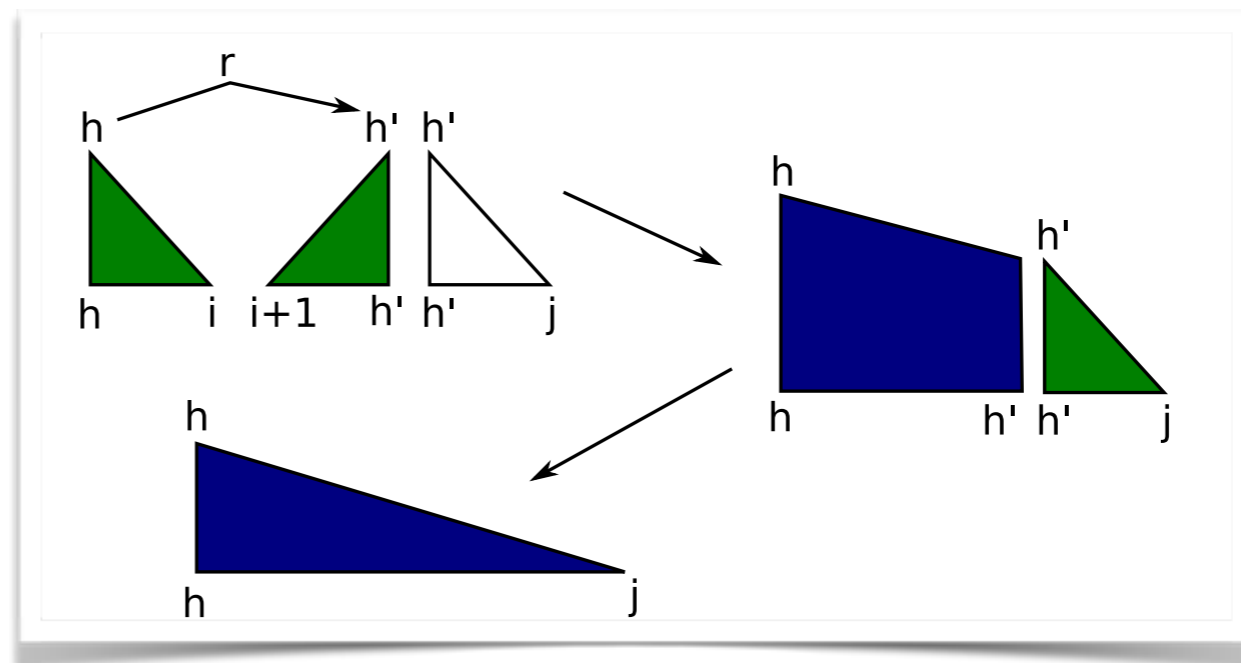
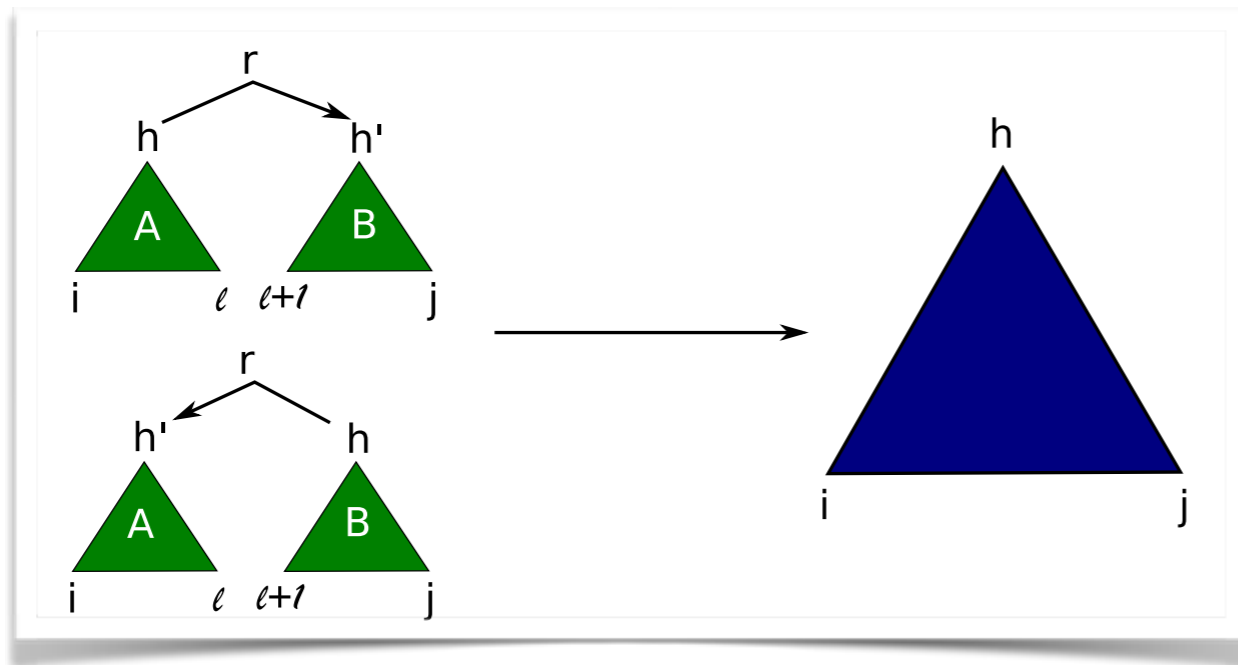
Basic idea



In Eisner's algorithm, the same thing is done in three steps, each one specified by 3 positions.



Comparison





Adding a left-to-right arc

```
for each [min, max] with max - min > 1 do
  double best- = score[min][max][left][-]
  double best+ = score[min][max][left][+]
  for each mid from min + 1 to max - 1 do
    t1 = score[min][mid][left][+]
    t2 = score[mid][max][right][+]
    double current- = t1 + t2 + score(min → max)
    if current- > best- then
      best- = current-
    t3 = score[min][mid+1][left][-]
    t4 = score[mid][max][left][+]
    double current+ = t3 + t4
    if current+ > best+ then
      best+ = current+
  score[min][max][left][-] = best-
  score[min][max][left][+] = best+
```



Adding a right-to-left arc

```
for each [min, max] with max - min > 1 do
  double best- = score[min][max][right][-]
  double best+ = score[min][max][right][+]
  for each mid from min + 1 to max - 1 do
    t1 = score[min][mid][left][+]
    t2 = score[mid][max][right][+]
    double current- = t1 + t2 + score(max → min)
    if current- > best- then
      best- = current-
    t3 = score[min][mid][right][+]
    t4 = score[mid-1][max][right][-]
    double current+ = t3 + t4
    if current+ > best+ then
      best+ = current+
  score[min][max][right][-] = best-
  score[min][max][right][+] = best+
```




Finishing up

- Find mid (from 1 to n) that maximizes:
$$\text{score}[0][\text{mid}][\text{right}][+] + \text{score}[\text{mid}-1][n][\text{left}][+]$$
- Left and right half-trees are built independently, combined only in the last step



Summary

- Eisner's algorithm is an improvement over Collin's algorithm that runs in time $O(|w|^3)$.
- The same scoring model can be used.
- The same technique for extending the parser to labeled parsing can be used.
- Eisner's algorithm is the basis of current arc-factored dependency parsers.