# Using N-Gram LMs with SCFG TMs
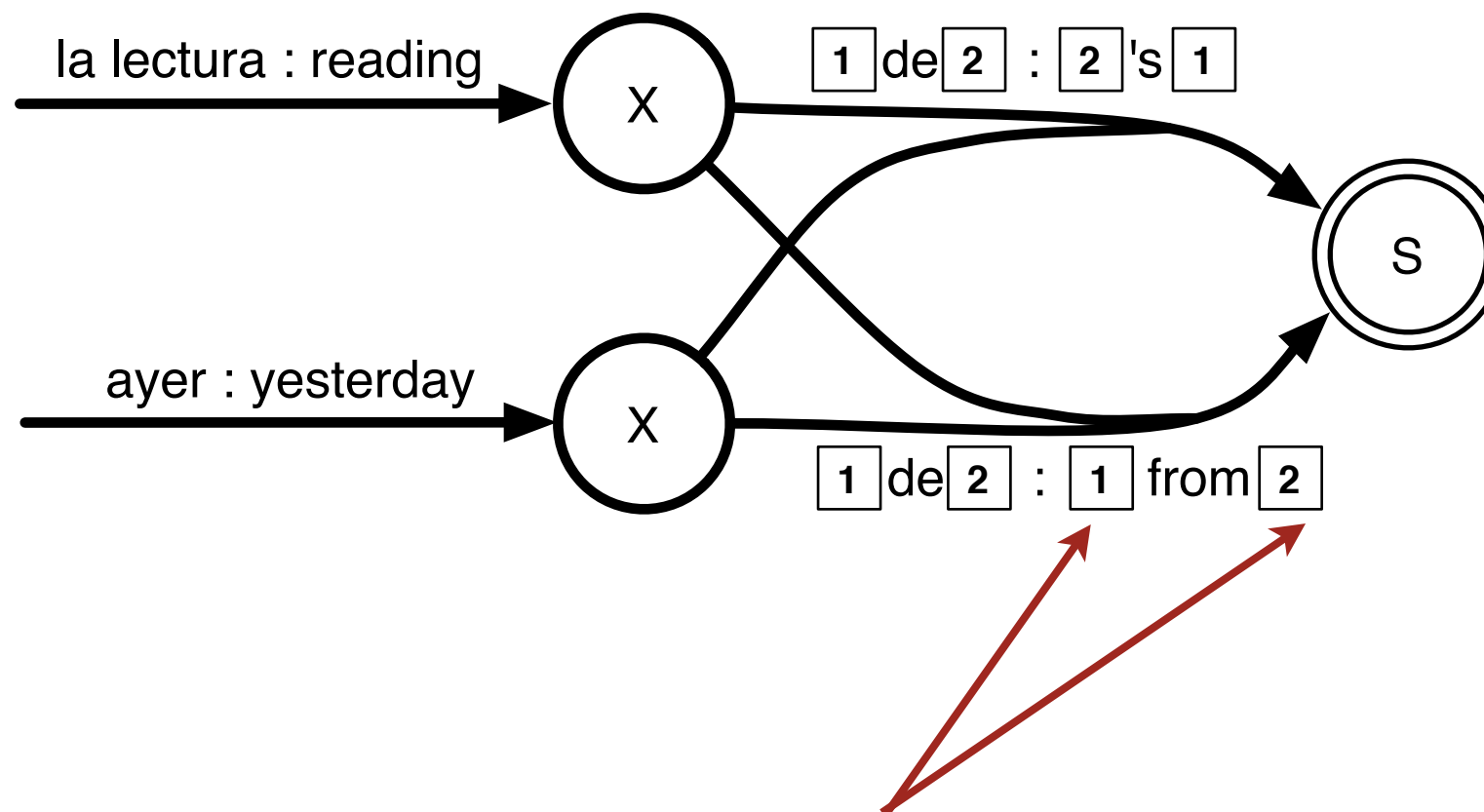
March 20, 2013

# Hypergraph review
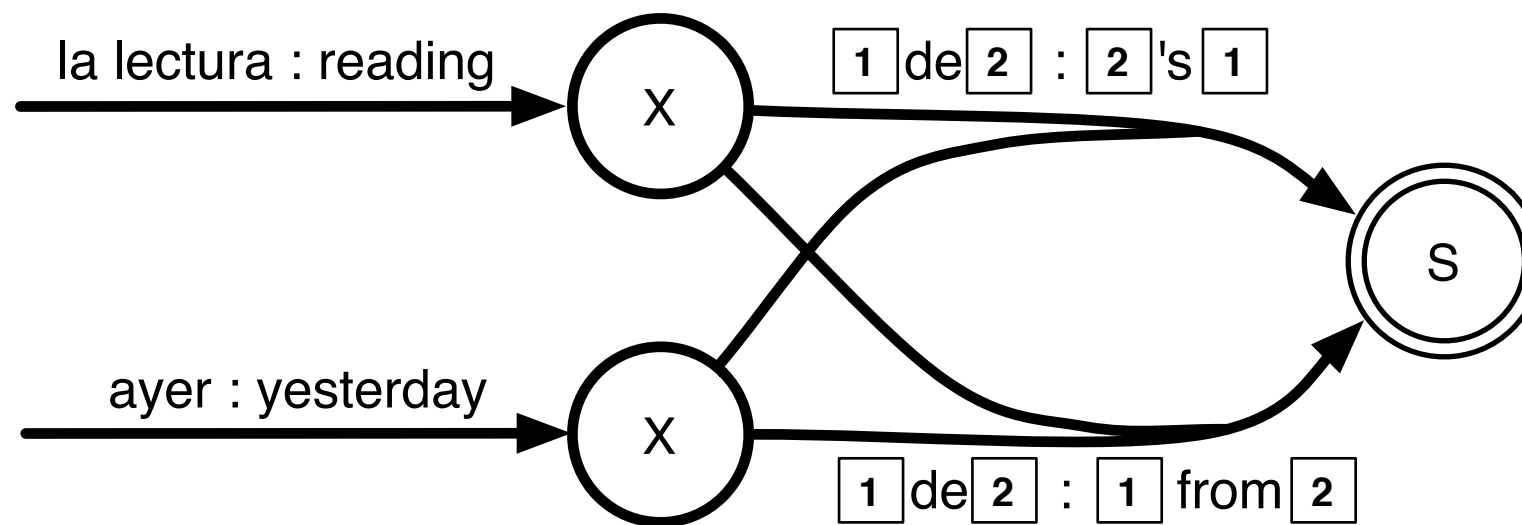
la lectura : reading ⟶ X

$\boxed{1}$ de $\boxed{2}$ : $\boxed{2}$ 's $\boxed{1}$

ayer : yesterday ⟶ X

$\boxed{1}$ de $\boxed{2}$ : $\boxed{1}$ from $\boxed{2}$

S

Source label

Target label

Goal node

# Hypergraph review



la lectura : reading

ayer : yesterday

X

X

S

**1** de **2** : **2** 's **1**

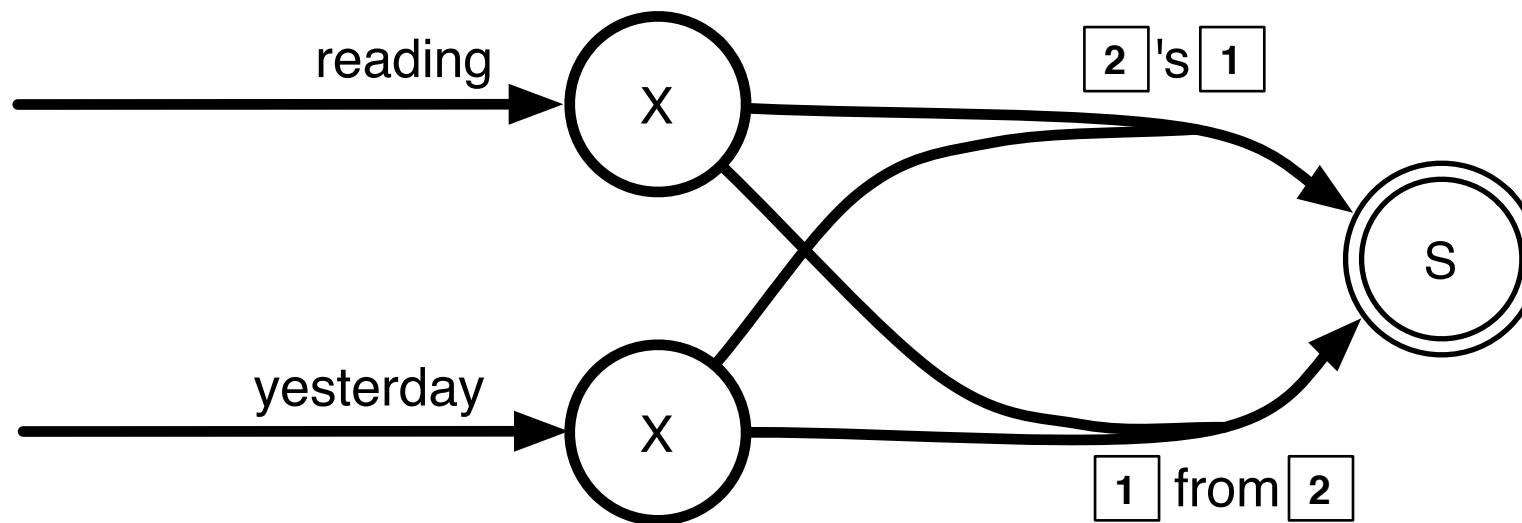**1** de **2** : **1** from **2**

Substitution sites / variables / non-terminals
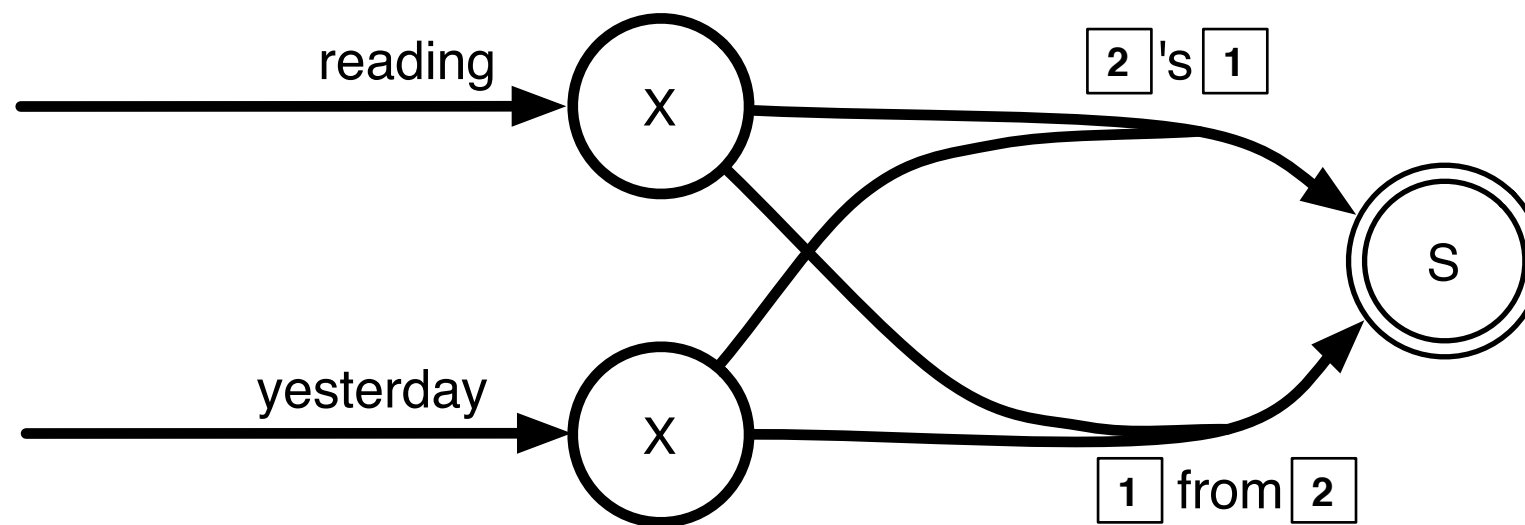
# Hypergraph review



For LM integration, we ignore the source!

# Hypergraph review



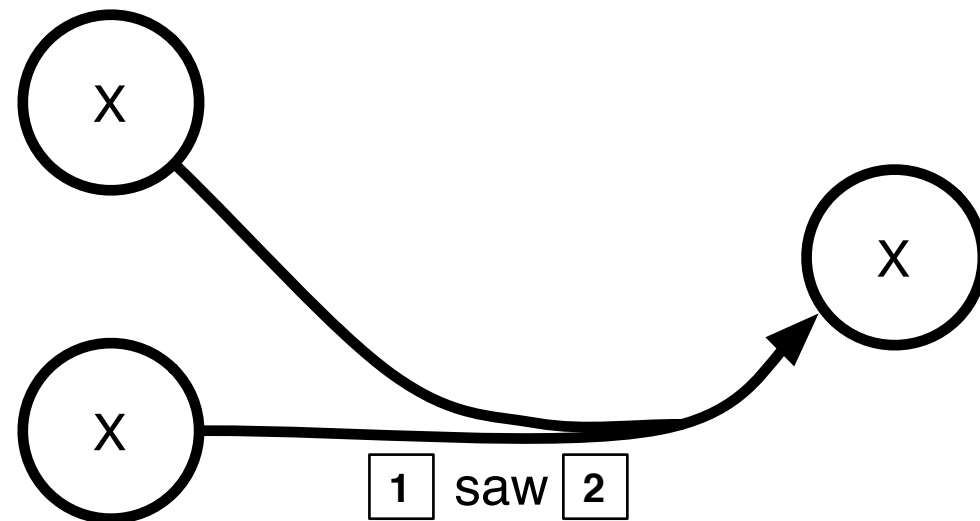**For LM integration, we ignore the source!**

# Hypergraph review



$$\{(\ yesterday\ 's\ reading\ ),$$
$$(\ reading\ from\ yesterday\ )\}$$

**How can we add the LM score to each string derived by the hypergraph?**

# LM Integration

- If LM features were purely local …

  - "Unigram" model

  - Discriminative LM

- … integration would be a breeze

  - Add an "LM feature" to every edge

- But, LM features are non-local!

# Why is it hard?



Two problems:

1. What is the content of the variables?

# Why is it hard?



Two problems:

   1. What is the content of the variables?
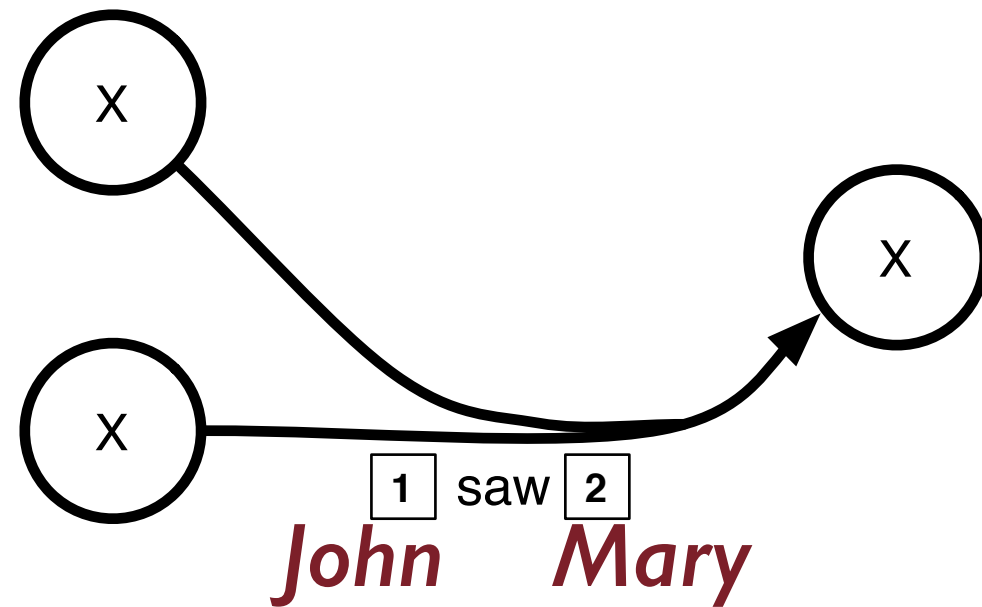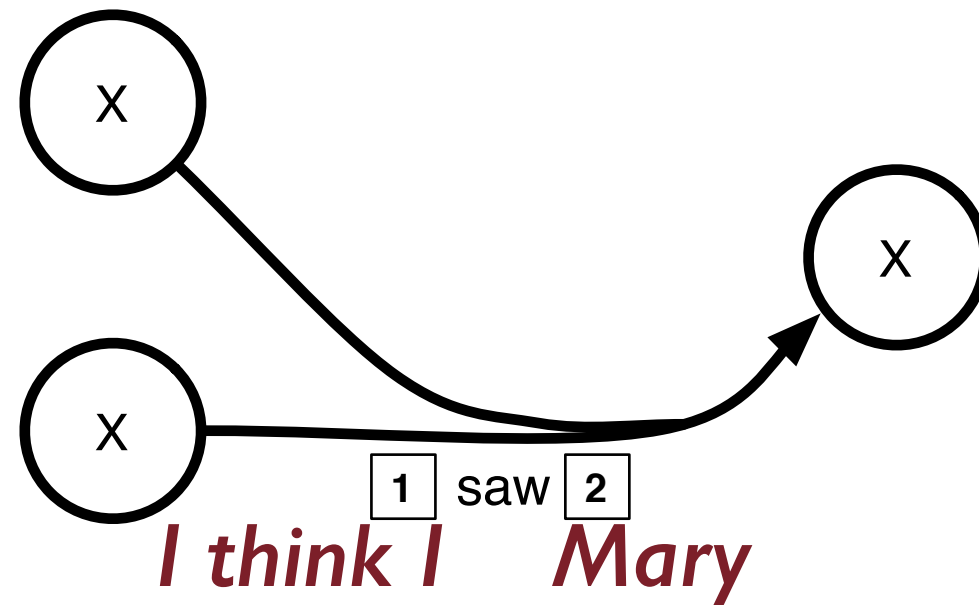
# Why is it hard?



Two problems:

1. What is the content of the variables?

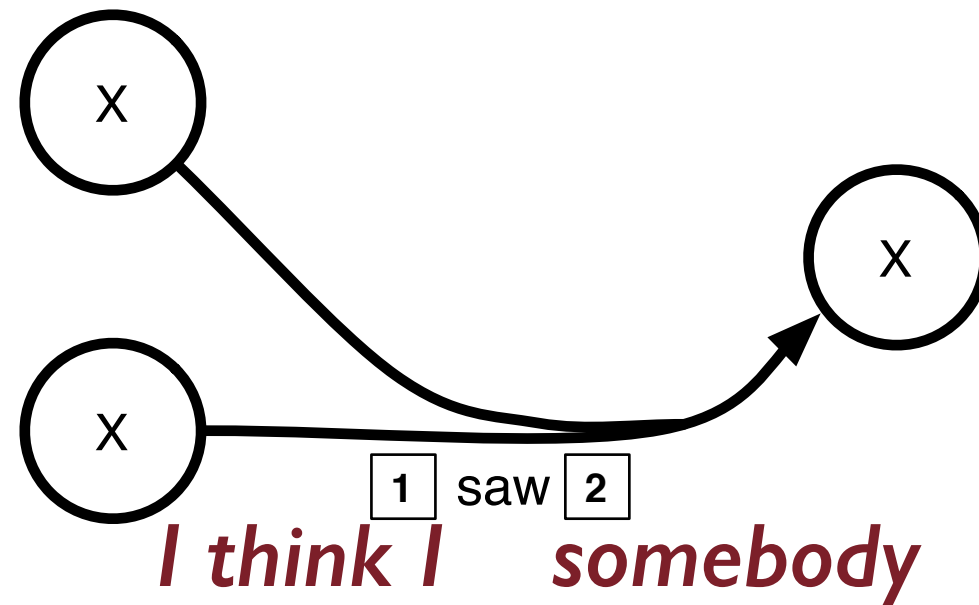# Why is it hard?



Two problems:

1. What is the content of the variables?

# Why is it hard?



Two problems:

1. What is the content of the variables?

2. What will be the **left context** when this string is substituted somewhere?

# Why is it hard?



Two problems:

1. What is the content of the variables?

2. What will be the **left context** when this string is substituted somewhere?

# Why is it hard?



Two problems:

1. What is the content of the variables?

2. What will be the **left context** when this string is substituted somewhere?

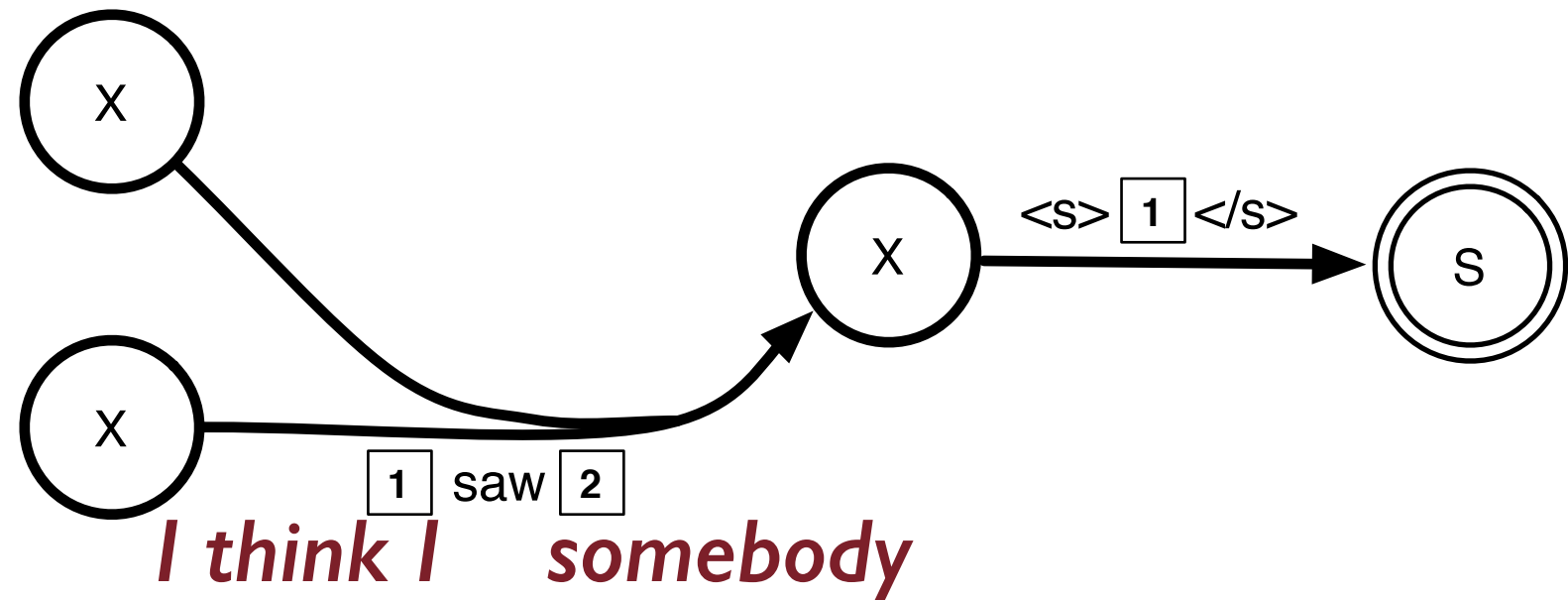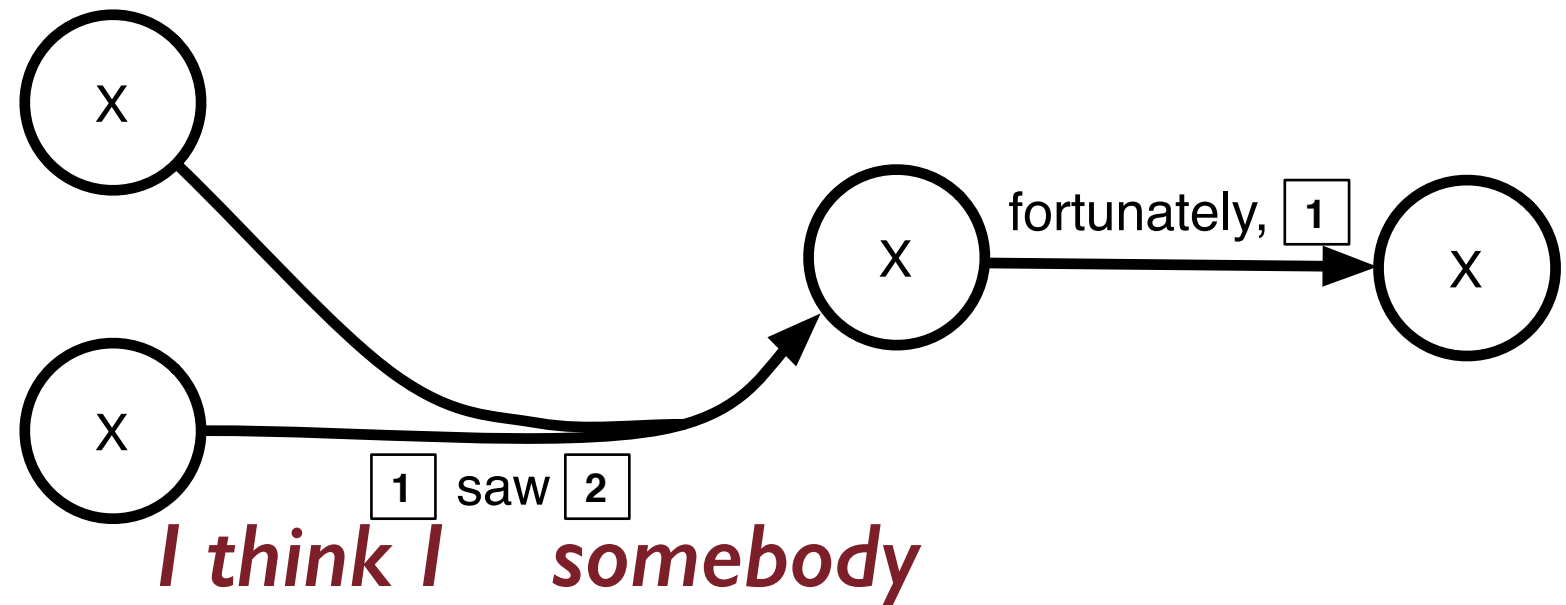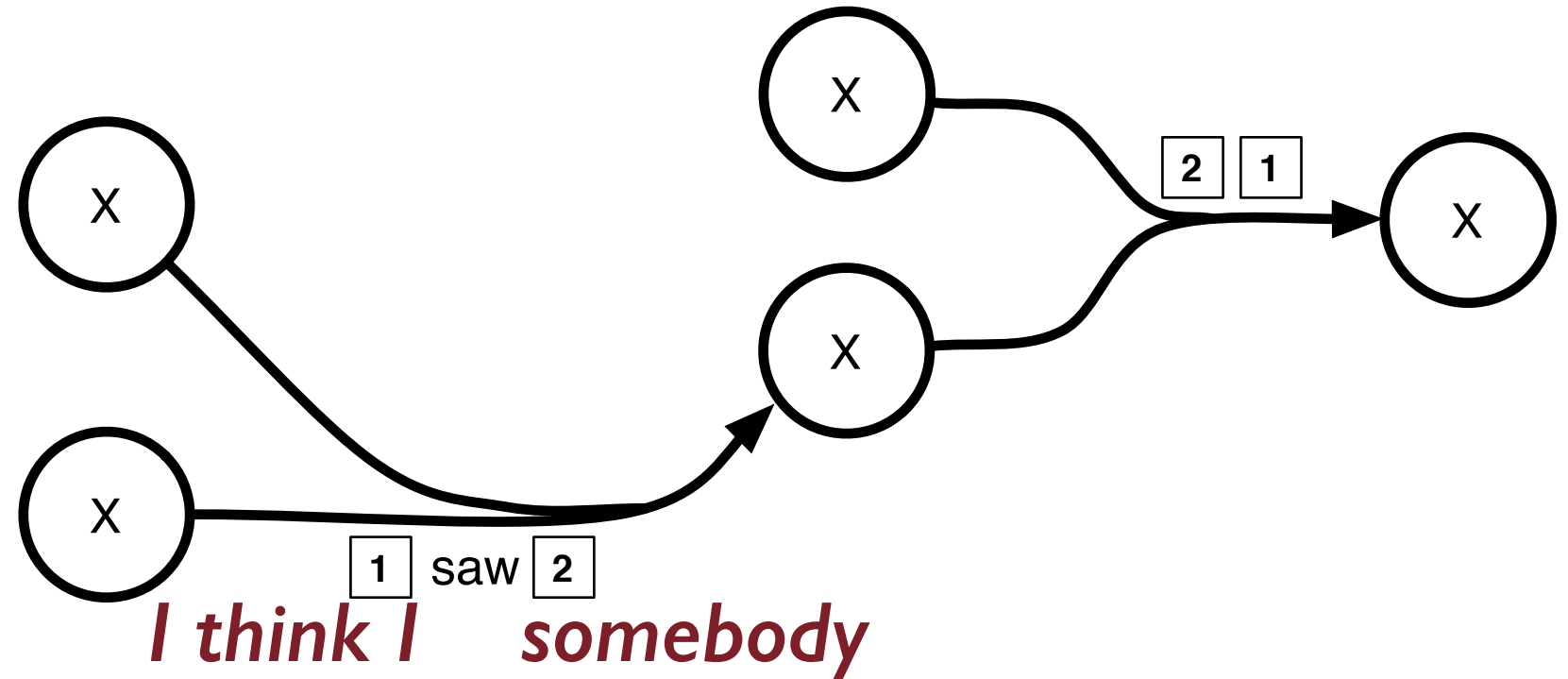# Why is it hard?



Two problems:

1. What is the content of the variables?

2. What will be the **left context** when this string is substituted somewhere?

# Naive solution

- Extract the all (k-best?) translations from the translation model

- Score them with an LM

- What's the problem with this?

# Outline of DP solution

- Use $n$-order Markov assumption to help us

  - In an $n$-gram LM, words more than $n$ words away will not affect the local (conditional) probability of a word in context

  - **This is not generally true, just the Markov assumption!**

- General approach

  - Restructure the hypergraph so that LM probabilities decompose along edges.

  - Solves both "problems"

    - we will not know the full value of variables, but we will know "enough".

    - defer scoring of left context until the context is established.

# Hypergraph restructuring

- Note the following three facts:

  - If you know $n$ or more consecutive words, the conditional probabilities of the $n$th, $(n+1)$th, ... words can be computed.

    - Therefore: add a feature weight to the edge for words.

  - $(n-1)$ words of context to the **left** is enough to determine the probability of any word

    - Therefore: split nodes based on the $(n-1)$ words on the **right** side of the span dominated by every node

  - $(n-1)$ words on the **left** side of a span cannot be scored with certainty because the context is not known

    - Therefore: split nodes based on the $(n-1)$ words on the **left** side of the span dominated by every node

# Hypergraph restructuring

- Note the following three facts:

  - If you know $n$ or more consecutive words, the conditional probabilities of the $n$th, $(n+1)$th, ... words can be computed.
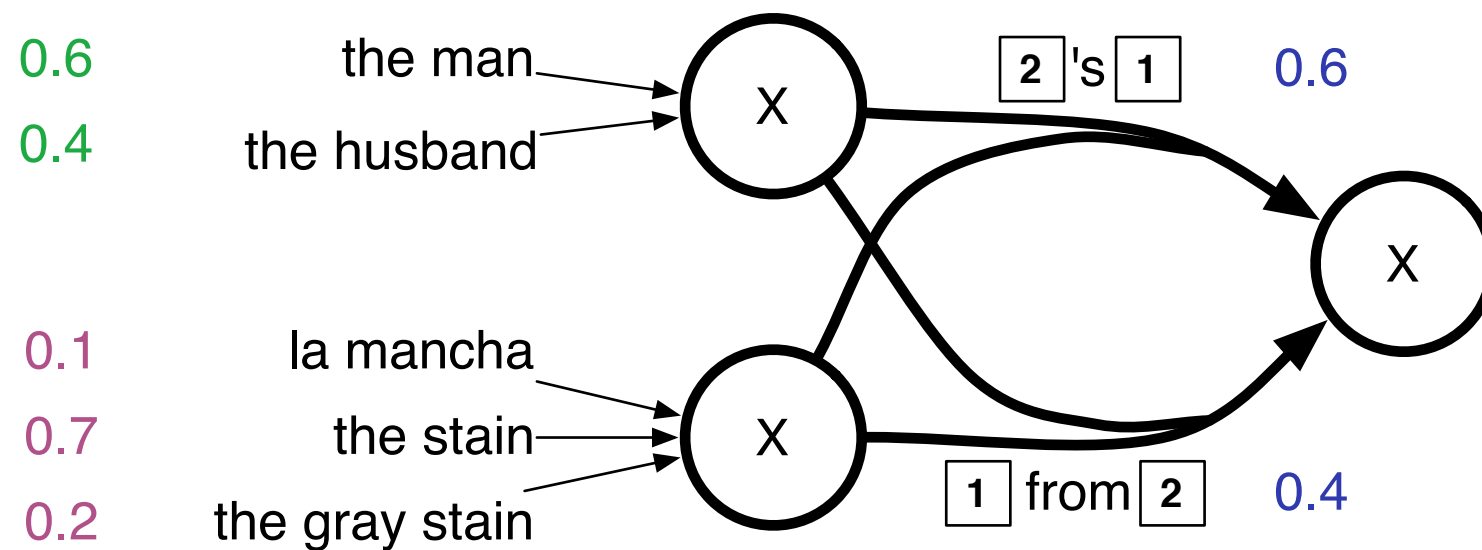
Split nodes by the (n-1) words on both sides of the convergent edges.

  - $(n-1)$ words on the **left** side of a span cannot be scored with certainty because the context is not known

    - Therefore: split nodes based on the $(n-1)$ words on the **left** side of the span dominated by every node

# Hypergraph restructuring

- Algorithm ("cube intersection"):

  - For each node $v$ (proceeding in **topological order** through the nodes)

    - For each edge $e$ with head-node $v$, compute the ($n$-1) words on the left and right; call this $q_e$

      - Do this by substituting the ($n$-1)x2 word string from the tail node corresponding to the substitution variable

      - If node $vq_e$ does not exist, create it, duplicating all outgoing edges from $v$ so that they also proceed from $vq_e$

      - Disconnect $e$ from $v$ and attach it to $vq_e$

  - Delete $v$

# Hypergraph restructuring

# Hypergraph restructuring



-LM Viterbi:

the stain's the man

# Hypergraph restructuring

## Let's add a bi-gram language model!

# Hypergraph restructuring

## Let's add a bi-gram language model!

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

# Hypergraph restructuring

0.6     the man

X
the man

Every node "remembers" enough
for edges to compute LM costs

X
the stain

# Complexity

- What is the run-time of this algorithm?

# Complexity

- What is the run-time of this algorithm?

$$O(|V||E||\Sigma|^{2(n-1)})$$

Going to longer n-grams is exponentially expensive!

# Cube pruning

- Expanding every node like this exhaustively is impractical

  - Polynomial time, but really, really big!

- Cube pruning: minor tweak on the above algorithm

  - Compute the k-best expansions at each node

  - Use an **estimate** (usually a unigram probability) of the unscored left-edge to rank the nodes

# Cube pruning

- Why "cube" pruning?

  - Cube-pruning only involves a "cube" when arity-2 rules are used!

  - More appropriately called "square" pruning with arity-1

  - Or "hypercube" pruning with arity > 2!

# Cube Pruning



monotonic grid?

|  | $\left(PP^{\ with\ \star\ Sharon}_{1,3}\right)$ 1.0 | $\left(PP^{\ along\ \star\ Sharon}_{1,3}\right)$ 3.0 | $\left(PP^{\ with\ \star\ Shalong}_{1,3}\right)$ 8.0 |
|---|---|---|---|
| $\left(VP^{\ held\ \star\ meeting}_{3,6}\right)$ 1.0 | 2.0 | 4.0 | 9.0 |
| $\left(VP^{\ held\ \star\ talk}_{3,6}\right)$ 1.1 | 2.1 | 4.1 | 9.1 |
| $\left(VP^{\ hold\ \star\ conference}_{3,6}\right)$ 3.5 | 4.5 | 6.5 | 11.5 |

# Cube Pruning



VP$_{1,6}$

PP$_{1,3}$  VP$_{3,6}$

non-monotonic grid
due to LM combo costs

|  | $\left(PP_{1,3}^{\,with\,\star\,Sharon}\right)$ | $\left(PP_{1,3}^{\,along\,\star\,Sharon}\right)$ | $\left(PP_{1,3}^{\,with\,\star\,Shalong}\right)$ |
|---|---|---|---|
|  | 1.0 | 3.0 | 8.0 |
| $\left(VP_{3,6}^{\,held\,\star\,meeting}\right)$ 1.0 | 2.0 + 0.5 | 4.0 + 5.0 | 9.0 + 0.5 |
| $\left(VP_{3,6}^{\,held\,\star\,talk}\right)$ 1.1 | 2.1 + 0.3 | 4.1 + 5.4 | 9.1 + 0.3 |
| $\left(VP_{3,6}^{\,hold\,\star\,conference}\right)$ 3.5 | 4.5 + 0.6 | 6.5 + 10.5 | 11.5 + 0.6 |

# Cube Pruning



VP$_{1,6}$

PP$_{1,3}$  VP$_{3,6}$

bigram (meeting, with)

**non-monotonic grid due to LM combo costs**

| | $\left(PP^{with}_{1,3} \star Sharon\right)$ | $\left(PP^{along}_{1,3} \star Sharon\right)$ | $\left(PP^{with}_{1,3} \star Shalong\right)$ |
|---|---|---|---|
| | **1.0** | **3.0** | **8.0** |
| $\left(VP^{held}_{3,6} \star meeting\right)$ **1.0** | 2.0 + 0.5 | 4.0 + 5.0 | 9.0 + 0.5 |
| $\left(VP^{held}_{3,6} \star talk\right)$ **1.1** | 2.1 + 0.3 | 4.1 + 5.4 | 9.1 + 0.3 |
| $\left(VP^{hold}_{3,6} \star conference\right)$ **3.5** | 4.5 + 0.6 | 6.5 + 10.5 | 11.5 + 0.6 |

# Cube Pruning



VP$_{1,6}$

PP$_{1,3}$   VP$_{3,6}$

non-monotonic grid
due to LM combo costs

| | $\left(\text{PP}^{\text{with} \star \text{Sharon}}_{1,3}\right)$ | $\left(\text{PP}^{\text{along} \star \text{Sharon}}_{1,3}\right)$ | $\left(\text{PP}^{\text{with} \star \text{Shalong}}_{1,3}\right)$ |
|---|---|---|---|
| | 1.0 | 3.0 | 8.0 |
| $\left(\text{VP}^{\text{held} \star \text{meeting}}_{3,6}\right)$  1.0 | 2.5 | 9.0 | 9.5 |
| $\left(\text{VP}^{\text{held} \star \text{talk}}_{3,6}\right)$  1.1 | 2.4 | 9.5 | 9.4 |
| $\left(\text{VP}^{\text{hold} \star \text{conference}}_{3,6}\right)$  3.5 | 5.1 | 17.0 | 12.1 |

# Cube Pruning

**k-best parsing**
(Huang and Chiang, 2005)

- a priority queue of candidates
- extract the best candidate

|  | $\left(\text{PP}_{1,3}^{\text{with} \star \text{Sharon}}\right)$ | $\left(\text{PP}_{1,3}^{\text{along} \star \text{Sharon}}\right)$ | $\left(\text{PP}_{1,3}^{\text{with} \star \text{Shalong}}\right)$ |
|---|---|---|---|
|  | 1.0 | 3.0 | 8.0 |
| $\left(\text{VP}_{3,6}^{\text{held} \star \text{meeting}}\right)$   1.0 | 2.5 | 9.0 | 9.5 |
| $\left(\text{VP}_{3,6}^{\text{held} \star \text{talk}}\right)$   1.1 | 2.4 | 9.5 | 9.4 |
| $\left(\text{VP}_{3,6}^{\text{hold} \star \text{conference}}\right)$   3.5 | 5.1 | 17.0 | 12.1 |

# Cube Pruning

*k*-best parsing
(Huang and Chiang, 2005)

- a priority queue of candidates
- extract the best candidate
- push the two successors

|  | $\left(PP_{1,3}^{\text{with } \star \text{ Sharon}}\right)$ 1.0 | $\left(PP_{1,3}^{\text{along } \star \text{ Sharon}}\right)$ 3.0 | $\left(PP_{1,3}^{\text{with } \star \text{ Shalong}}\right)$ 8.0 |
|---|---|---|---|
| $\left(VP_{3,6}^{\text{held } \star \text{ meeting}}\right)$ 1.0 | 2.5 | 9.0 | 9.5 |
| $\left(VP_{3,6}^{\text{held } \star \text{ talk}}\right)$ 1.1 | 2.4 | 9.5 | 9.4 |
| $\left(VP_{3,6}^{\text{hold } \star \text{ conference}}\right)$ 3.5 | 5.1 | 17.0 | 12.1 |

# Cube Pruning

**k-best parsing**
(Huang and Chiang, 2005)

- a priority queue of candidates
- extract the best candidate
- push the two successors

|  | $(\text{PP}^{\text{with} \star \text{Sharon}}_{1,3})$ | $(\text{PP}^{\text{along} \star \text{Sharon}}_{1,3})$ | $(\text{PP}^{\text{with} \star \text{Shalong}}_{1,3})$ |
|---|---|---|---|
|  | 1.0 | 3.0 | 8.0 |
| $(\text{VP}^{\text{held} \star \text{meeting}}_{3,6})$ 1.0 | 2.5 | 9.0 | 9.5 |
| $(\text{VP}^{\text{held} \star \text{talk}}_{3,6})$ 1.1 | 2.4 | 9.5 | 9.4 |
| $(\text{VP}^{\text{hold} \star \text{conference}}_{3,6})$ 3.5 | 5.1 | 17.0 | 12.1 |

# Cube pruning

- Widely used for phrase-based and syntax-based MT

- May be applied in conjunction with a bottom-up decoder, or as a second "rescoring" pass

  - Nodes may also be grouped together (for example, all nodes corresponding to a certain source span)

- Requirement for topological ordering means translation hypergraph may not have cycles

# LM Integration

| Method | Settings | Time | BLEU |
|---|---|---|---|
| rescore | $k = 10^4$ | 16 | 33.31 |
| rescore | $k = 10^5$ | 139 | 33.33 |
| intersect* | | 1455 | 37.09 |
| cube prune | $\varepsilon = 0$ | 23 | 36.14 |
| cube prune | $\varepsilon = 0.1$ | 35 | 36.77 |
| cube prune | $\varepsilon = 0.2$ | 111 | 36.91 |