

Decoding for SMT

Wilker Aziz

Universiteit van Amsterdam
w.aziz@uva.nl

April 11, 2016

Decoding for SMT

Decoding for SMT

We will be reasoning

Decoding for SMT

We will be reasoning

making decisions

Decoding for SMT

We will be reasoning

making decisions

about solutions

Decoding for SMT

We will be reasoning

making decisions

about solutions

structures (translations)

Decoding for SMT

We will be reasoning

making decisions

about solutions

structures (translations)

under a statistical model

Decoding for SMT

We will be reasoning

making decisions

about solutions

structures (translations)

under a statistical model

a function / a partial ordering

Keywords

We need to better characterise

Keywords

We need to better characterise

- ① space of solutions
get to know the terrain

Keywords

We need to better characterise

- ① space of solutions
get to know the terrain
- ② statistical model
how to assess the importance of things

Keywords

We need to better characterise

- ① space of solutions
get to know the terrain
- ② statistical model
how to assess the importance of things
- ③ decision rules
know what you want

Terrain



Lucas Leuzinger, December 2012

Terrain



not everywhere is a good place to stand

The importance (or cost) of things



German xmas market

The importance (or cost) of things



Brazilian street market

The importance (or cost) of things

German xmas market



Price of product depends on

The importance (or cost) of things

German xmas market



Price of product depends on

- the product

The importance (or cost) of things

German xmas market



Price of product depends on

- the product
- brand

The importance (or cost) of things

German xmas market



Price of product depends on

- the product
- brand
- quality of service

The importance (or cost) of things

Brazilian street market



Price of product depends on

- the product
- brand
- quality of service

The importance (or cost) of things

Brazilian street market



Price of product depends on

- the product
- brand
- quality of service
- seller's mood

The importance (or cost) of things

Brazilian street market



Price of product depends on

- the product
- brand
- quality of service
- seller's mood
- your look

The importance (or cost) of things

Brazilian street market



Price of product depends on

- the product
- brand
- quality of service
- seller's mood
- your look
- products you acquired from other stalls

The importance (or cost) of things

Brazilian street market



Price of product depends on

- the product
- brand
- quality of service
- seller's mood
- your look
- products you acquired from other stalls
- from whom you bought those

Know what you want



what do you think when you see this picture?

Know what you want



make a selfie (and lose a couple of fingers in the process)?

Know what you want



French goodies in a cozy chalet?

Content

- ① [Space of translations](#)
- ② [Formal devices](#)
- ③ [Linear models](#)
- ④ [Decision rules](#)
- ⑤ [Decoding](#)

Model of translational equivalences

Describes the process of generating translations of a given input

- constrains and characterises
the set of possible translation derivations

Model of translational equivalences

Describes the process of generating translations of a given input

- constrains and characterises
the set of possible translation derivations

Phrase-based MT

we observe an input, segment it into phrases, permute the phrases into target language word-order, and finally, translate segments independently

Model of translational equivalences

Describes the process of generating translations of a given input

- constrains and characterises the set of possible translation derivations

Phrase-based MT

we observe an input, segment it into phrases, permute the phrases into target language word-order, and finally, translate segments independently

Hierarchical MT

we parse the input with a CFG, then translate (using synchronous rules) each and every edge independently

CFGs and FSAs

Compactly represent the set of translations

- keep the representation cost a tractable (polynomial) function of the input length

Phrase-based MT $O(n^2 2^d)$

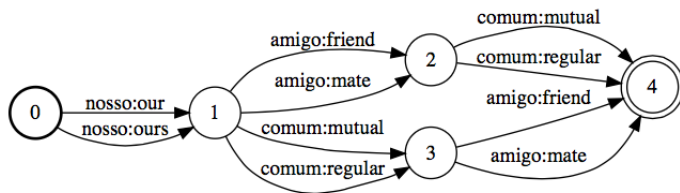
Hierarchical MT $O(n^3)$

Independence assumptions

Translation rules (flat or CFG) are applied independently

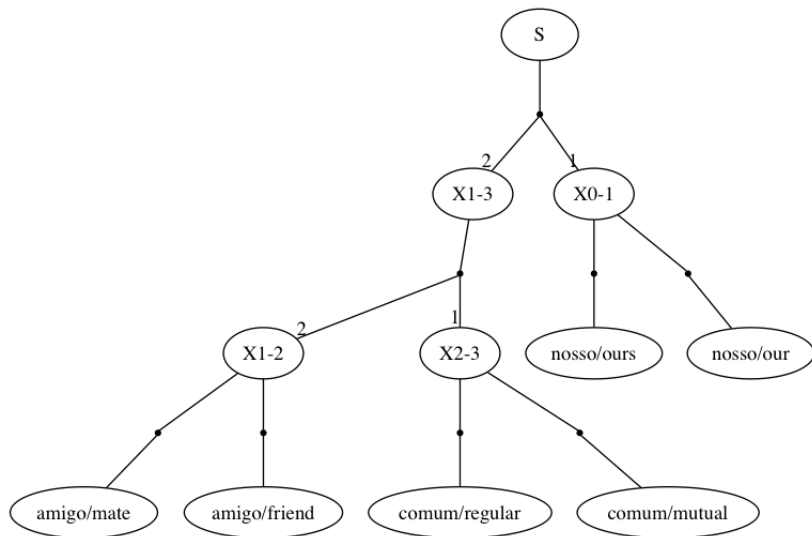
Independence assumptions

Translation rules (flat or CFG) are applied independently



Independence assumptions

Translation rules (flat or CFG) are applied independently



Directed B-hypergraphs

A hypergraph $\langle V, E \rangle$ consists of

- a set of nodes V
- a set of edges E
- an edge e has
 - a head node $\text{head}(e) \in V$
 - a tail $\text{tail}(e) \in V^*$ (sequence of nodes)

Directed B-hypergraphs

A hypergraph $\langle V, E \rangle$ consists of

- a set of nodes V
- a set of edges E
- an edge e has
 - a head node $\text{head}(e) \in V$
 - a tail $\text{tail}(e) \in V^*$ (sequence of nodes)

CFGs

- nonterminal \rightarrow node
- terminal \rightarrow terminal node
- rule \rightarrow edge
- LHS \rightarrow head
- RHS \rightarrow tail

Directed B-hypergraphs

A hypergraph $\langle V, E \rangle$ consists of

- a set of nodes V
- a set of edges E
- an edge e has
 - a head node $\text{head}(e) \in V$
 - a tail $\text{tail}(e) \in V^*$ (sequence of nodes)

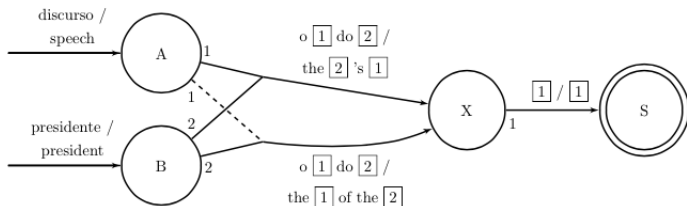
CFGs

- nonterminal \rightarrow node
- terminal \rightarrow terminal node
- rule \rightarrow edge
- LHS \rightarrow head
- RHS \rightarrow tail

FSAs

- state \rightarrow node
- symbol \rightarrow terminal node
- transition \rightarrow edge
- origin \rightarrow tail node
- destination \rightarrow head

A forest as a hypergraph



(a) Hypergraph

LHS	RHS _i	RHS _o
S →	X	[1]
X →	o A do B	the [2] 's [1]
X →	o A do B	the [1] of the [2]
A →	discurso	speech
B →	presidente	president

(b) Synchronous rules

Weighted sets

A weighted set $\langle \mathcal{D}, \omega \rangle$ consists of

- a set of structures (e.g. hyperpaths/derivations)
- a function $w : \mathcal{D} \rightarrow \mathcal{K}$

Let us focus on weighted sets whose weight functions factorise

$$w(\mathbf{d}) = \bigotimes_{e \in \mathbf{d}} w(e)$$

Often the structure is just a means to an end (the yield)

$$w(\mathbf{y}) = \bigoplus_{\mathbf{d} \in \mathcal{D}_{\mathbf{y}}} w(\mathbf{d})$$

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)
- \oplus and \otimes are binary operators

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)
- \oplus and \otimes are binary operators
- \oplus is commutative and has identity $\bar{0}$
 $a \oplus b = b \oplus a$ and $\bar{0} \oplus a = a \oplus \bar{0} = a$

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)
- \oplus and \otimes are binary operators
- \oplus is commutative and has identity $\bar{0}$
 $a \oplus b = b \oplus a$ and $\bar{0} \oplus a = a \oplus \bar{0} = a$
- \otimes is associative and has identity $\bar{1}$
 $(a \otimes b) \otimes c = a \otimes (b \otimes c)$ and $\bar{1} \otimes a = a \otimes \bar{1} = a$

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)
- \oplus and \otimes are binary operators
- \oplus is commutative and has identity $\bar{0}$
 $a \oplus b = b \oplus a$ and $\bar{0} \oplus a = a \oplus \bar{0} = a$
- \otimes is associative and has identity $\bar{1}$
 $(a \otimes b) \otimes c = a \otimes (b \otimes c)$ and $\bar{1} \otimes a = a \otimes \bar{1} = a$
- \otimes left distributes over \oplus
 $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$

Semirings

An algebraic structure $\mathcal{K} = \langle \mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1} \rangle$

- \mathbb{K} is a set (e.g. $\mathbb{N}, \mathbb{R}, \{0, 1\}$)
- \oplus and \otimes are binary operators
- \oplus is commutative and has identity $\bar{0}$
 $a \oplus b = b \oplus a$ and $\bar{0} \oplus a = a \oplus \bar{0} = a$
- \otimes is associative and has identity $\bar{1}$
 $(a \otimes b) \otimes c = a \otimes (b \otimes c)$ and $\bar{1} \otimes a = a \otimes \bar{1} = a$
- \otimes left distributes over \oplus
 $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$
- $\bar{0}$ is the \otimes -annihilator
 $\bar{0} \otimes a = a \otimes \bar{0} = \bar{0}$

Examples of semirings

Name	\mathbb{K}	\oplus	\otimes	$\bar{0}$	$\bar{1}$
BINARY	$\{0, 1\}$	\vee	\wedge	0	1
COUNTING	\mathbb{N}	$+$	\times	0	1
PROB	$[0, 1] \subset \mathbb{R}$	$+$	\times	0	1
LOGPROB	$\mathbb{R} \cup \{-\infty\}$	\oplus_{\log}	$+$	$-\infty$	0
VITERBI	$\mathbb{R} \cup \{-\infty\}$	\max	$+$	$-\infty$	0

where $a \oplus_{\log} b = \log(\exp(a) + \exp(b))$

Linear models

$$f(\mathbf{d}) = \mathbf{w}^\top \Phi(\mathbf{d})$$

where

- $\mathbf{w} \in \mathbb{R}^m$
- $\Phi(\mathbf{d}) = \langle \Phi_1(\mathbf{d}), \dots, \Phi_m(\mathbf{d}) \rangle$
- $\Phi_i(\mathbf{d}) \in \mathbb{R}$ is a feature function
- w_i is the relative contribution of the i th feature

Linear models and independence assumptions

$$f(\mathbf{d}) = \mathbf{w}^\top \Phi(\mathbf{d}) \quad (1)$$

$$= \sum_{i=1}^m w_i \Phi_i(\mathbf{d}) \quad (2)$$

$$= \sum_{i=1}^m w_i \prod_{e \in \mathbf{d}} \phi_i(e) \quad (3)$$

$$= \prod_{e \in \mathbf{d}} \sum_{i=1}^m w_i \phi_i(e) \quad (4)$$

$$= \prod_{e \in \mathbf{d}} \mathbf{w}^\top \phi(e) \quad (5)$$

Assumption

- $\Phi_i(\mathbf{d})$ factorises over edges
 $\phi_i(e)$ is a local feature function

Linear models and CFGs

Linear models can be expressed through hypergraphs using an appropriate semiring

Decision rules

Best translation (MAP)

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \sum_{d \in \mathcal{D}_{\mathbf{y}}} f(\mathbf{d})$$

Decision rules

Best translation (MAP)

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \sum_{d \in \mathcal{D}_{\mathbf{y}}} f(\mathbf{d})$$

Best derivation (Viterbi)

$$\mathbf{y}^* \approx \text{yield} \left\{ \arg \max_{\mathbf{d}} f(\mathbf{d}) \right\}$$

- less disambiguation power
- VITERBI semiring

Other decision rules?

Minimum Bayes risk (MBR)

$$\mathbf{y}^* = \arg \min_{\mathbf{y}'} \langle L(\mathbf{y}', \mathbf{y}) \rangle_{p(\mathbf{y})}$$

- requires the underlying model to have a probabilistic interpretation
- can be estimated through sampling

Other decision rules?

Minimum Bayes risk (MBR)

$$\mathbf{y}^* = \arg \min_{\mathbf{y}'} \langle L(\mathbf{y}', \mathbf{y}) \rangle_{p(\mathbf{y})}$$

- requires the underlying model to have a probabilistic interpretation
- can be estimated through sampling

Log-linear models

$$p(\mathbf{d}) = \frac{\exp(f(\mathbf{d}))}{\sum_{\mathbf{d}'} \exp(f(\mathbf{d}'))} \propto \exp \left(\sum_{e \in \mathbf{d}} \mathbf{w}^\top \phi(e) \right) = \prod_{e \in \mathbf{d}} \exp(\mathbf{w}^\top \phi(e))$$

LOGPROB semiring

Decoding

In SMT, decoding typically means the Viterbi approximation

$$\mathbf{d}^* = \arg \max_{\mathbf{d}} f(\mathbf{d})$$

Decoding

In SMT, decoding typically means the Viterbi approximation

$$\mathbf{d}^* = \arg \max_{\mathbf{d}} f(\mathbf{d})$$

If the statistical model $f(\mathbf{d})$ does not violate the independence assumptions of the model of translational equivalences

- steps in a derivation are weighted independently

Decoding

In SMT, decoding typically means the Viterbi approximation

$$\mathbf{d}^* = \arg \max_{\mathbf{d}} f(\mathbf{d})$$

If the statistical model $f(\mathbf{d})$ does not violate the independence assumptions of the model of translational equivalences

- steps in a derivation are weighted independently

there is a straightforward (tractable) decomposition of $f(\mathbf{d})$

- wFSA (phrase-based MT)
- wCFG (hierarchical MT)

Inside

The INSIDE recursion can be generalised to an arbitrary semiring

$$\beta(v) = \begin{cases} \bar{1} & \text{if } BS(v) = \emptyset \\ \bigoplus_{e \in BS(v)} w(e) \bigotimes_{u \in \text{tail}(e)} \beta(u) & \text{otherwise} \end{cases}$$

Inside

The INSIDE recursion can be generalised to an arbitrary semiring

$$\beta(v) = \begin{cases} \bar{1} & \text{if } BS(v) = \emptyset \\ \bigoplus_{e \in BS(v)} w(e) \bigotimes_{u \in \text{tail}(e)} \beta(u) & \text{otherwise} \end{cases}$$

- efficient bottom-up dynamic program $O(|G|)$
 $|G|$ is the size of the graphical representation of $f(\mathbf{d})$
 - a lattice (phrase-based MT)
 - a forest (hierarchical MT)

Inference

Viterbi derivation

- 1 start from the goal (root)
- 2 recursively rewrite every symbol v by solving

$$e = \arg \max_{e \in BS(v)} w(e) \bigotimes_{u \in \text{tail}(e)} \beta(u)$$

Inference

Viterbi derivation

- ❶ start from the goal (root)
- ❷ recursively rewrite every symbol v by solving

$$e = \arg \max_{e \in BS(v)} w(e) \bigotimes_{u \in \text{tail}(e)} \beta(u)$$

Sampling

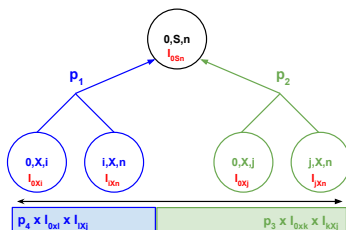
- ❶ start from the goal (root)
- ❷ recursively rewrite every symbol v by solving

$$e \sim p(e \in BS(v)|v) = \frac{w(e) \bigotimes_{u \in \text{tail}(e)} \beta(u)}{\beta(v)}$$

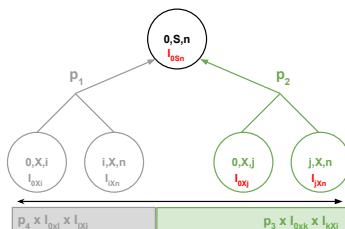
Viterbi



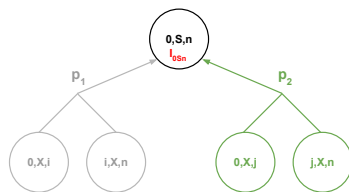
Viterbi



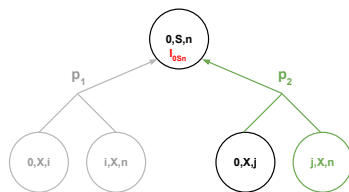
Viterbi



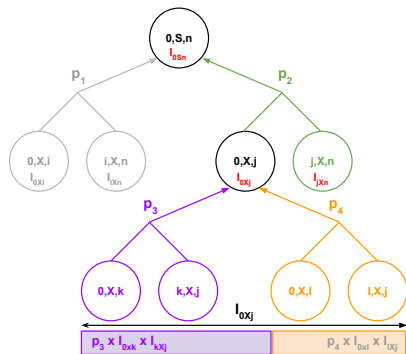
Viterbi



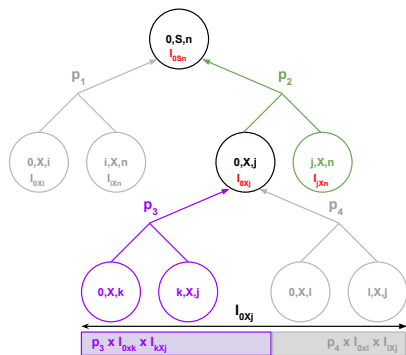
Viterbi



Viterbi



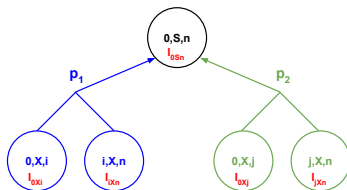
Viterbi



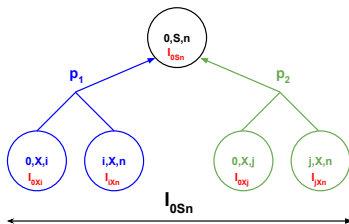
Sampling



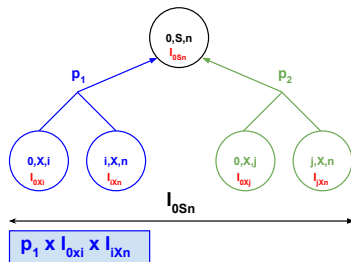
Sampling



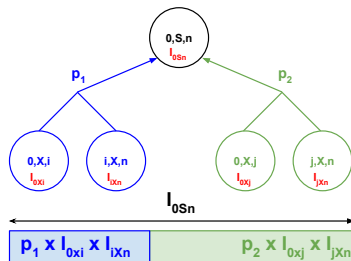
Sampling



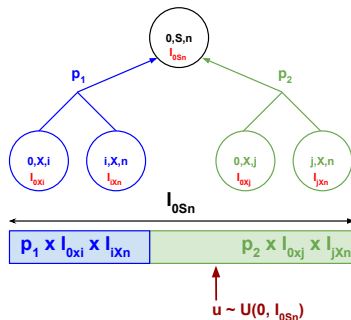
Sampling



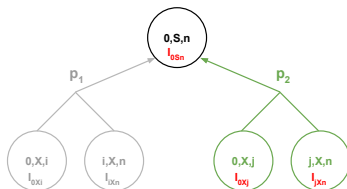
Sampling



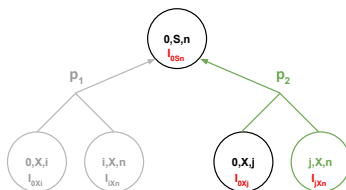
Sampling



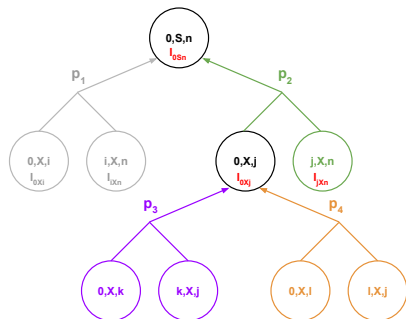
Sampling



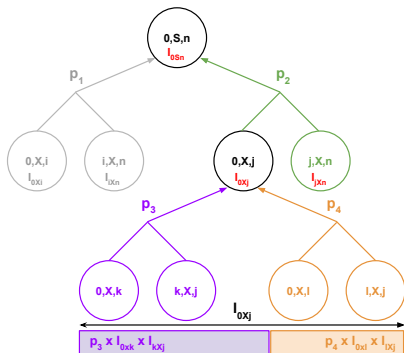
Sampling



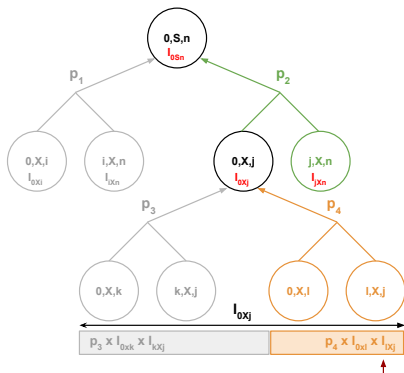
Sampling



Sampling



Sampling



An example for hierarchical models

Model $f(\mathbf{d}) = \sum_i \varphi(e_i)$

where $\varphi(e_i)$ is a weighted combination of local features

Grammar

$X \rightarrow \langle \text{a, the} \rangle$

$X \rightarrow \langle \text{luz, light} \rangle$

$X \rightarrow \langle \text{apague } X_1, \text{switch } X_1 \text{ off} \rangle$

$X \rightarrow \langle X_1 \text{ por favor, please, } X_1 \rangle$

$X \rightarrow \langle X_1 X_2, X_1, X_2 \rangle$

$S \rightarrow \langle \vdash X_1 \neg, \vdash X_1 \neg \rangle$

Input: **apague a luz por favor**

Reference: **please, switch the light off**

Decoding with local features

[illegible]

Decoding with local features

[illegible]

Decoding with local features

[illegible]

Decoding with local features

[illegible]

Decoding with local features

[illegible]

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$			
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					
		$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$				

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					
		$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$				
$X_{1,5}$				$e_8 : \frac{X_{1,3} por favor}{please, X_{1,3}}$		

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$			
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					
		$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$				
$X_{1,5}$		$e_8 : \frac{X_{1,3} por favor}{please, X_{1,3}}$				
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{off}}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$			
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{off}}$					
		$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$				
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$				
		$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$				
$X_{0,5}$		$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$				

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$			
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
	$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$					
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$				
		$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$				
$X_{0,5}$	$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$					
	$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$					

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$				
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$			
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
	$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$					
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$				
		$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$				
$X_{0,5}$		$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$				
		$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$				
		$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$				

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} por favor}{please, X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{apague X_{1,5}}{switch X_{1,5} off}$			
			$e_{12} : \frac{X_{0,3} por favor}{please, X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{the}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{luz}{light}$			
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch X_{1,2} off}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} por favor}{please, X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch X_{1,3} off}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} por favor}{please, X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{apague X_{1,5}}{switch X_{1,5} off}$			
			$e_{12} : \frac{X_{0,3} por favor}{please, X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		$w(e_5)\beta(X_{2,3})$
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					
			$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$			
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		$w(e_5)\beta(X_{2,3})$
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					$w(e_6)\beta(X_{1,3}) \oplus$ $w(e_7)\beta(X_{0,2})\beta(X_{2,3})$
$X_{1,5}$				$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$		
				$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$		
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		$w(e_5)\beta(X_{2,3})$
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					$w(e_6)\beta(X_{1,3}) \oplus$
		$e_7 : \frac{X_{0,2} X_{2,3}}{X_{0,2} X_{2,3}}$				$w(e_7)\beta(X_{0,2})\beta(X_{2,3})$
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			$w(e_8)\beta(X_{1,3}) \oplus$
			$e_9 : \frac{X_{1,2} X_{2,5}}{X_{1,2} X_{2,5}}$			$w(e_9)\beta(X_{1,2})\beta(X_{2,5})$
$X_{0,5}$		$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$				
		$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$				
		$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$				
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		$w(e_5)\beta(X_{2,3})$
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					$w(e_6)\beta(X_{1,3}) \oplus$ $w(e_7)\beta(X_{0,2})\beta(X_{2,3})$
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			$w(e_8)\beta(X_{1,3}) \oplus$ $w(e_9)\beta(X_{1,2})\beta(X_{2,5})$
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			$w(e_{10})\beta(X_{0,2})\beta(X_{2,5}) \oplus$ $w(e_{11})\beta(X_{1,5}) \oplus$ $w(e_{12})\beta(X_{0,3})$
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			

Decoding with local features

NODE	apague	a	luz	por	favor	INSIDE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				$w(e_1)$
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			$w(e_2)$
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{X_{1,2} X_{2,3}}$			$w(e_3)\beta(X_{1,2})\beta(X_{2,3})$
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } X_{1,2} \text{ off}}$					$w(e_4)\beta(X_{1,2})$
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } X_{2,3}}$		$w(e_5)\beta(X_{2,3})$
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } X_{1,3} \text{ off}}$					$w(e_6)\beta(X_{1,3}) \oplus$ $w(e_7)\beta(X_{0,2})\beta(X_{2,3})$
$X_{1,5}$			$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } X_{1,3}}$			$w(e_8)\beta(X_{1,3}) \oplus$ $w(e_9)\beta(X_{1,2})\beta(X_{2,5})$
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{X_{0,2} X_{2,5}}$			$w(e_{10})\beta(X_{0,2})\beta(X_{2,5}) \oplus$ $w(e_{11})\beta(X_{1,5}) \oplus$ $w(e_{12})\beta(X_{0,3})$
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } X_{1,5} \text{ off}}$			
			$e_{12} : \frac{X_{0,3} \text{ por favor}}{\text{please, } X_{0,3}}$			
$S_{0,5}$			$e_{13} : \frac{\vdash X_{0,5} \dashv}{\vdash X_{0,5} \dashv}$			$w(e_{13})\beta(X_{0,5})$

The problem

Most interesting models employ nonlocal features!

- reordering model: previously translated span
- language model: generated strings

The problem

Most interesting models employ nonlocal features!

- reordering model: previously translated span
- language model: generated strings

Example

$$f(\mathbf{d}) = \psi(\text{yield}(\mathbf{d})) + \sum_i \varphi(e_i)$$

where $\psi(\mathbf{y}) = w_\psi \log p_{\text{LM}}(\mathbf{y})$

and $p_{\text{LM}}(\mathbf{y}) = \prod_i p(y_i | y_{i-n+1}^{i-1})$ is an n -gram LM

The problem

Most interesting models employ nonlocal features!

- reordering model: previously translated span
- language model: generated strings

Example

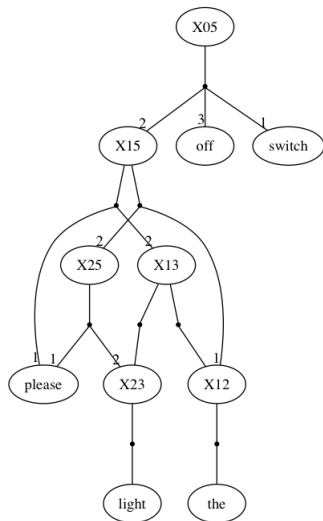
$$f(\mathbf{d}) = \psi(\text{yield}(\mathbf{d})) + \sum_i \varphi(e_i)$$

where $\psi(\mathbf{y}) = w_\psi \log p_{\text{LM}}(\mathbf{y})$

and $p_{\text{LM}}(\mathbf{y}) = \prod_i p(y_i | y_{i-n+1}^{i-1})$ is an n -gram LM

- p_{LM} violates independence assumptions

Illustration of the problem



How do we score the top edge?

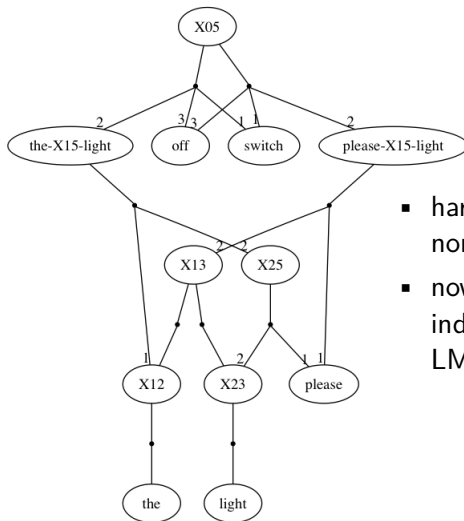
- [switch [please [[the][light]]] off]
- [switch [[the][please [light]]] off]

The solution

“Hard-code” structural dependencies

- disambiguate nodes w.r.t. the context they offer to feature functions
- intuition: we will be “splitting” nodes
- more intuition: nodes must memorise how to complete boundary n -grams

The intuition



- hard-code dependencies through nonterminals
- now we can weight edges independently even with a bigram LM

An example for hierarchical models

Model $f(\mathbf{d}) = \psi(\text{yield}(\mathbf{d})) \sum_i \varphi(e_i)$

where $\varphi(e_i)$ is a weighted combination of local features

$\psi(\text{yield}(\mathbf{d}))$ contains a 3-gram LM

i.e, $p_{\text{LM}_3}(\mathbf{y}) = \prod_i p(y_i | y_{i-2} y_{i-1})$

Grammar

$X \rightarrow \langle \text{a, the} \rangle$

$X \rightarrow \langle \text{luz, light} \rangle$

$X \rightarrow \langle \text{apague } X_1, \text{switch } X_1 \text{ off} \rangle$

$X \rightarrow \langle X_1 \text{ por favor, please, } X_1 \rangle$

$X \rightarrow \langle X_1 X_2, X_1, X_2 \rangle$

$S \rightarrow \langle \vdash X_1 \neg, \vdash X_1 \neg \rangle$

Input: **apague a luz por favor**

Reference: **please, switch the light off**

	apague a luz por favor	LEFT	RIGHT	NODE
				1
				2
				3
				4
				5
				6
				7
				8
				9
				10
				11
				12
				13
				14
				15
				16
				17
				18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{the}$				the	the	1
								2
								3
								4
								5
								6
								7
								8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{the}$				the	the	1
$X_{2,3}$			$e_2 : \frac{luz}{light}$			light	light	2
								3
								4
								5
								6
								7
								8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
								4
								5
								6
								7
								8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

[illegible]

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{the}$				the	the	1
$X_{2,3}$			$e_2 : \frac{luz}{light}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\textcolor{red}{apague } X_{1,2}}{\textcolor{blue}{switch } (1) \textcolor{teal}{off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \textcolor{red}{por favor}}{\textcolor{blue}{please}, (2)}$			please ,	, light	5
								6
								7
								8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

[illegible]

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$			the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
								8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
		$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$				switch the	off light	7
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$				please ,	the light	8
								9
								10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
		$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$				switch the	off light	7
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$				please ,	the light	8
			$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$			the please	, light	9
								10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$				please ,	the light	8
		$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$				the please	, light	9
$X_{0,5}$	$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$					switch the	, light	10
								11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$			the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$		please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$				$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$		please ,	the light	8
				$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$		the please	, light	9
$X_{0,5}$				$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$		switch the	, light	10
				$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$		switch please	light off	11
								12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$			the light	the light	3
$X_{0,2}$		$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$				switch the	the off	4
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$		please ,	, light	5
$X_{0,3}$		$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$				switch the	light off	6
			$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$			switch the	off light	7
$X_{1,5}$				$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$		please ,	the light	8
				$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$		the please	, light	9
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$			switch the	, light	10
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$			switch please	light off	11
			$e_{12} : \frac{\text{apague } X_{1,5}}{\text{switch } (9) \text{ off}}$			switch the	light off	12
								13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{the}$				the	the	1
$X_{2,3}$			$e_2 : \frac{luz}{light}$			light	light	2
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$			the light	the light	3
$X_{0,2}$	$e_4 : \frac{apague X_{1,2}}{switch (1) off}$					switch the	the off	4
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{please, (2)}$		please ,	, light	5
$X_{0,3}$	$e_6 : \frac{apague X_{1,3}}{switch (3) off}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$				$e_8 : \frac{X_{1,3} \text{ por favor}}{please, (3)}$		please ,	the light	8
				$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$		the please	, light	9
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$			switch the	, light	10
			$e_{11} : \frac{apague X_{1,5}}{switch (8) off}$			switch please	light off	11
			$e_{12} : \frac{apague X_{1,5}}{switch (9) off}$			switch the	light off	12
			$e_{13} : \frac{X_{0,3} \text{ por favor}}{please, (6)}$			please ,	light off	13
								14
								15
								16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$			$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$			the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$				$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$		please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$				$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$		please ,	the light	8
				$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$		the please	, light	9
$X_{0,5}$			$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$			switch the	, light	10
			$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$			switch please	light off	11
			$e_{12} : \frac{\text{apague } X_{1,5}}{\text{switch } (9) \text{ off}}$			switch the	light off	12
			$e_{13} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (6)}$			please ,	light off	13
			$e_{14} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (7)}$			please ,	off light	14
								15
								16
								17
								18

[illegible]

[illegible]

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$	$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$					please ,	the light	8
	$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$					the please	, light	9
$X_{0,5}$	$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$					switch the	, light	10
	$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$					switch please	light off	11
	$e_{12} : \frac{\text{apague } X_{1,5}}{\text{switch } (9) \text{ off}}$					switch the	light off	12
	$e_{13} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (6)}$					please ,	light off	13
	$e_{14} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (7)}$					please ,	off light	14
$S_{0,5}$	$e_{15} : \frac{\vdash X_{0,5} \neg}{\vdash (10) \neg}$					\vdash switch	light \neg	15
	$e_{16} : \frac{\vdash X_{0,5} \neg}{\vdash (11) \neg}$ or $e_{17} : \frac{\vdash X_{0,5} \neg}{\vdash (12) \neg}$					\vdash switch	off \neg	16
								17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$				please ,	the light	8
		$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$				the please	, light	9
$X_{0,5}$	$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$					switch the	, light	10
	$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$					switch please	light off	11
	$e_{12} : \frac{\text{apague } X_{1,5}}{\text{switch } (9) \text{ off}}$					switch the	light off	12
	$e_{13} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (6)}$					please ,	light off	13
	$e_{14} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (7)}$					please ,	off light	14
$S_{0,5}$	$e_{15} : \frac{\vdash X_{0,5} \neg}{\vdash (10) \neg}$					\vdash switch	light \neg	15
	$e_{16} : \frac{\vdash X_{0,5} \neg}{\vdash (11) \neg} \text{ or } e_{17} : \frac{\vdash X_{0,5} \neg}{\vdash (12) \neg}$					\vdash switch	off \neg	16
	$e_{18} : \frac{\vdash X_{0,5} \neg}{\vdash (13) \neg}$					\vdash please	off \neg	17
								18

	apague	a	luz	por	favor	LEFT	RIGHT	NODE
$X_{1,2}$		$e_1 : \frac{a}{\text{the}}$				the	the	1
$X_{2,3}$			$e_2 : \frac{\text{luz}}{\text{light}}$			light	light	2
$X_{1,3}$		$e_3 : \frac{X_{1,2} X_{2,3}}{(1) (2)}$				the light	the light	3
$X_{0,2}$	$e_4 : \frac{\text{apague } X_{1,2}}{\text{switch } (1) \text{ off}}$					switch the	the off	4
$X_{2,5}$			$e_5 : \frac{X_{2,3} \text{ por favor}}{\text{please, } (2)}$			please ,	, light	5
$X_{0,3}$	$e_6 : \frac{\text{apague } X_{1,3}}{\text{switch } (3) \text{ off}}$					switch the	light off	6
	$e_7 : \frac{X_{0,2} X_{2,3}}{(4) (2)}$					switch the	off light	7
$X_{1,5}$		$e_8 : \frac{X_{1,3} \text{ por favor}}{\text{please, } (3)}$				please ,	the light	8
		$e_9 : \frac{X_{1,2} X_{2,5}}{(1) (5)}$				the please	, light	9
$X_{0,5}$	$e_{10} : \frac{X_{0,2} X_{2,5}}{(4) (5)}$					switch the	, light	10
	$e_{11} : \frac{\text{apague } X_{1,5}}{\text{switch } (8) \text{ off}}$					switch please	light off	11
	$e_{12} : \frac{\text{apague } X_{1,5}}{\text{switch } (9) \text{ off}}$					switch the	light off	12
	$e_{13} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (6)}$					please ,	light off	13
	$e_{14} : \frac{X_{0,3} \text{ por favor}}{\text{please, } (7)}$					please ,	off light	14
$S_{0,5}$	$e_{15} : \frac{\vdash X_{0,5} \neg}{\vdash (10) \neg}$					\vdash switch	light \neg	15
	$e_{16} : \frac{\vdash X_{0,5} \neg}{\vdash (11) \neg}$		or $e_{17} : \frac{\vdash X_{0,5} \neg}{\vdash (12) \neg}$			\vdash switch	off \neg	16
	$e_{18} : \frac{\vdash X_{0,5} \neg}{\vdash (13) \neg}$					\vdash please	off \neg	17
	$e_{19} : \frac{\vdash X_{0,5} \neg}{\vdash (14) \neg}$					\vdash please	light \neg	18

The problem with the solution

The problem with the solution

Computational complexity!

The problem with the solution

Computational complexity!

- ❶ it seems like the underlying grammar is growing
- ❷ there are way too many n -grams leading to way too many nonterminals
- ❸ the graphical representation (forest) is growing

What is really going on?

We are transferring memory from an automaton to the forest

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

A nonterminal in the forest yields a set of strings

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

A nonterminal in the forest yields a set of strings

- strings project onto paths in the LM automaton

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

A nonterminal in the forest yields a set of strings

- strings project onto paths in the LM automaton
- paths are weighted

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

A nonterminal in the forest yields a set of strings

- strings project onto paths in the LM automaton
- paths are weighted

Nonterminals must be aware of (parts of) the strings they yield

What is really going on?

We are transferring memory from an automaton to the forest

- n -gram LMs can be thought of as an automaton where each state q uniquely represents a k -gram prefix α_q ($k < n$)
- a transition from a state q labelled with word w is weighted by the probability $p(w|\alpha_q)$

A nonterminal in the forest yields a set of strings

- strings project onto paths in the LM automaton
- paths are weighted

Nonterminals must be aware of (parts of) the strings they yield

- they must be annotated with states of the automaton

How hard is it?

Weighted intersection between a wCFG and a wFSA

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata
- weighted sets

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata
- weighted sets

Complexity

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata
- weighted sets

Complexity

- input: $X_0 \rightarrow X_1 X_2 \dots X_a$ where $X_i \in N$

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata
- weighted sets

Complexity

- input: $X_0 \rightarrow X_1 X_2 \dots X_a$ where $X_i \in N$
- output: $X_0^{(q_1, q_a)} \rightarrow X_1^{(q_1, q_2)} X_2^{(q_2, q_3)} \dots X_a^{(q_{a-1}, q_a)}$
where $q_j \in Q$

How hard is it?

Weighted intersection between a wCFG and a wFSA

Generalisation of parsing for

- arbitrary automata
- weighted sets

Complexity

- input: $X_0 \rightarrow X_1 X_2 \dots X_a$ where $X_i \in N$
- output: $X_0^{(q_1, q_a)} \rightarrow X_1^{(q_1, q_2)} X_2^{(q_2, q_3)} \dots X_a^{(q_{a-1}, q_a)}$
where $q_j \in Q$
- complexity: $O(|N||Q|^{a+1})$

Solution

The usual suspect

- pruning

Solution

The usual suspect

- pruning

Alternatives

- local search (greedy methods)
- relaxation techniques
- sampling

Pruning: beam search

Approximate intersection by budgeting the combination of
“comparable” nodes

Pruning: beam search

Approximate intersection by budgeting the combination of “comparable” nodes

- nodes that share structure

Pruning: beam search

Approximate intersection by budgeting the combination of “comparable” nodes

- nodes that share structure
 - phrase-based: coverage vector
 - hierarchical: input spans

Pruning: beam search

Approximate intersection by budgeting the combination of “comparable” nodes

- nodes that share structure
 - phrase-based: coverage vector
 - hierarchical: input spans
- heuristic view of interaction with LM
 - phrase-based: approximate future cost
 - local approximation based on limited context

Naive beam search

- ① enumerate combinations
- ② sort and prune all but the k best

Naive beam search

(a)

$[X, 6, 8; \text{the scheme}]$
 $[X, 6, 8; \text{the plan}]$
 $[X, 6, 8; \text{the project}]$

		1	4	7	
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1	8.2	$[X, 5, 8; \text{from the } \star \text{ the scheme}] : 2.1$
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5	8.5	11.5	$[X, 5, 8; \text{from the } \star \text{ the plan}] : 5.1$
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6	7.7	10.6	13.1	$[X, 5, 8; \text{from the } \star \text{ the scheme}] : 5.5$
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10	11.1	14.3	17.3	$[X, 5, 8; \text{since the } \star \text{ the scheme}] : 7.7$
					\vdots

Cube pruning

An agenda for pruning [Chiang, 2007]

- tries to enumerate combinations in best-first order
- stops after k items have been enumerated
- inspiration: product of sorted lists
- heuristic: assumes the LM a monotone function over edges

Cube pruning

(b)

(b)

			[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
			1	4	7
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1		
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5			
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6				
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10				

			[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
			1	4	7
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1	8.2	
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5	8.5		
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6				
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10				

			[X, 6, 8; the scheme]	[X, 6, 8; the plan]	[X, 6, 8; the project]
			1	4	7
$X \rightarrow \langle \text{cong } X_{\square}, \text{from } X_{\square} \rangle$	1	2.1	5.1	8.2	
$X \rightarrow \langle \text{cong } X_{\square}, \text{from the } X_{\square} \rangle$	2	5.5	8.5		
$X \rightarrow \langle \text{cong } X_{\square}, \text{since } X_{\square} \rangle$	6	7.7			
$X \rightarrow \langle \text{cong } X_{\square}, \text{through } X_{\square} \rangle$	10				

Problem with pruning

- ① unbounded approximation
- ② approximating the Viterbi solution
- ③ incompatible with models which have a probabilistic interpretation
- ④ cannot handle arbitrarily nonlocal dependencies

Beyond beam search

Local search [Hardmeier et al., 2012]

- computationally cheap
- unbounded approximation
- approximate Viterbi
- can handle arbitrarily nonlocal dependencies
- too local view of the distribution (bad for tuning)

Beyond beam search

Relaxation methods

[Chang and Collins, 2011, Rush and Collins, 2011]

- computationally expensive
- bounded approximation
- (approximate) Viterbi
- may handle arbitrarily nonlocal dependencies

Beyond beam search

Sampling [Arun et al., 2009, Aziz et al., 2013, Aziz, 2014]

- (bounded) approximation
- (approximate) Viterbi, expectations
- handle arbitrarily nonlocal dependencies
- in principle ideal for tuning (global view of distribution)
- potentially computationally expensive

Questions?

References I

Abhishek Arun, Chris Dyer, Barry Haddow, Phil Blunsom, Adam Lopez, and Philipp Koehn. Monte Carlo inference and maximization for phrase-based translation. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*, CoNLL '09, pages 102–110, Stroudsburg, PA, USA, 2009. Association for Computational Linguistics. ISBN 978-1-932432-29-9. URL

<http://dl.acm.org/citation.cfm?id=1596374.1596394>.

Wilker Aziz, Marc Dymetman, and Sriram Venkatapathy. Investigations in exact inference for hierarchical translation. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 472–483, Sofia, Bulgaria, August 2013. Association for Computational Linguistics. URL

<http://www.aclweb.org/anthology/W13-2260>.

References II

Wilker Ferreira Aziz. *Exact Sampling and Optimisation in Statistical Machine Translation*. PhD thesis, University of Wolverhampton, 2014.

Yin-Wen Chang and Michael Collins. Exact decoding of phrase-based translation models through Lagrangian relaxation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, EMNLP '11, pages 26–37, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics. ISBN 978-1-937284-11-4.

David Chiang. Hierarchical phrase-based translation. *Computational Linguistics*, 33(2):201–228, June 2007. ISSN 0891-2017. doi: 10.1162/coli.2007.33.2.201. URL <http://dx.doi.org/10.1162/coli.2007.33.2.201>.

References III

Christian Hardmeier, Joakim Nivre, and Jörg Tiedemann.
Document-wide decoding for phrase-based statistical machine translation. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1179–1190, Jeju Island, Korea, July 2012. Association for Computational Linguistics. URL
<http://www.aclweb.org/anthology/D12-1108>.

Alexander M. Rush and Michael Collins. Exact decoding of syntactic translation models through Lagrangian relaxation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, HLT '11, pages 72–82, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics. ISBN 978-1-932432-87-9.