

ESSCaSS'15 day 2 Nelijärve, Estonia 2015.08.19

AI = Learning to Translate Language, Music, and Creativity

Dekai Wu
dekai@cs.ust.hk <http://www.cs.ust.hk/~dekai>

HKUST
Human Language Technology Center
Department of Computer Science and Engineering
University of Science and Technology, Hong Kong




Where does it go?
该放哪儿?
應該擺哪邊呀?
ये कहाँ जाता है?
Bu nereye?
أين يذهب؟
Taruh di mana?
어디 뭘 올까?
Saan mo ilalagay ito?
それはどこに置くのでしょうか？
Où est-ce que ça va?
Vart ska det?

Where does it go?
它在哪里去了?
它在哪儿去了?
यह कहाँ है?
Nereye gidiyor?
أين يذهب؟
Di mana ia pergi?
가 어디에 있습니까?
Saan pumunta?
それはどこに行くのでしょうか？
Où faut-il aller?
Vart tar dot vägen?

we all use
different
frames of reference

Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

word alignment

Where is the Secretary of Finance when needed ?

財政 司 有需要 时 在 那里 ?

Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

word alignment

aim
automatically map tokens between input and output sequences

input

- bisequence (input-output pair of sequences)
- token translation lexicon
- language-independent **bracketing ITG** (only 1 generic nonterminal A)

output

- aligned bisequence

Implementation

- biparser

(Wu, IJCAI 1995)

Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

BITG (bracketing ITG)

$$A \xrightarrow{a} [A A]$$

$$A \xrightarrow{a} \langle A A \rangle$$

$$A \xrightarrow{b_{ij}} u_i/v_j \quad \text{for all } i, j \text{ English-Chinese lexical translations}$$

$$A \xrightarrow{b_{i\epsilon}} u_i/\epsilon \quad \text{for all } i \text{ English vocabulary}$$

$$A \xrightarrow{b_{\epsilon j}} \epsilon/v_j \quad \text{for all } j \text{ Chinese vocabulary}$$

Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

word alignment

Where is the Secretary of Finance when needed ?

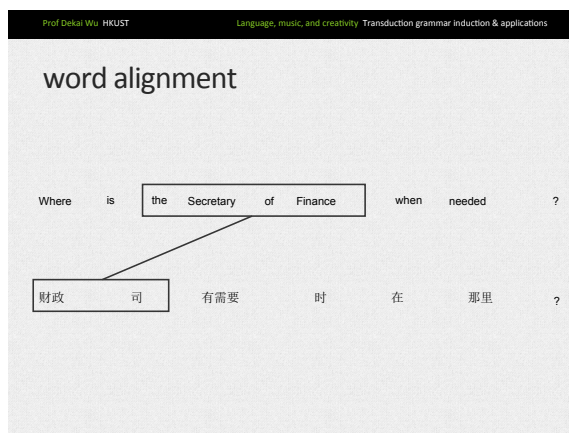
財政 司 有需要 时 在 那里 ?

Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

word alignment

Where is the Secretary of Finance when needed ?

財政 司 有需要 时 在 那里 ?

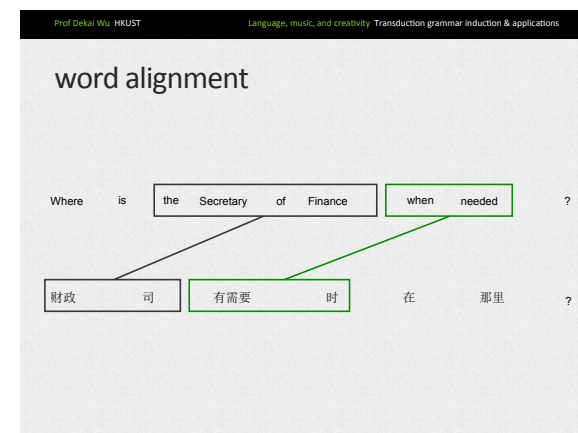


Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

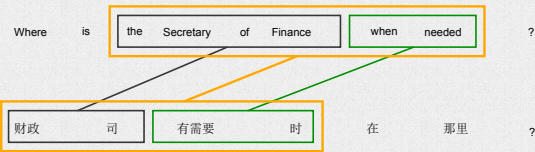
word alignment

Where is the Secretary of Finance when needed ?

財政 司 有需要 时 在 那里 ?



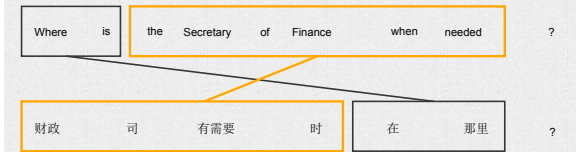
word alignment



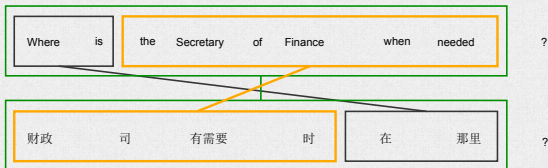
word alignment



word alignment



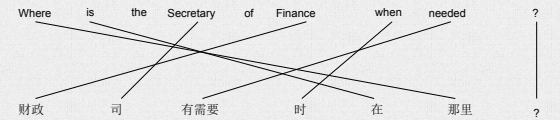
word alignment



word alignment



word alignment



biparsing

aim

automatically parse an input/output sequence pair

input

- bisequence (input-output pair of sequences)
- token translation lexicon
- language-independent **bracketing ITG** (only 1 generic nonterminal **A**)

output

- aligned bisequence

Implementation

- biparser

biparsing

$$\begin{aligned}
 &1. \text{ Initialization} \\
 &\quad A_{1,1:n-1,i,j} = A_{1,1:n-1,i,j} \quad 1 \leq i \leq T \quad (1) \\
 &\quad A_{1,1:n-1,i,j} = A_{1,1:n-1,i,j} \quad 1 \leq i \leq T \quad (2) \\
 &\quad A_{1,1:n-1,i,j} = A_{1,1:n-1,i,j} \quad 1 \leq i \leq T \quad (3) \\
 &2. \text{ Recursion} \\
 &\quad \text{For all } i, j, k, l \text{ such that } \begin{cases} 1 \leq i \leq T \\ 1 \leq j \leq V \\ 1 \leq k \leq T \\ 1 \leq l \leq V \end{cases} \\
 &\quad \quad f_{i,j,k,l} = \max_{1 \leq p \leq T} \{ f_{i,j,p,q} \cdot A_{p,q,k,l} \} \quad (4) \\
 &\quad \quad f_{i,j,k,l} = \begin{cases} 1 & \text{if } f_{i,j,p,q} \geq f_{i,j,k,l} \\ 0 & \text{otherwise} \end{cases} \quad (5) \\
 &\text{where} \\
 &\quad \hat{f}_{i,j,k,l} = \max_{1 \leq p \leq T} \{ \hat{f}_{i,j,p,q} \cdot A_{p,q,k,l} \} \quad (6) \\
 &\quad \begin{bmatrix} \hat{f}_{i,j,k,l} \\ \hat{f}_{i,j,k,l} \\ \hat{f}_{i,j,k,l} \end{bmatrix} = \max_{1 \leq p \leq T} \{ \hat{f}_{i,j,p,q} \cdot A_{p,q,k,l} \} \quad (7) \\
 &\quad \hat{f}_{i,j,k,l} = \max_{1 \leq p \leq T} \{ \hat{f}_{i,j,p,q} \cdot A_{p,q,k,l} \} \quad (8) \\
 &\quad \begin{bmatrix} \hat{f}_{i,j,k,l} \\ \hat{f}_{i,j,k,l} \\ \hat{f}_{i,j,k,l} \end{bmatrix} = \max_{1 \leq p \leq T} \{ \hat{f}_{i,j,p,q} \cdot A_{p,q,k,l} \} \quad (9)
 \end{aligned}$$

2-normal form for ITGs

Lemma 1

For any inversion transduction grammar G , there exists an equivalent inversion transduction grammar G' where $T(G) = T(G')$, such that:

- If $\epsilon \in L_1(G)$ and $\epsilon \in L_2(G)$, then G' contains a single production of the form $S' \rightarrow \epsilon/\epsilon$, where S' is the start symbol of G' and does not appear on the right-hand side of any production of G' ;
- otherwise G' contains no productions of the form $A \rightarrow \epsilon/\epsilon$.

Lemma 2

For any inversion transduction grammar G , there exists an equivalent inversion transduction grammar G' where $T(G) = T(G')$, such that G' does not contain any productions of G' contains either a single terminal-pair or a list of nonterminals.

Lemma 3

For any inversion transduction grammar G , there exists an equivalent inversion transduction grammar G' where $T(G) = T(G')$, such that G' does not contain any productions of the form $A \rightarrow B$.

Theorem 1

For any inversion transduction grammar G , there exists an equivalent inversion transduction grammar G' in which every production takes one of the following forms:

$$\begin{aligned}
 S &\rightarrow \epsilon/\epsilon & A &\rightarrow x/\epsilon & A &\rightarrow [B C] \\
 A &\rightarrow x/y & A &\rightarrow \epsilon/y & A &\rightarrow (B C)
 \end{aligned}$$

translation-driven segmentation

aim

- avoid premature segmentation errors
- optimize chunk boundaries using bilingual/bimodal clues

input

- bisequence (input-output pair of sequences)
- token translation lexicon

output

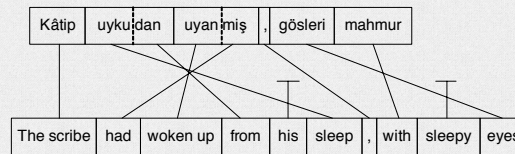
- segmented bisequence

implementation

- segmental biparser, instead of token-based biparser
- integrates segmentation decisions into dynamic programming

(Wu 1997)

translation-driven segmentation



bracketing

aim

- bootstrap grammar induction, with **zero** syntactic knowledge
- use bilingual/bimodal clues to label sequences with tree structure

input

- bisequence (input-output pair of sequences)
- token translation lexicon
- language-independent **bracketing ITG** (only 1 generic nonterminal A)

output

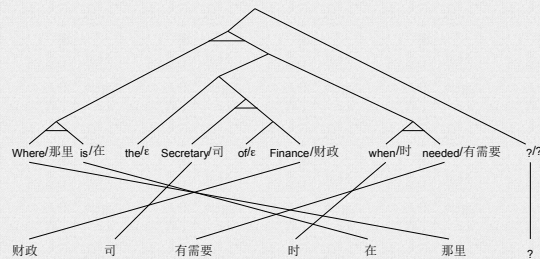
- bracketing tree structure for **both** sequences

implementation

- biparser (choose either token-based or segmental version)

(Wu, ACL 1995)

bracketing



bracketing

[These/這些 arrangements/安排 will/會 enhance/加強 our/我們 (the/的 ability/能力) [to/在/日後 maintain/維持 monetary/金融 stability/穩定 in the years to come/在/] ./.]
[The/他 Authority/管理局 will/將會 (be/他 accountable/負責) [to the/他 (向 Financial/財政 Secretary/司)) ./.]
[They/他們 (are/在 right/正確 enough/十分 to/在 do/做 this/這樣 so/他) ./.]
[([Even/在 more/更 important/重要) [(/在 however/但) [(/在 the/的, is/是 to make the very best of our/在 (善用香港 own/本身 the/的 talent/人才) ./.]
[I/我 hope/在 (望 employers/僱主 will/會 make full/在 (充分善 use/用 [ot/在 those/那些] (the/的工 who/人) [have acquired/在 (學到 new/新 skills/技能) through/透過 this/這個 programme/計劃] ./.]
[I/我 have/已 (at/在 length/詳細 (on/在 how/怎樣 we/我們 (講述) [can/可以 boost/在 (促進 our/本港 the/的 prosperity/繁榮] ./.]

coercion

aim

- bootstrap grammar induction, knowing only **some other** language
- use an English grammar to parse Chinese
- use bilingual/bimodal clues to coerce Chinese into English tree structure

input

- bisequence (input-output pair of sequences)
- token translation lexicon
- input language CFG, with each rule mirrored into straight & inverted ITG rules

output

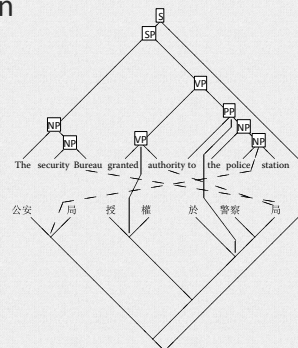
- labeled tree structure for **output** language sequence

implementation

- biparser (choose either token-based or segmental version)

(Wu, WVL 1995)

coercion



projection (bilingual constraint transfer)

aim

- bootstrap grammar induction, projecting constraints from another language
- use an English parse to constrain Chinese tree structure
- use bilingual/bimodal clues to coerce Chinese into English tree structure

input

- bisequence, where input sequence has been parsed
- token translation lexicon
- language-independent **bracketing ITG** (only 1 generic nonterminal A)

output

- labeled tree structure for **output** language sequence

Implementation

- biparser, incorporating constraints into dynamic programming (hard or soft)

learning phrasal translation lexicons

aim

- bootstrap grammar induction, projecting labels from some other language
- use an English parse to constrain Chinese tree structure
- use bilingual/bimodal clues to coerce Chinese into English tree structure

input

- bisequence
- token translation lexicon
- language-independent **bracketing ITG** (only 1 generic nonterminal A)

output

- labeled tree structure for **output** language sequence

Implementation

- biparser, incorporating constraints into dynamic programming (hard or soft)

(Wu, TMI 1995)

learning phrasal translation lexicons

1 % in real	1%的實質
Would you	你是否
an acceptable starting point for this new policy	是可接受為這項新政策的起點
are about 3.5 million	大概有350萬
born in Hong	在香港出生
for Hong	為香港
have the right to decide our	有權決定我
in what way the Government would increase	政府如何增加他們的就業機會及
their job opportunities ; and	
last month	上個月
never to say " never "	不要說“永不”
reserves and surpluses	儲備和盈餘
starting point for this new policy	為這項新政策的起點
there will be many practical difficulties in terms	實行時會有很多實際困難
of implementation	
year ended 31 March 1991	截至一九九一年三月三十一日

tree alignment

aim

automatically map constituents between input and output sequences

Input

- bisequence, where input and output sequences have been parsed
- token translation lexicon
- language-independent **bracketing** ITG (only 1 generic nonterminal A)

output

- recursively aligned bisequence

Implementation

- biparser

transduction (“decoding”)

aim

translate a sequence from an input representation to an output representation, with **zero** syntactic knowledge

Input

- input language sequence
- token translation lexicon
- language-independent **bracketing** ITG (only 1 generic nonterminal A)

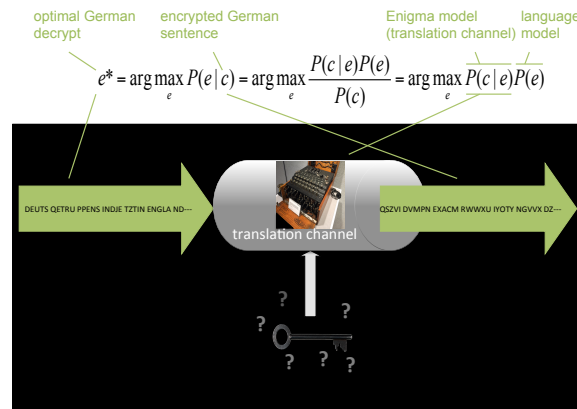
Output

- output language sequence that is a translation of the input sequence

Implementation

- CKY style BITG transducer with output language n-gram model

(Wu, ACL 1996)



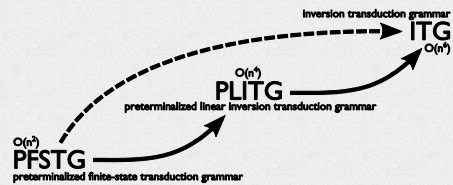
an effective language is **elegant**

an effective representation is **elegant**

Turing machines are
too inefficient to learn



bootstrapping



(Saers, Addanki & Wu, Coling 2012)

wait...

you wanna translate with finite-state transducers?!

- no!
 - FSTs are useless for anything but the most simple, monotone translation tasks
- but...
 - complex translation tasks are made up of simple, monotone translation tasks

review

a **grammar** generates a **language**
... which is a set of *sentences*

a **transduction grammar** generates a **transduction**
... which is a set of *sentence pairs*

Monolingual Languages		SEGMENTAL	Bilingual Transductions	
regular or finite-state languages FSA or CFG that is right or left linear or regular	$O(n^2)$		regular or finite-state transductions FST or SDTG that is right or left linear or regular	$O(n^4)$
linear languages LG or CFG that is linear or unary	$O(n^2)$		linear transductions LTG or SDTG that is linear or unary	$O(n^4)$
context-free languages CFG	$O(n^3)$		inversion transductions ITG or SDTG that is binary or ternary or inverting	$O(n^6)$
			syntax-directed transductions SDTG (or synchronous CFG)	$O(n^{2n+2})$

finite-state

■ finite-state grammar

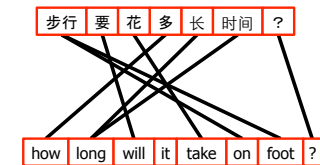
- $S \rightarrow A$
- $A \rightarrow \epsilon$
- $A \rightarrow e B$

■ finite-state transduction grammar

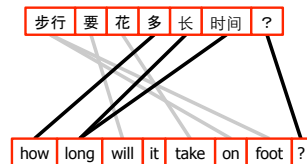
- $S \rightarrow A$
- $A \rightarrow \epsilon/\epsilon$
- $A \rightarrow e/f B$
- $A \rightarrow e/\epsilon B$
- $A \rightarrow \epsilon/f B$

substitution
insertion
deletion

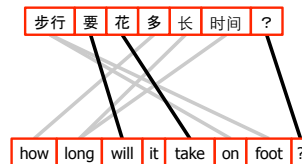
example sentence



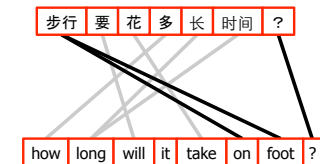
FSTG alignments



FSTG alignments



FSTG alignments



Linear

Linear grammar

- $S \rightarrow A$
- $A \rightarrow \epsilon$
- $A \rightarrow e B$
- $A \rightarrow B e$

Linear transduction grammar

- $S \rightarrow A$ $A \rightarrow [e/f B]$ $A \rightarrow [\epsilon/f B]$ $A \rightarrow \langle \epsilon/f B \rangle$
- $A \rightarrow \epsilon/\epsilon$ $A \rightarrow \langle e/f B \rangle$ $A \rightarrow [e/\epsilon B]$ $A \rightarrow \langle e/\epsilon B \rangle$
- ~~$A \rightarrow \epsilon/f B$~~ $A \rightarrow [B e/f]$ $A \rightarrow [B \epsilon/f]$ $A \rightarrow \langle B \epsilon/f \rangle$
- ~~$A \rightarrow B e/f$~~ $A \rightarrow \langle B e/f \rangle$ $A \rightarrow [B e/\epsilon]$ $A \rightarrow \langle B e/\epsilon \rangle$

Lots of rule forms...

- ... for the same token pairs (biterminals)
- Risks diluting the statistics
- Introduce preterminalized transduction grammars

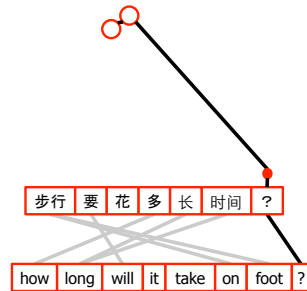
Linear transduction grammar

- $S \rightarrow A$ $S \rightarrow A$
- $A \rightarrow \epsilon/\epsilon$ $A \rightarrow \epsilon/\epsilon$
- $A \rightarrow [e/f B]$ $A \rightarrow [P B], P \rightarrow e/f$
- $A \rightarrow \langle e/f B \rangle$ $A \rightarrow \langle P B \rangle, P \rightarrow e/f$
- $A \rightarrow [B e/f]$ $A \rightarrow [B P], P \rightarrow e/f$
- $A \rightarrow \langle B e/f \rangle$ $A \rightarrow \langle B P \rangle, P \rightarrow e/f$

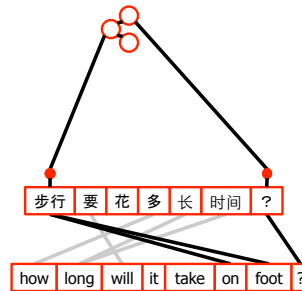
PLITG parsing



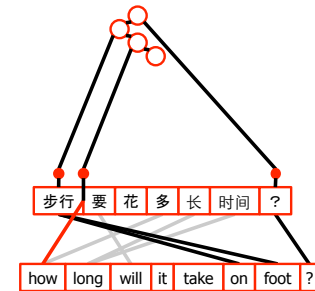
PLITG parsing



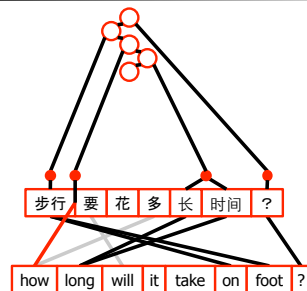
PLITG parsing



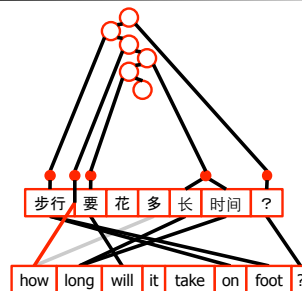
PLITG parsing



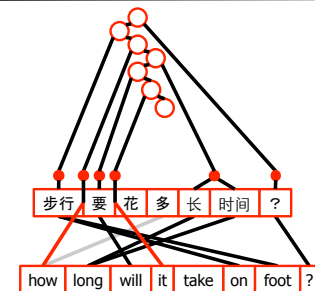
PLITG parsing



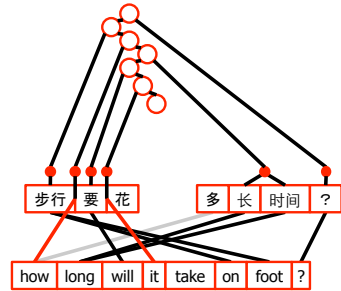
PLITG parsing



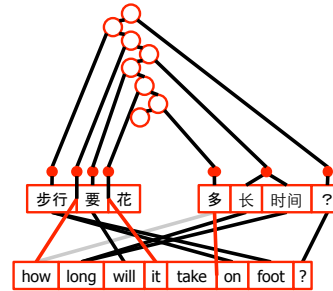
PLITG parsing



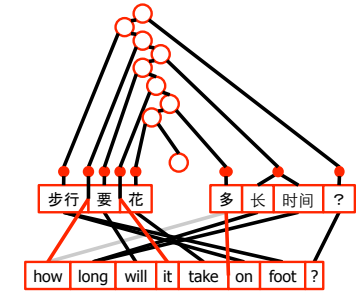
PLITG parsing



PLITG parsing



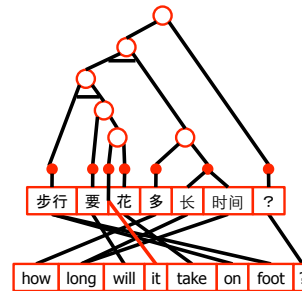
PLITG parsing



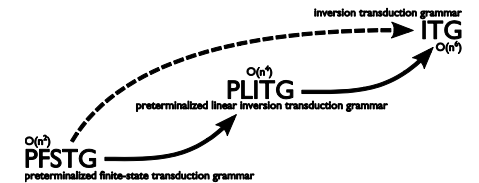
ITG

- context-free grammar (2-normal form)
 - $S \rightarrow A$
 - $A \rightarrow B C$
 - $A \rightarrow e$
- inversion transduction grammar
 - $S \rightarrow A$
 - $A \rightarrow [B C]$
 - $A \rightarrow \langle B C \rangle$
 - $A \rightarrow e/f$
 - $A \rightarrow e/\epsilon$
 - $A \rightarrow \epsilon/f$

ITG parsing



Roadmap



Grammar conversion: PFSTG-PLITG

- | PFSTG | PLITG |
|-----------------------------------|-------------------------------------|
| $S \rightarrow A$ | $S \rightarrow A$ |
| $A \rightarrow \epsilon/\epsilon$ | $A \rightarrow \epsilon/\epsilon$ |
| $A \rightarrow P B$ | $A \rightarrow [P B]$ |
| | $A \rightarrow \langle P B \rangle$ |
| | $A \rightarrow [B P]$ |
| | $A \rightarrow \langle B P \rangle$ |
| $P \rightarrow e/f$ | $P \rightarrow e/f$ |

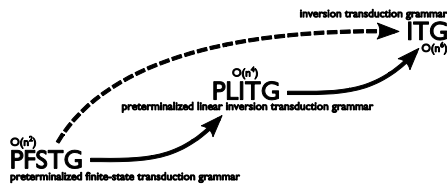
Grammar conversion: PLITG-ITG

- Not as straight forward
- Idea: "promote" preterminals to be proper nonterminals
 - $P \rightarrow e/f$ becomes $P \rightarrow e/f$ and $P \rightarrow A$
 - $p'(e/f|P) = \alpha p(e/f|P)$,
 $p'(A|P) = (1 - \alpha) \beta(A|P) p(e/f|P)$
 - $\alpha = 0.5$,
 - $\beta(A|P) = \text{uniform over nonterminals}$
 - Perform standard grammar normalization to eliminate nullary and unary rules

Grammar conversion: PFSTG-ITG

- Cheat!
- Compose the previous conversions
 PFSTG-PLITG • PLITG-ITG
 gives
 PFSTG-ITG
- No training at the PLITG stage

roadmap



Fun with simple grammars

- Splitting
 - Split one nonterminal or preterminal symbol into n new symbols
 - Apply controlled perturbation to split the probability mass
 - (See the paper for details)

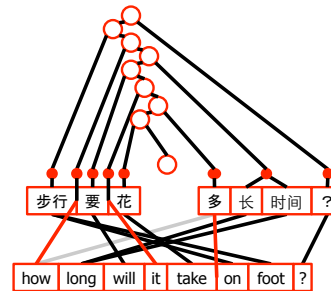
Splitting helps!

- Assume two rules:
 - $A \rightarrow [P A]$
 - $P \rightarrow \text{will/要}$
- We want to:
 - Split A into A and B
 - Split P into P and Q
- Resulting rules:
 - $A \rightarrow [P A], A \rightarrow [P B], A \rightarrow [Q A], A \rightarrow [Q B]$
 - $B \rightarrow [P A], B \rightarrow [P B], B \rightarrow [Q A], B \rightarrow [Q B]$
 - $P \rightarrow \text{will/要}, Q \rightarrow \text{will/要}$

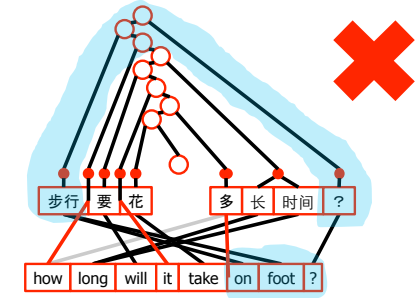
Fun with simple grammars

- Chunking
 - Any sequence of two contiguous terminal productions found in the parse forest could be one production
 - Each round of chunking doubles the potential phrase length
 - Time consuming on large forests
 - (See our EAMT and Interspeech paper from last year for details)

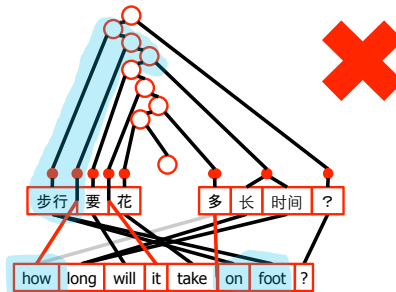
Chunking helps!



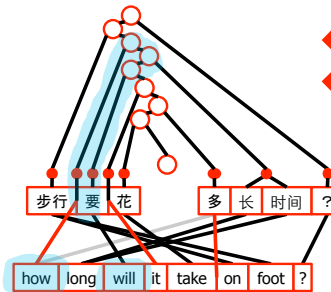
Chunking helps!



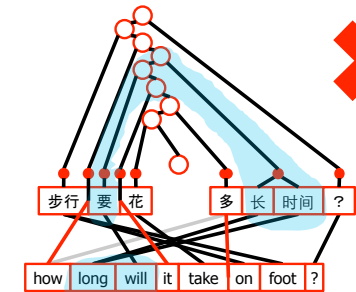
Chunking helps!



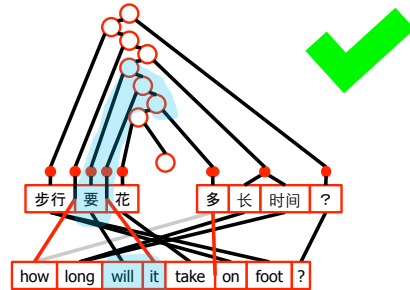
Chunking helps!



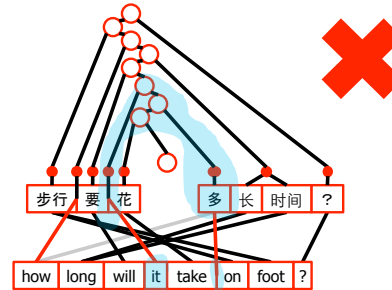
Chunking helps!



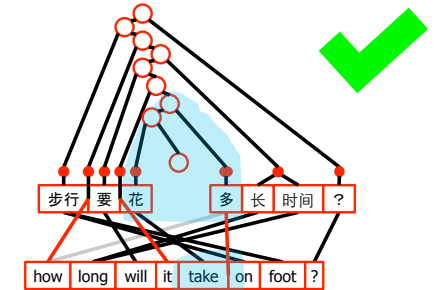
Chunking helps!



Chunking helps!



Chunking helps!



experimental setup

- Initialize a PFSTG from a parallel corpus
 - IWSLT07 Chinese-English
- Chunk and split to your heart's desire
- Move to PLITG then to ITG
 - or move straight to ITG
 - Traverse the roadmap
- Train at each stage
- Measure cross-entropy

why is cross-entropy on the training set important?

- Low cross-entropy may be indicative of over-fitting
- Over-fitting to a lower cross-entropy is indicative of model fit
- PFSTGs are unable to over-fit to the same level that ITGs are able to
- Indicates that ITGs as a model is a better fit to the problem

lessons ways to lower the entropy

- Baseline:** cross-entropy 110.2 (PFSTG)
- Best:** cross-entropy 32.8
PFSTG-chunk-chunk-split-split-split-ITG
 - Nice: still fits in 16Gb RAM
- Promising candidate:** cross-entropy 35.9
PFSTG-chunk-chunk-chunk-ITG
 - After only 1 chunk: 60.7
 - After chunk-chunk: 36.5 (already mostly there!)
 - Unfortunately, couldn't split after chunk-chunk-chunk within 16Gb RAM

lessons ways to lower the entropy

- In general
 - ITG is better than PFSTG and PLITG
 - Chunking helps
 - Splitting helps

remember what we want

- big parallel corpus → small transduction grammar
- learn lexical phrase translations (i.e. a segmental transduction grammar)
- compact generalization of the translation knowledge encoded in the corpus
- unsupervised learning of transduction grammar rules without Giza, Moses, parsers, or anything else

why?

- rearchitecting the SMT core: "Machine Learning 101"
 - Do training and testing on the **same** model
 - Get the **inductive bias** right: core internal representation designed from the start for learning semantic frame generalizations
 - Emphasis on **generalizing** rather than memorizing
 - Minimum description length / MAP → **Occam's razor** for model size
- evaluated in pure, unadulterated form
 - Not as a preprocessing subroutine (eg, for word alignment) within an off-the-shelf "stack-of-hacks" SMT spaghetti architecture
 - ITG decoder matched to ITG learner
 - We'd rather see lower BLEU scores temporarily, so we can better understand transduction grammar induction behavior
 - Don't obscure your model by burying it within a big "stack-of-hacks"!

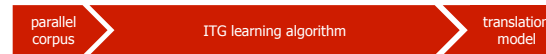
common SMT training pipeline



- Long pipeline propagates errors: risk of premature commitment
 - no way to recover from mistakes in earlier steps
- No credit/blame assignment during learning!
- To try to compensate:
 - massively over-generate phrases
 - result: heavy bias toward memorizing huge corpus instead of learning the right abstract generalizations
- Same pipeline problem for Hierarchical/Syntactic {tree, string} to {tree, string}
 - parse parallel corpus where the trees go
 - replace "phrases" with TG rules



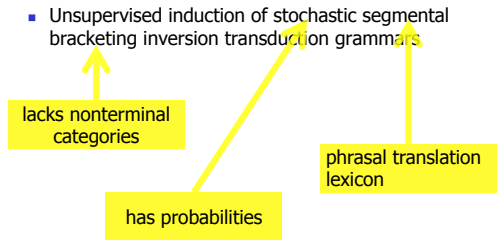
our training "pipeline"



- No pipeline
- No risk of premature commitment
- Replaces many intermediate learning steps
 - ... which, admittedly, have been engineered with heuristic tweaks over a long time
 - worth it to carefully understand correct unsupervised learning of generalizations



our specific SSBITG model



- Unsupervised induction of stochastic segmental bracketing inversion transduction grammars

Our specific SSBITG model

- Unsupervised induction of stochastic segmental bracketing inversion transduction grammars
 - Compact generalization of the translation knowledge encoded in the corpus

Our specific SSBITG model

- Unsupervised induction of stochastic segmental bracketing inversion transduction grammars
 - Compact generalization of the translation knowledge encoded in the corpus
- Bayesian learning objective
 - Closely related to minimum description length
 - Structural prior based on minimum description length
 - Dirichlet parameter prior
 - Fixed model type prior

Our specific SSBITG model

- Unsupervised induction of stochastic segmental bracketing inversion transduction grammars
 - Compact generalization of the translation knowledge encoded in the corpus
- Bayesian learning objective
 - Closely related to minimum description length
 - Structural prior based on minimum description length
 - Dirichlet parameter prior
 - Fixed model type prior
- Structural search based on iteratively segmenting sentence pairs

How to induce transduction rules?

- Can we beat our COLING 2012 bottom-up chunking method?
- Test by using induced transduction grammar directly to translate unseen sentences

Opposite search strategy "directions" for learning transduction rules

- Chunk rules bottom-up (COLING 2012)
 - Start with very small lexical equivalences (biterminals)
 - Look for promising chunks during biparsing
 - Make the chunks explicit biterminals
 - Repeat
- Segment rules top-down (IJCNLP 2013)
 - Start with all sentence pairs as biterminals
 - Segment the existing biterminals into smaller chunks
 - Make the smaller chunks explicit biterminals
 - Repeat

Bottom-up rule chunking strategy

- initialize a token-based FSTG
- parse the bicorpus
- assume that any two adjacent lexical productions could have been generated with a single **chunked** rule, and add them
- train the FSTG using EM
- transform the FSTG into an LITG and train using EM
- transform the LITG into an ITG and train using EM

Opposite search strategy "directions" for learning transduction rules

- **Chunk** rules bottom-up (COLING 2012)
 - Start with very small lexical equivalences (biterminals)
 - Look for promising chunks during biparsing
 - Make the chunks explicit biterminals
 - Repeat
- **Segment** rules top-down (IJCNLP 2013)
 - Start with all sentence pairs as biterminals
 - Segment the existing biterminals into smaller chunks
 - Make the smaller chunks explicit biterminals
 - Repeat

Opposite search strategy "directions" for learning transduction rules

- **Chunk** rules bottom-up (COLING 2012)
 - Start with very small lexical equivalences (biterminals)
 - Look for promising chunks during biparsing
 - Make the chunks explicit biterminals
 - Repeat
- **Segment** rules top-down (IJCNLP 2013)
 - Start with all sentence pairs as biterminals
 - Segment the existing biterminals into smaller chunks
 - Make the smaller chunks explicit biterminals
 - Repeat

Segmentation, intuitively

- **five thousand yen** is my limit
我最多出**五千日元**
- the total fare is **five thousand yen**
总共的费用是**五千日元**

Segmentation, intuitively

- is my limit
我最多出
- the total fare is
总共的费用是
- five thousand yen
五千日元

SHORTER

New search strategy Top-down rule-segmenting strategy

1. **initialize** an ITG containing each sentence pair as a biterminal
2. **foreach** "frequently shared biaffix"
3. **hypothesize** the set of segmentations that it suggests
4. **evaluate** the delta in ITG "goodness" the set would cause*
5. **commit** greedily to good hypothesis sets
6. **goto** 2

* for present purposes, "goodness" = description length

Iterative rule segmentation as phrasal translation lexicon search

Q. How to suggest possible ways to segment rules?

A. Look for frequently shared biaffixes

- Each biaffix suggests a set of rule segmentation hypotheses
 - 4 types of rule segmentations can be suggested



- Can be efficiently computed
- Estimate delta in objective function for each

"Goodness" of an ITG If objective function is description length...

- #bits needed to encode the model $DL(\phi)$
 - ϕ = model = ITG
- plus...
- #bits needed for the model to encode the data $DL(D|\phi)$
 - D = data = parallel training corpus

Top-down rule segmentation

31
symbols

$S \rightarrow A$
 $A \rightarrow$ five thousand yen is my limit/我最多出五千日元

$A \rightarrow$ the total fare is five thousand yen/总共的费用是五千日元

Top-down rule segmentation

31
symbols

$S \rightarrow A$
 $A \rightarrow$ five thousand yen is my limit/我最多出五千日元
 $A \rightarrow \langle AA \rangle$
 $A \rightarrow$ five thousand yen/五千日元
 $A \rightarrow$ is my limit/我最多出
 $A \rightarrow$ the total fare is five thousand yen/总共的费用是五千日元
 $A \rightarrow [AA]$
 $A \rightarrow$ the total fare is/总共的费用是
 $A \rightarrow$ five thousand yen/五千日元



Top-down rule segmentation

24
symbols

$S \rightarrow A$

$A \rightarrow \langle AA \rangle$

$A \rightarrow \text{five thousand yen/五千日元}$

$A \rightarrow \text{is my limit/我最多出}$

$A \rightarrow [AA]$

$A \rightarrow \text{the total fare is/总共的费用是}$

$A \rightarrow \text{five thousand yen/五千日元}$



Description length of an ITG

- Serialize grammar into a message
- Measure number of bits needed to encode the message
- Example:

$S \rightarrow A \quad A \rightarrow \langle AA \rangle \quad A \rightarrow [AA]$

$A \rightarrow \text{have/有} \quad A \rightarrow \text{yes/有} \quad A \rightarrow \text{yes/是}$

- Becomes the message:

$S \langle \rangle A \langle \rangle A \langle \rangle \text{have} \langle \rangle \text{yes} \langle \rangle \text{yes}$

- Assume each symbol requires $-\lg 1/N$ bits (where N is the total number of symbols)
- The above message contains 8 unique symbols \rightarrow 3 bits each
- The message is 23 symbols long, and needs $(23 \cdot 3 =)$ 69 bits to encode



Minimum Description Length objective

- Want: $\underset{\Phi}{\operatorname{argmin}} DL(\Phi, D)$
- $DL(\Phi, D) \propto DL(D|\Phi) + DL(\Phi)$
- $DL(D|\Phi) = -\lg P(D|\Phi)$
- $DL(\Phi) \approx \text{\#bits needed to encode model}$



MAP vs. MDL

- Maximum *a posteriori* probability
- Minimum description length
- Relationship:
 - $DL(x) = -\lg P(x)$
 - $P(x) = 2^{-DL(x)}$
- Enables:
 - probabilistic formulation of description length search
 - description length formulation of probabilistic search



Bayesian search (MAP)

- Maximize *a posteriori* probability of the model given the parallel corpus:

$$P(\Phi|D) = \frac{P(\Phi) P(D|\Phi)}{P(D)}$$

- Decompose the model prior such that:

$$P(\Phi) = P(\Phi_G) P(\Phi_S|\Phi_G) P(\theta_\Phi|\Phi_S, \Phi_G)$$

$$P(\Phi_G) = \text{bracketing inversion transduction grammar}$$

$$P(\Phi_S|\Phi_G) = 2^{-DL(\Phi_S|\Phi_G)}$$

$$P(\theta_\Phi|\Phi_S, \Phi_G) = \text{symmetric Dirichlet distribution } (\alpha=2)$$



Bayesian search (MAP)

- Full search problem:

$$\underset{\Phi_G, \Phi_S, \theta_\Phi}{\operatorname{argmax}} P(\Phi_G) \times P(\Phi_S|\Phi_G) \times P(\theta_\Phi|\Phi_S, \Phi_G) \times P(D|\Phi_G, \Phi_S, \theta_\Phi)$$

- In MDL:

$$\underset{\Phi_G, \Phi_S, \theta_\Phi}{\operatorname{argmin}} DL(\Phi_G) + DL(\Phi_S|\Phi_G) + DL(\theta_\Phi|\Phi_S, \Phi_G) + DL(D|\theta_\Phi, \Phi_S, \Phi_G)$$



evaluating the delta in $P(D|\Phi)$

- Requires biparsing for every hypothesized new Φ
- Intractable
- ... must approximate



estimating the delta in $P(D|\Phi)$

- Assume that r_0 is segmented into r_1, r_2 and r_3
- Approximation assumption:

$$\frac{P(D|\Phi'_S, \Phi_G, \theta_{\Phi'})}{P(D|\Phi_S, \Phi_G, \theta_\Phi)} = \frac{\hat{p}'(r_1) \hat{p}'(r_2) \hat{p}'(r_3)}{\hat{p}(r_0)}$$

- The new rule probability function \hat{p}' will be:

$$\hat{p}'(r_0) = 0$$

$$\hat{p}'(r_1) = \hat{p}(r_1) + \frac{1}{3} \hat{p}(r_0)$$

$$\hat{p}'(r_2) = \hat{p}(r_2) + \frac{1}{3} \hat{p}(r_0)$$

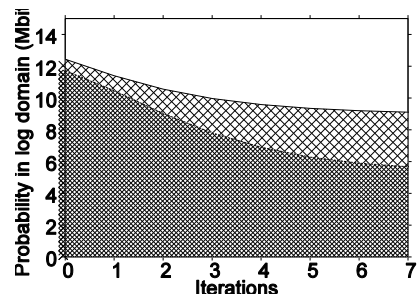
$$\hat{p}'(r_3) = \hat{p}(r_3) + \frac{1}{3} \hat{p}(r_0)$$



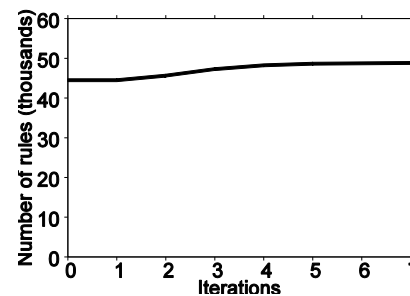
transduction grammar induction by top-down segmenting transduction rules

```
G // The grammar
biaffixes_to_rules // Maps biaffixes to the rules they occur in
biaffixes_delta = [] // Hypothesized biaffixes' impact on P(D|G)
for each biaffix b:
    delta = eval_dl(b, biaffixes_to_rules[b], G)
    if (delta < 0)
        biaffixes_delta.push(b, delta)
sort_by_delta(biaffixes_delta)
for each b:delta pair in biaffixes_delta:
    real_delta = eval_dl(b, biaffixes_to_rules[b], G)
    if (real_delta < 0)
        G = make_segmentations(b, biaffixes_to_rules[b], G)
```

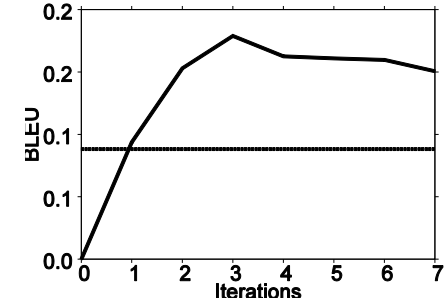

impact of model structure a posteriori probability during learning (log domain)



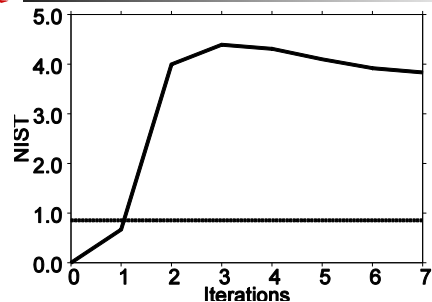
Impact of model structure rule count during learning (log domain)



BLEU score



NIST score



(in table form)

System	NIST	BLEU
baseline	0.8554	8.83
initial	0.0000	0.00
iteration 1	0.6686	9.38
iteration 2	3.9976	15.30
iteration 3	4.3928	17.89
iteration 4	4.3122	16.26
iteration 5	4.0981	16.10
iteration 6	3.9191	15.97
iteration 7	3.8338	15.06

MDL/MAP transduction grammar induction

- Bayesian MAP quite promising for driving ITG induction via top-down segmentation of rules
 - closely related to MDL (DL is more natural for ITG structure prior)
- Beats bilingual lexical chunking driven by ML
 - learns a much smaller ITG...
 - ... that performs better on held-out test data
- New!** even better results combining chunking and segmentation
- Rearchitecting the SMT core: "Machine Learning 101"
 - inductive bias** – internal representation is set up from the start for learning semantic frame generalizations
 - learns small models** – less reliance on memorizing huge corpora
 - rapidly improving** – newer results already into mid-20s BLEU range without Giza, Moses, parsers, or anything else
 - representational transparency** – error analyses to understand learning properties

escaping from blocks world

Can the same exact **core capability** be used not only for conventional AI tasks, but also for all sorts of creative tasks?

transduction grammars simultaneously model Boden's (1992) three types of creativity

combinational new combinations of familiar ideas

exploratory generation of new ideas by exploration of a space of concepts

transformational involves a transformation of the search space so new kinds of ideas can be generated

what
makes
music
music
?

the many languages of music

lyrics	sequences of words
melodies	sequences of notes
chord progressions	sequences of chords
rhythms	sequences of percussive hits
meters	sequences of pulses
ostinatos	sequences of repetitions
stanzas	sequences of lines
verses	sequences of stanzas
songs	sequences of verses
dynamics	sequences of volumes

the many languages of music

lyrics	sequences of words
melodies	sequences of notes
chord progressions	sequences of chords
rhythms	sequences of percussive hits
meters	sequences of pulses
ostinatos	sequences of repetitions
stanzas	sequences of lines
verses	sequences of stanzas
songs	sequences of verses
dynamics	sequences of volumes

if you master one of
these languages,
are you a musician?

it's about the relationships

- the relationships between multiple different musical languages differentiate *aesthetically pleasing music* from jarring noise
- internalizing these relationships accounts for
 - expectation
 - surprise
 - resolution

learning the relationships between different languages

technically: learning the transduction
that relates two languages

Why apply NLP to music?

- Music
 - is a form of language
 - has had major impact across all human cultures
 - emerges from similar cognitive process as speech and written language ([McMullen and Saffran 2004](#))
- Applying NLP to music
 - similar generalizations to be captured in music and natural language
 - distinguishes well-motivated learning methods from language specific fixes
 - adaptation of statistical NLP models presents interesting challenges



Hip hop is spoken language.

- rap** is one of the world's most popular forms of spoken language (for decades!) – arguably spoken language's most significant development in ages
- yet inexplicably ignored in language research
- musical lyrics** way more challenging than classical poetry due to absence of traditional constraints
 - far fewer meter restrictions
 - variable rhyme schemes
 - unusual vocabulary (*Bill Gates* \approx *40 ill dates*)
- user generated content available online
- but off-the-shelf NLP tools not suitably designed



freestyling



freestyling

- freestyle**: improvisational style of rap performed "off the top of the head" with no previously composed lyrics
- freestyle battle**: contest of rappers dueling by challenging and responding using improvised lyrics



freestyling

- freestyle**: improvisational style of rap performed "off the top of the head" with no previously composed lyrics
- freestyle battle**: contest of rappers dueling by challenging and responding using improvised lyrics

Steve Jobs vs. Bill Gates

(source: Epic Rap Battles)

Jobs A man uses the machines you built to sit down and pay his taxes
A man uses the machines I built to listen to the Beatles while he relaxes

Gates Well Steve, you steal all the credit for work that other people do
Did your fat beard Wozniak write these raps for you too?

verse vs. stanza vs. line

All right stop collaborate and **listen**
Ice is back with my brand new **invention**
Something grabs a hold of me **tightly**
Flow like a harpoon daily and **nightly**
Will it ever stop **yo**, I don't **know**
Turn off the lights, and I'll **glow**
To the extreme I rock a mic like a **vandal**
Light up a stage and wax a chump like a **candle**
Dance go rush to the speaker that **booms**
I'm killing your brain like a poisonous **mushroom**
Deadly when I play a dope **melody**
Anything less than the best is a **felony**
Love it or leave it you better gain **weight**
You better hit bull's eye the kid don't **play**
If there was a problem, yo I'll solve **it**
Check out the hook while my DJ revolves **it**

verse

Hong Kong University of Science & Technology

verse vs. stanza vs. line

All right stop collaborate and **listen** stanza AA
Ice is back with my brand new **invention**
Something grabs a hold of me **tightly**
Flow like a harpoon daily and **nightly**
Will it ever stop **yo**, I don't **know** stanza AABA
Turn off the lights, and I'll **glow**
To the extreme I rock a mic like a **vandal**
Light up a stage and wax a chump like a **candle**
Dance go rush to the speaker that **booms**
I'm killing your brain like a poisonous **mushroom**
Deadly when I play a dope **melody**
Anything less than the best is a **felony**
Love it or leave it you better gain **weight**
You better hit bull's eye the kid don't **play**
If there was a problem, yo I'll solve **it**
Check out the hook while my DJ revolves **it**

verse

Hong Kong University of Science & Technology

verse vs. stanza vs. line

line All right stop collaborate and **listen** stanza AA
Ice is back with my brand new **invention**
Something grabs a hold of me **tightly**
Flow like a harpoon daily and **nightly**
Will it ever stop **yo**, I don't **know** stanza AABA
Turn off the lights, and I'll **glow**
To the extreme I rock a mic like a **vandal**
Light up a stage and wax a chump like a **candle**
Dance go rush to the speaker that **booms**
I'm killing your brain like a poisonous **mushroom**
Deadly when I play a dope **melody**
Anything less than the best is a **felony**
Love it or leave it you better gain **weight**
You better hit bull's eye the kid don't **play**
If there was a problem, yo I'll solve **it**
Check out the hook while my DJ revolves **it**

verse

Hong Kong University of Science & Technology



freestyling as transduction

(Wu, Addanki, Beloucif, Saers; EMNLP 2013, Interspeech 2013, LREC 2014, IJCAI 2015, ICMC 2015)



FREESTYLE model hip hop challenge-response as MT

- **approach** – learn to “translate” any line of hip hop to produce improvised, rhyming lyrics
- **requirements** – the underlying model must have
 - strong enough inductive bias for learning to generate responses even without any linguistic or phonetic knowledge
 - sufficient expressive capacity to represent structural relationship between lyrical lines
- **model** – bracketing inversion transduction grammar (BITG) induction
 - rule learning via chunking vs segmentation
 - data selection via adjacent lines vs rhyme scheme detection
 - for English vs Maghrebi French hip hop

Hong Kong University of Science & Technology

140



freestyling by learning transduction grammars



FREESTYLE model learning to transduce challenges into responses

- **transduction** – a relation between languages
- **transduction grammar** – ideal for representing structural relationships between lyrical lines
 - model associations between words/chunks to generate fluent and rhyming responses
- **transduction grammar induction** – learn relation between challenges and responses
 - good inductive bias for generating lyrics without any prior linguistic or phonetic knowledge
- **transduce** – use induced transduction grammar to respond to any challenge rap with an improvised, rhyming response

Hong Kong University of Science & Technology



experiments



transducing challenges to responses

- HKUST ITG decoder
 - bottom-up CKY style parsing algorithm with cube pruning (Chiang 2008)
- monotonic constraint
 - bias to match the rhyming order of challenge

Hong Kong University of Science & Technology



transducing challenges to responses

- HKUST ITG decoder
 - Earley style transducer supports arbitrary non-normal-form, non-binarized ITGs of any rank with mixed terminals/nonterminals in rules (Wu & Wong 1998)
 - segmental ITGs support any length lexical phrase translations (Wu 1997)
 - cube pruning (Chiang 2008)
- penalize reflexive rules like $A \rightarrow \text{yeah} / \text{yeah}$
 - challenge and response identical in ~5% of training instances (chorus lines)
 - danger: reflexive rules that map a line back to itself get induced
 - so decoder penalizes responses too similar to the challenges
 - also decoder can penalize inverted rules as a bias to match rhyming order

challenge	response
hello hello hello hello hello	hello hello hello hello hello
yeah yeah yeah yeah yeah	yeah yeah yeah yeah yeah
they call me superman...	they call me superman...
now can i get it yo	now can i get it yo

Hong Kong University of Science & Technology



issues for evaluation

- systems to be judged on quality of responses
- but what is "quality of responses"?
 - absence of improvised "references"
 - extremely subjective akin to "translation"
 - no automatic metrics like BLEU exist
 - multiple evaluation criterion
 - larger context necessary for evaluation



evaluation criterion

- human evaluators
 - native English speakers
 - frequent hip hop listeners
- humans asked to evaluate system responses on
 - **rhyming** – degree of rhyming with challenge
 - **fluency** – "sing"-ability of the challenge-response pair
- each system response was rated *good*, *acceptable* or *bad* for both criterion



the good, the bad, and the acceptable

challenge	response	fluency	rhyming
cause you and your friends ain't nothing but fiends	you know that you me can you a dream	acceptable	bad
cause you ain't going home till the early morn	and the you this alone i i gotta on	bad	acceptable
created different elements and they travelled abroad	so you rather to the well land	good	acceptable
man i die to see em all thun i just don't care	in the sky and me the in polia and the you there	acceptable	good
what would i do	just me and you	good	good
almost a thought	what is this	good	bad
we fell off into a club to try to pop it to daz	you was a of the love i like a stop the the i	bad	good



corpus & contrastive baseline

- data
 - 52,000 lyrics, 260,000 verses
 - 4.2M tokens and 153,000 token types
 - Small fraction of Arabic, French and Spanish lyrics
- phrase based SMT baseline (PBSMT)
 - evaluate out-of-the-box SMT performance
 - standard Moses baseline
 - 4-gram LM trained on all the lyrics



data selection?



disfluency in hip hop lyrics

- ~10% of data had successive repetitions of words like *the* and *i*
- disfluencies typically result from repetitive chants, exclamations, and interjections in lyrics

chorus style lyrics	"hypeman" style backing vocal lyrics
i i i i	hey hey like like like that like that
oh oh oh oh ahhh	
rock rock rock the boat	
i i i i ice-t ice-t	
yes yes yes a yes yes y'all	

- compare two disfluency handling strategies
 - **filtering** – remove lines with disfluencies
 - **correction** – replace all successive repetitions (*the the the* → *the*)



creating training data

how do we select challenge-response pairs?

- need – lots of training examples consisting of
 - a line of rap, with
 - a fluent and salient rhyming response
- naïve approach – all pairs of lines in the same stanza
 - explodes the training data size
 - does not capture rhyming dependencies
- better approach 1 – all successive line pairs in the same stanza
 - keeps training set size proportional to rap corpus size
 - but still does not ensure that training examples rhyme
- better approach 2 ("**RS**") – only successive line pairs that rhyme
 - but how can we know which line pairs actually rhyme?



rhyme scheme detection

rhyme scheme detection

(Addanki & Wu, SLSP 2013)

- in keeping with our linguistics-lite model – we don't use a pronunciation dictionary
 - hip hop rhyming often defies mainstream pronunciations
 - want a language-independent model
- instead – identify rhyming lines using a **rhyme scheme detector** (Addanki & Wu, SLSP 2013)

generative model for verses

- generated by a fully connected HMM

verse vs. stanza vs. line

line [All right stop collaborate and listen
Ice is back with my brand new invention
Something grabs a hold of me tightly
Flow like a harpoon daily and nightly
Will it ever stop yo, I don't know
Turn off the lights, and I'll glow
To the extreme I rock a mic like a vandal
Light up a stage and wax a chump like a candle
Dance go rush to the speaker that booms
I'm killing your brain like a poisonous mushroom
Deadly when I play a dope melody
Anything less than the best is a felony
Love it or leave it you better gain weight
You better hit bull's eye the kid don't play
If there was a problem, yo I'll solve it
Check out the hook while my DJ revolves it

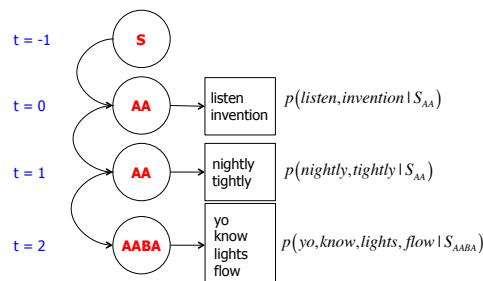
stanza AA
stanza AABA
verse

Hong Kong University of Science & Technology

Hong Kong University of Science & Technology

rhyme scheme detection

(Addanki & Wu, SLSP 2013)



Hong Kong University of Science & Technology

generative model for verses

- generated by a fully connected HMM
- each state S_t is a stanza with a particular rhyme scheme r
- emissions are a sequence of final tokens $x_{1...n}$ in each line of the stanza
- a single textual line of lyrics might contain two lyrical lines separated by a comma

Hong Kong University of Science & Technology

rhyme scheme detection

(Addanki & Wu, SLSP 2013)

- uses an HMM based generative model for verses
- no prior linguistic or phonetic information
- partitions each verse into stanzas with different rhyme schemes
- f-score – 44.06% as evaluated by humans

result — data selection

disfluency correction + rhyme scheme detection

model+disfluency strat.	fluency (good)	fluency (acceptable)	rhyming (good)	rhyming (acceptable)
PBSMT+filtering	4.3%	13.72%	3.53%	7.06%
PBSMT+correction	3.14%	4.70%	1.57%	4.31%
PBSMT+RS+filtering	31.76%	43.91%	12.15%	21.17%
PBSMT+RS+correction	30.59%	43.53%	1.96%	9.02%
TG+filtering	17.25%	46.27%	18.04%	33.33%
TG+correction	21.18%	54.51%	23.53%	39.21%
TG+RS+filtering	28.63%	56.86%	14.90%	34.51%
TG+RS+correction	34.12%	60.39%	20.00%	42.74%

- disfluency correction is better than disfluency filtering
 - improves both fluency and (surprisingly?) rhyming
- rhyme scheme detection
 - improves fluency for both TG and PBSMT models
 - improves the fraction of sentences with \geq acceptable rhyming
 - similar results can be observed for ISTG models also (coming up)

Hong Kong University of Science & Technology

160



chunking vs segmentation?

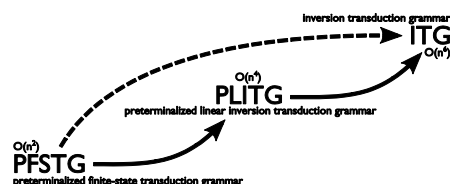
transduction grammar induction

transduction grammar induction

- bracketing inversion transduction grammars (BITG) (Wu 1997)
 - empirically high coverage, accuracy across various NLP tasks
 - sufficient expressiveness to handle word associations between lines
 - efficient induction and decoding algorithms
- compare two approaches, trained on same amount of data
 - TG token-based BITG
 - ISTG interpolated segmental BITG
- token-based BITG captures word associations better
 - efficient induction algorithm using
 - EM
 - beam pruning
 - bootstrapped induction
- segmental BITG captures phrasal associations better
 - efficient induction algorithm using
 - greedy iterative segmentation of transduction rules
 - MDL driven induction

(Tsai, Addanki & Wu, COLING 2012)

(Saeedi, Addanki & Wu, SSST 2013)



- transduction grammar rules can contain
 - rules** $A \rightarrow [B A \text{ "long"/"wrong"}]$
 $A \rightarrow \langle A \text{ "felt bad"/"what I really had"} B \rangle$
 - lexical rules** $A \rightarrow \text{"long"/"wrong"}$
 $A \rightarrow \text{"felt bad"/"what I really had"}$
 - structural rules** $A \rightarrow [A A]$
 $A \rightarrow \langle A B \rangle$

- transduction grammar rules can contain
 - rules** $A \rightarrow [B A \text{ "long"/"wrong"}]$
 $A \rightarrow \langle A \text{ "felt bad"/"what I really had"} B \rangle$
 - lexical rules** $A \rightarrow \text{"long"/"wrong"}$
 $A \rightarrow \text{"felt bad"/"what I really had"}$
 - structural rules** $A \rightarrow [A A]$
 $A \rightarrow \langle A B \rangle$
- token-based** transduction grammars
 - biterminals contain at most one token in each language
 $A \rightarrow e/f \mid e/\varepsilon \mid \varepsilon/f$
 - simple and efficient learning algorithms
 - suffer from a lack of fluency in the output

- transduction grammar rules can contain
 - rules** $A \rightarrow [B A \text{ "long"/"wrong"}]$
 $A \rightarrow \langle A \text{ "felt bad"/"what I really had"} B \rangle$
 - lexical rules** $A \rightarrow \text{"long"/"wrong"}$
 $A \rightarrow \text{"felt bad"/"what I really had"}$
 - structural rules** $A \rightarrow [A A]$
 $A \rightarrow \langle A B \rangle$

- token-based** transduction grammars
 - biterminals contain at most one token in each language
 $A \rightarrow e/f \mid e/\varepsilon \mid \varepsilon/f$
 - simple and efficient learning algorithms
 - suffer from a lack of fluency in the output

- transduction grammar rules can contain
 - rules** $A \rightarrow [B A \text{ "long"/"wrong"}]$
 $A \rightarrow \langle A \text{ "felt bad"/"what I really had"} B \rangle$
 - lexical rules** $A \rightarrow \text{"long"/"wrong"}$
 $A \rightarrow \text{"felt bad"/"what I really had"}$
 - structural rules** $A \rightarrow [A A]$
 $A \rightarrow \langle A B \rangle$

- segmental** transduction grammars
 - biterminals can contain multiple tokens in each language
 $A \rightarrow e_{0.1}/f_{0.2} \mid e_{0.1}/\varepsilon \mid \varepsilon/f_{0.2}$ (Wu 1997)
 - models phrases for more fluent output
 - but unsupervised learning algorithms not as well developed...

- unfortunately there still aren't many algorithms for segmental transduction grammar induction that
 - are purely unsupervised
 - assume the same matched training and testing models
- most SMT approaches instead
 - use heuristics to extract segments from token alignments
 - typ. justified by BLEU motivations not applicable to our task

- segmental ITG learned using a top-down rule segmenting approach driven by a MAP/minimum description length objective (Saers, Addanki & Wu, IJCNLP 2013)

- segmental ITG learned using a top-down rule segmenting approach driven by a MAP/minimum description length objective (Saers, Addanki & Wu, IJCNLP 2013)
- aim: balance fluency and expressivity by interpolating the token-based ITG G_a with a segmental ITG G_b
 - interpolation parameter α controls mixture
 $p_{a+b}(r) = \alpha p_a(r) + (1 - \alpha) p_b(r)$

challenge	man i die to see em all thun i just don't care
FREESTYLE	in the sky and me the in pollan and the you there
PBSMT	man i live to see em all i just don't care
challenge	did a twelve year bid in the streets and held it down
FREESTYLE	to the girls here kid the and to the thought the now
PBSMT	did a year in the streets and it down
challenge	oh i believe in yesterday
FREESTYLE	can you see the day
PBSMT	oh i believe in tomorrow
challenge	what would i do
FREESTYLE	just me and you
PBSMT	what would you do
challenge	cause you ain't going home till the early morn
FREESTYLE	and the you this alone i i gotta on
PBSMT	cause you and your friends aint nothing but

transduction grammar rule	log prob.
$A \rightarrow \text{long/wrong}$	-11.6747
$A \rightarrow \text{rhyme/time}$	-11.6604
$A \rightarrow \text{felt bad/couldn't see what i really had}$	-11.3196
$A \rightarrow \text{matter what you say/leaving anyway}$	-11.8792
$A \rightarrow \text{arhythmic/this rhythm is sick}$	-12.3492

model	fluency ($\geq \text{good}$)	fluency ($\geq \text{acceptable}$)	rhyming ($\geq \text{good}$)	rhyming ($\geq \text{acceptable}$)
PBSMT	3.14%	4.70%	1.57%	4.31%
TG	21.18%	54.51%	23.53%	39.21%
ISTG	26.27%	57.64%	27.45%	48.23%
PBSMT+RS	30.59%	43.53%	1.96%	9.02%
TG+RS	34.12%	60.39%	20.00%	42.74%
ISTG+RS	30.98%	61.18%	30.98%	53.72%

- interpolated segmenting TG (ISTG) produces more fluent responses than token-based TG on both data sets (more later)
- ISTG also produces better rhyming responses (surprising?)
- results also demonstrate that off-the-shelf phrase-based SMT systems (PBSMT) cannot be directly adopted for this task



Maghrebi French hip hop

- advantage of our linguistics-lite model is that it should work independent of language
- test if it can learn to generate response lyrics in languages other than English
- initial experiments on Maghrebi French hip hop lyrics
- our model performs surprisingly well despite
 - no special adaptation
 - much smaller training data size
 - diverse and noisy training data

- about 1300 hip hop song lyrics
- majority in **Maghrebi French**: French interspersed with romanized **Arabic**, **Berber** and **English** phrases
 - De la traversée du désert au bon couscous de **Yéma** (*Yéma = My mother*)
 - a **yemmi ino** = *my son* / a **thizizwith** = *a bee*
 - T'es **game over**, **game over** / Le son de **Chicken wings**
- linguistically complex – language dependent models will be hard to adapt
- 47,000 sentence pairs selected using rhyme scheme detection

- FREESTYLE produces responses rated
 - good* – 9.2% (fluency), 14.5% (rhyming)
 - acceptable* – 30.2% (fluency), 38% (rhyming)
- encouraging results despite small and diverse training corpus causing sparse data issues
- error analysis on responses indicates
 - realization of less common rhyme schemes like **ABAB**
 - responses with semantically related terms

transduction grammar rule	log prob.
$A \rightarrow \text{terre/la guerre}$	-9.4837
$A \rightarrow \text{haine/peine}$	-9.77056
$A \rightarrow \text{mal/pays natal}$	-10.6877
$A \rightarrow \text{je frissonne/mi corazon}$	-11.0931
$A \rightarrow \text{gratteurs/rappeurs}$	-11.7306

challenge	Si je me trompe
response	faut que je raconte
challenge	Un jour je suis un livre
response	et ce que je de vivre
challenge	Pacha mama ils ne voient pas ta souffrance
response	Combat ni leur de voulait de la décadence
challenge	le palestine n'était pas une terre sans peuple
response	le darfour d'autre de la guerre on est
challenge	Une banlieue qui meut
response	les yeux et

- FREESTYLE: first known model for learning how to rap battle
 - transduces challenge lyrics to improvised responses
 - transduction grammar induction is fully unsupervised
 - learns fluent, rhyming responses absent linguistic knowledge
- new MDL driven learning of hip hop ITG rules by segmentation
 - segmental grammars improve on token-based grammars
 - both outperform off-the-shelf PBSMT contrastive baseline
- hip hop domain specific models against very noisy training data yield more fluent and better rhyming responses
 - data selection via unsupervised rhyme scheme detection model
 - disfluency correction
- completely unsupervised generation of hip hop challenge responses without linguistic knowledge despite the noisy domain
- generalizes to non-English hip hop: encouraging results on Maghrebi French validate language independence assumptions



maximum *a posteriori* objective

$$\begin{aligned} & \underset{\Phi}{\operatorname{argmax}} P(\Phi) P(D | \Phi) \\ & \quad \text{prior} \quad \text{bilingual corpus likelihood} \\ & P(\Phi) = P(\Phi_G) P(\Phi_S | \Phi_G) P(\theta_\Phi | \Phi_S, \Phi_G) \\ & \quad \text{MTG} \quad \text{Dirichlet prior} \\ & \text{DL prior}^* \quad P(\theta_\Phi | \Phi_S, \Phi_G) = \prod_{i=0}^{N-1} \frac{1}{B(\alpha_0, \alpha_1, \dots, \alpha_{R_{n_i}-1})} \prod_{j=0}^{R_{n_i}-1} \theta_\Phi^{n_i}(j) \\ & -\log_2 (P(\Phi_S | \Phi_G)) \propto \text{DL}(\Phi_S) \end{aligned}$$

putting everything together, we want

$$\underset{\Phi_G, \Phi_S, \theta_\Phi}{\operatorname{argmax}} P(\Phi_G) P(\Phi_S | \Phi_G) P(\theta_\Phi | \Phi_S, \Phi_G) P(D | \Phi_G, \Phi_S, \theta_\Phi)$$

*prior constructed in terms of the grammar's description length



challenges with symbolic ITGs

- unsupervised ITG induction remains hard
 - extremely large model space!
- nonterminals don't capture context efficiently
 - intractable – each context has explicit nonterminal
 - solution: replace nonterminals with feature vectors
- new idea: apply **TRAAM**, a distributed representation for ITGs we began developing for statistical MT (Addanki & Wu 2014)
 - *bilingual* recursive neural network model
 - uses *both* input & output language contexts



modeling recursive structures

- TRAAM goes beyond neural network approaches that model **monolingual recursive structures**
- neural language models and SRNs (Bengio *et al.* 2003)
 - contextual history modeled by a RNN
- convolutional networks (Collobert & Weston 2008)
 - learn vector representations of words
 - used in NLP tasks such as POS tagging, chunking and SRL
- RAAM and recursive autoencoders (RAE) (Pollack 1990, Socher *et al.* 2011)
 - can be more flexible than convolutional networks
 - RAEs have been successfully applied to sentiment prediction



bilingual vector space models

- dearth of vector space models for **compositional** learning of bilingual relations
- predominantly augment “shake-n-bake” SMT modeling assumptions using feature vectors
- n-gram translation models (Son *et al.*, 2012)
 - bilingual generalization of class based n-grams using distributed representations
 - fails to model compositionality and cross-lingual reordering
- bilingual word embeddings (Zou *et al.*, 2013)
 - recurrent NNLM model with SMT word alignments
 - only learns non-compositional features



bilingual vector space models

- NNLMs + input language context (Devlin *et al.* 2014)
 - does not model input and output language features simultaneously
- recurrent probabilistic models (Kalchbrenner & Blunsom 2013)
 - generates an input sentence representation that generates an output sentence
 - lacks structural constraints and relies on a LM to reorder output
- reordering prediction using RAEs (Li *et al.* 2013)
 - **monolingual** RAEs to predict reordering in a maxent ITG model
 - uses only input language context

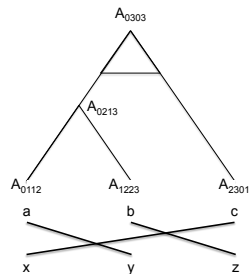


TRAAM model definition

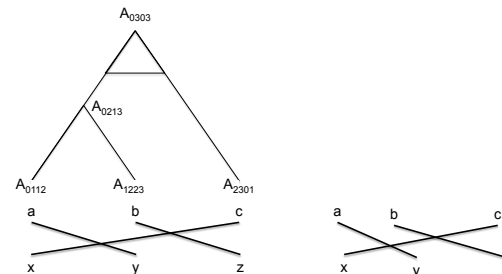
- bilingual recursive neural network
 - fully bilingual generalization of monolingual RAAM model
 - vectors represent **bilingual** constituents
- uniform feature vector dimension
 - task-dependent representation learning
 - similar biconstituents have similar vectors
 - feature vector clusters represent **soft bilingual categories**
- generates feature vectors recursively from smaller biconstituents
 - language bias via dimensionality reduction



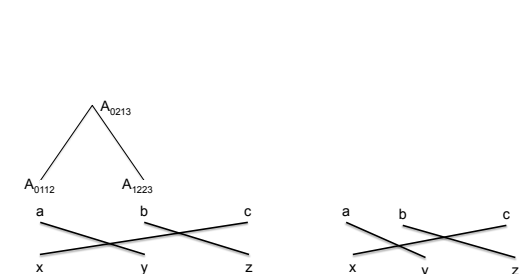
BITG → TRAAM



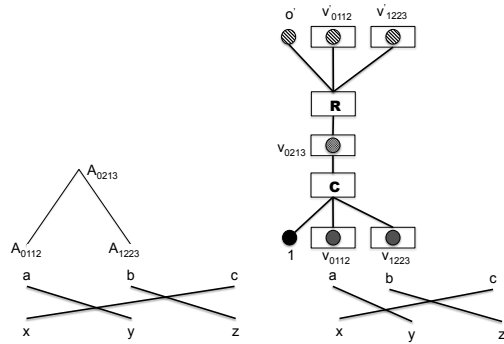
BITG → TRAAM



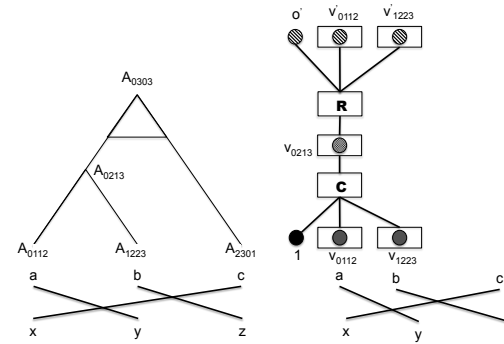
BITG → TRAAM



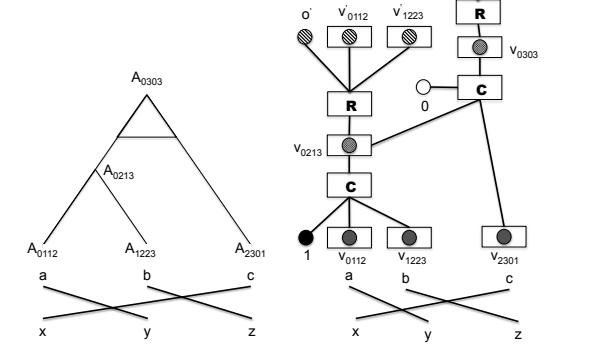
BITG → TRAAM



BITG → TRAAM



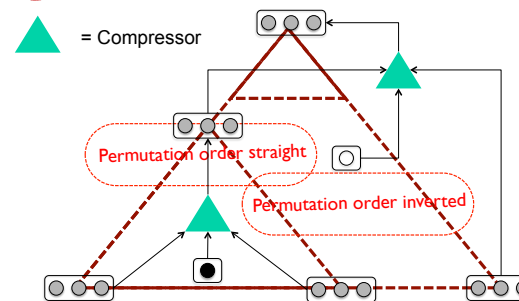
BITG → TRAAM



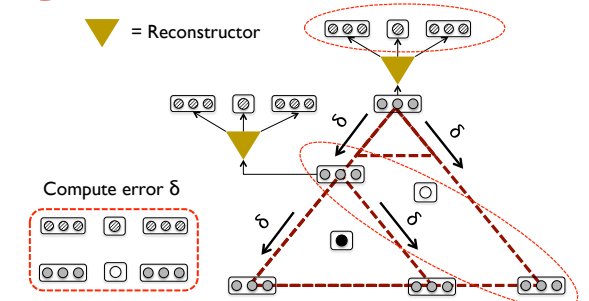
TRAAM model training

- TRAAM network contains a compressor and a reconstructor network
 - generalizes RAAM to represent **bilingual** sequences
 - for **bracketing** ITGs, reordering can be represented by a single bit (straight vs. inverted)
- biparses from a BITG are used to learn the network weights
 - compressor network generates feature vectors recursively
 - objective – to ensure the context of children is captured efficiently
 - reconstructor network provides the loss function
- backpropagation with structure is used to compute gradients
 - L-BFGS can be used to optimize network weights (Goller & Culcher 1996)

TRAAM forward propagation



TRAAM backpropagation



TRAAM improvisation using transduction engine

- reconstruction error used as a feature for transduction (decoding)
 - provides the context that symbolic ITGs lack
 - used in a log-linear combination with grammar and LM score
- transduction heuristics similar to Wu *et al.* (2013) also applied
 - rules mapping a surface form to itself were penalized
 - singleton rules were penalized
- our ITG transducer (decoder) is capable of handling features
 - cube pruning (Chiang 2008) used to generate k-best hypotheses
 - feature weights were determined via manual inspection on a small development set

challenges in flamenco

- no clear boundary between music and dance
- "constrained improvisation"
- regular and irregular hypermetrical structures
- rapid switching between 3/4 and 6/8 meters
- heavy syncopation
- sudden, misleading off-beat accents and patterns
- frequent eliding of downbeat accents (which humans and automatic meter-finding algorithms typically rely on)
- expert musicians rely on**
 - complex hypermetrical knowledge
 - syncopated meter patterns
 - irregular real
- to dynamically recognize when to switch meters/patterns**

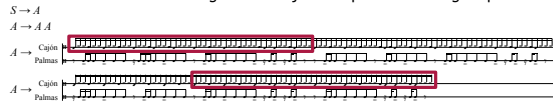
learning flamenco hypermetrical structure

- simultaneous learning of
 - metrical structure
 - hypermetrical structure
 - multipart structural relations
- learn the relationship between parallel frames of reference



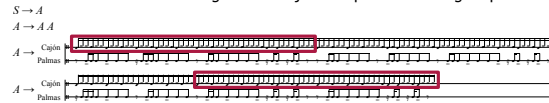
learning flamenco hypermetrical structure

initial transduction grammar is just the parallel training corpus



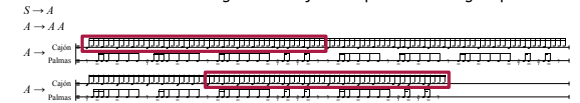
learning flamenco hypermetrical structure

initial transduction grammar is just the parallel training corpus



learning flamenco hypermetrical structure

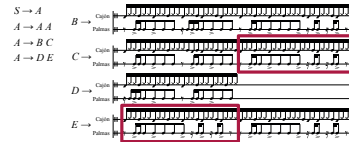
initial transduction grammar is just the parallel training corpus



learning flamenco hypermetrical structure



learning flamenco hypermetrical structure



learning flamenco hypermetrical structure



learning flamenco hypermetrical structure

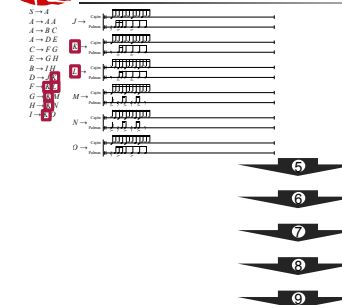
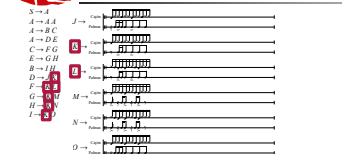
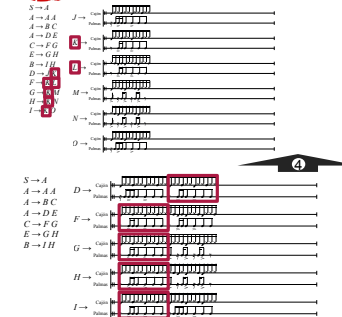
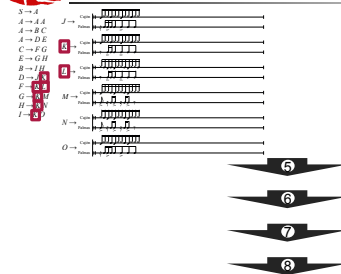
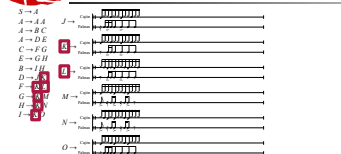
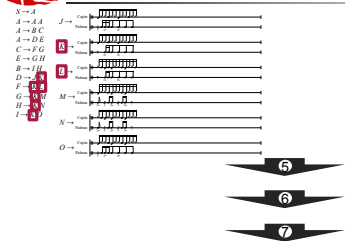
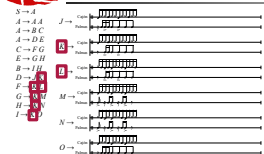
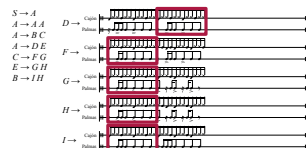


learning flamenco hypermetrical structure

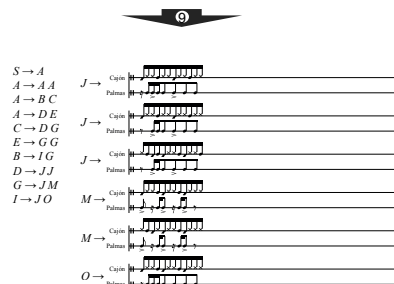


learning flamenco hypermetrical structure

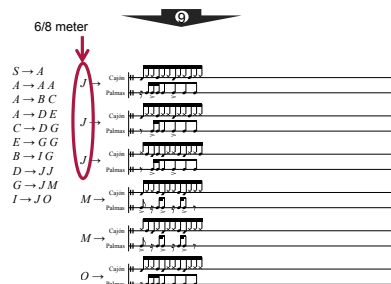




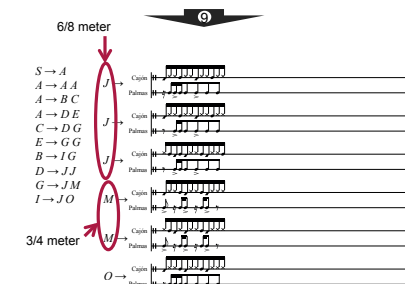
learning flamenco hypermetrical structure



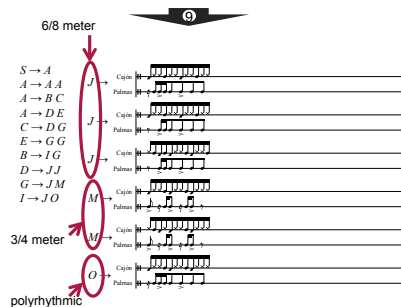
learning flamenco hypermetrical structure



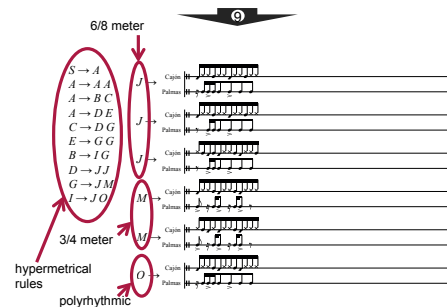
learning flamenco hypermetrical structure



learning flamenco hypermetrical structure



learning flamenco hypermetrical structure



Prof Dekai Wu HKUST Language, music, and creativity Transduction grammar induction & applications

escaping from blocks world

The same exact core capability of transduction grammar induction and application appears effective not only for conventional AI tasks, but also for all sorts of creative tasks.

ESSCaSS'15 day 2 Neljäärve, Estonia 2015.08.19

AI = Learning to Translate Language, Music, and Creativity

Dekai Wu
dekai@cs.ust.hk <http://www.cs.ust.hk/~dekai>

HKUST
Human Language Technology Center
Department of Computer Science and Engineering
University of Science and Technology, Hong Kong