

Combining corpus statistics and knowledge base to disambiguate and acquire verb frames

學生：高定慧

指導教授：張俊盛 博士

系所：資訊工程學系

時間：2013.07.09

Abstract

- Verb frames may play a key role in language learning
- Purpose
 - Automatically generate verb frames with semantic labels
- Method
 - Combining corpus statistics and knowledge base
 - Disambiguate semantic labels
- Results
 - We achieved quite satisfied performance comparing to a manually built gold standard

A verb frame

consists of one or more parts which express the requirements for their possible participants.

Tesnière (1953)

Verb frames in different forms

sentence: “I finally **abandon** the wild idea”

Linggle →	I	abandon	idea
VerbNet →	np	abandon	np
VerbNet →	subject	abandon	object
Dictionary →	somebody	abandon	something
FrameFinder →	person	abandon	cognition
CPA →	person	abandon	plan

Learning verb frames is important

For example, a language learner who is already familiar with:

“members + *establish* + friendship”

Learners tend to rephrase using
similar words

establish

建立

construct

建立

build

建立



synonym



share the same translation

Unfortunately, this often leads to a word choice error

members establish friendship

* members construct friendship

members build friendship

Dictionaries typically overly simplified semantic labels

1 to build something

construct something from/of/in something

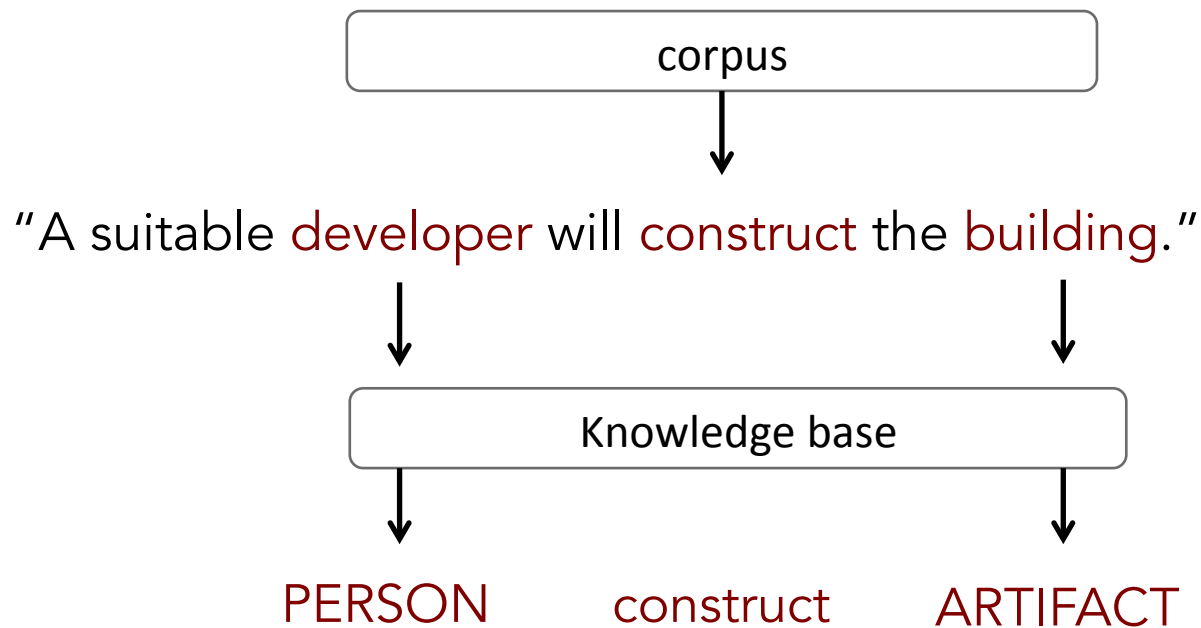
2 to form something

More specific semantic labels help one to use a verb more appropriately

PERSON establish RELATIONSHIP

PERSON construct ARTIFACT

We can easily derive such verb frames
from a corpus and a knowledge base



verb show instance > pattern %

1	PERSON construct ARTIFACT	27	29.35 %
	him construct kiln	●	Grumpily Dai Huang chose Li Lu , to help him construct a newer , bigger
	you construct shelter	●	There may have been so much conflict indoors that you had to construct
	I construct bomb	●	As for Meredith-Lee 's death , if you 're asking me did I construct a bomb answer is no .
2	COGNITION construct ARTIFACT	18	19.57 %
	intention construct road	●	Taken in conjunction the written statement and the key diagram indicating an A fifty nine relief road , passing to the north of Harrogate and Knaresborough
	techniques construct houses	●	For the first time , archaeologists have been able to study in detail the medieval builders to construct the typical cob houses for which the West
	example construct building	●	Would a trust , for example to construct a public building or to set up a be enforced as such ?
3	PERSON construct COGNITION	13	14.13 %
	I construct tactic	●	I 've always been too shy to construct a tactic in order to attract women
	I construct idea	●	In a split-second I would build on this particle of noise and construct a that could produce such a phenomenon .
	he construct understanding	●	Unlike the Chicago School , for example , he was not attempting to construct how distinct forms of natural area were resulting from people 's search
4	PERSON construct COMMUNICATION	8	8.7 %
	you construct movies	●	Flip through the wacky but informative manual which casts you as a man that the first three programs let you construct 20-frame icon-sized movies
	we construct database	●	This we have used to construct and maintain a database of management
	he construct discourse	●	Furthermore , he was attempting to construct a discourse which was in attack then being launched from inside and outside the discipline .
5	COGNITION construct COMMUNICATION	6	6.52 %

Ambiguity existed in a verb frame.

→ *Word Sense Disambiguation*

"workers **abandon** plants"



Many resources provide verb frames information in various forms

Linggle →	I	abandon	idea
VerbNet →	np	abandon	np
VerbNet →	subject	abandon	object
Dictionary →	somebody	abandon	something
FrameFinder →	person	abandon	cognition
CPA →	person	abandon	plan

VerbNet

FRAMES

REFKEY

NP V NP

EXAMPLE"David constructed a house."

SYNTAXAGENT V RESULT

SEMANTICSNOT(EXIST(START(E), RESULT)) EXIST(RESULT(E), RESULT) CAUSE(AGENT, E)

NP V NP PP.MATERIAL

EXAMPLE"David constructed a house out of sticks."

SYNTAXAGENT V RESULT {FROM OUT_OF} MATERIAL

SEMANTICSNOT(EXIST(START(E), RESULT)) EXIST(RESULT(E), RESULT) MADE_OF(RESULT(E), RESULT, MATERIAL) CAUSE(AGENT, E)

NP V NP PP.BENEFICIARY

EXAMPLE"David dug a hole for me."

SYNTAXAGENT V RESULT {FOR} BENEFICIARY

SEMANTICSNOT(EXIST(START(E), RESULT)) EXIST(RESULT(E), RESULT) CAUSE(AGENT, E) BENEFIT(E, BENEFICIARY)

NP V NP PP.ATTRIBUTE

EXAMPLE"They designed the Westinghouse-Mitsubishi venture as a non-equity transaction."

SYNTAXAGENT V RESULT (AS) ATTRIBUTE

SEMANTICSNOT(EXIST(START(E), RESULT)) EXIST(RESULT(E), RESULT) CAUSE(E, AGENT)

[next page](#)

No.	Chunks	Frequency	Examples (from BNC)	Parent (more general versions)
1	[noun sg] construct a [noun sg] of	20		
2	[noun] construct a [noun sg] of	34		
3	[noun] be constructed from [noun]	15		
4	[noun] construct a [noun sg]	91		
5	[noun sg] construct a [noun sg]	57		
6	[noun] be constructed in [noun]	9		
7	[noun] be constructed [prep] the [noun]	28		
8	[noun] be constructed by [noun]	12		
9	[noun sg] of constructing a [noun sg]	16		
10	the [noun sg] of constructing [noun]	7		
11	[noun] constructed by the [noun]	10		
12	[noun] construct a new [noun sg]	12		
13	[noun] be constructed [prep] the [noun sg]	22		

Corpus Pattern Analysis (CPA)

Human **construct** Artifact | Building

Human **construct** Theory | Hypothesis

Human **build** Relationship

Human **build** Building | Machine

Requires expert linguists to manually
derive patterns from a corpus

1. time-consuming
2. might not achieve high coverage for many
verbs a learner has to master

<u>students</u> construct meaning	8 %	1,000	+
<u>students</u> construct knowledge	8 %	1,000	+
<u>technologies</u> construct detail	4 %	500	+
<u>Students</u> construct meaning	4 %	500	+
<u>building</u> construct forklifts	4 %	500	+
<u>living</u> construct subtype	4 %	500	+
<u>attributes</u> construct error	3 %	400	+

Linggle: a Web-scale Linguistic Search Engine for Words in Context
 - Boisson et al. 2013

MEANING (7)

BUILDINGS (5)

MAPS (3)

GRAPHS (3)

ROADS (4)

NARRATIVES (3)

COMPUTER (3)

LINES (2)

BAR (2)

BUILDINGS

construct buildings	11,400	2 %	+
construct bridges	2,000	< 1%	+
construct dams	1,500	< 1%	+
construct facilities	9,200	2 %	+
construct schools	1,200	< 1%	+
construct school	900	< 1%	+
construct homes	4,500	< 1%	+
construct houses	3,600	< 1%	+
construct house	700	< 1%	+
construct complex	7,500	2 %	+

(more 10...)

Learners need more complete verb frames or patterns

Linggle →

1 PERSON apologize

I apologize

juvenile apologize

2 PERSON apologize to PERSON

you apologize to me

I apologize to him

3 PERSON apologize for ACT

I apologize for needing

I apologize for using

-FrameFinder

An automatically generating verb frames approach

1. Extract verb arguments - corpus
2. Obtain probable semantic roles for each argument - Knowledge base
3. Disambiguate senses of an argument - WSD
4. Tally and output verb frames

Step 1. Extracting verb arguments based on grammatical relations

"The deer would eat your plants"

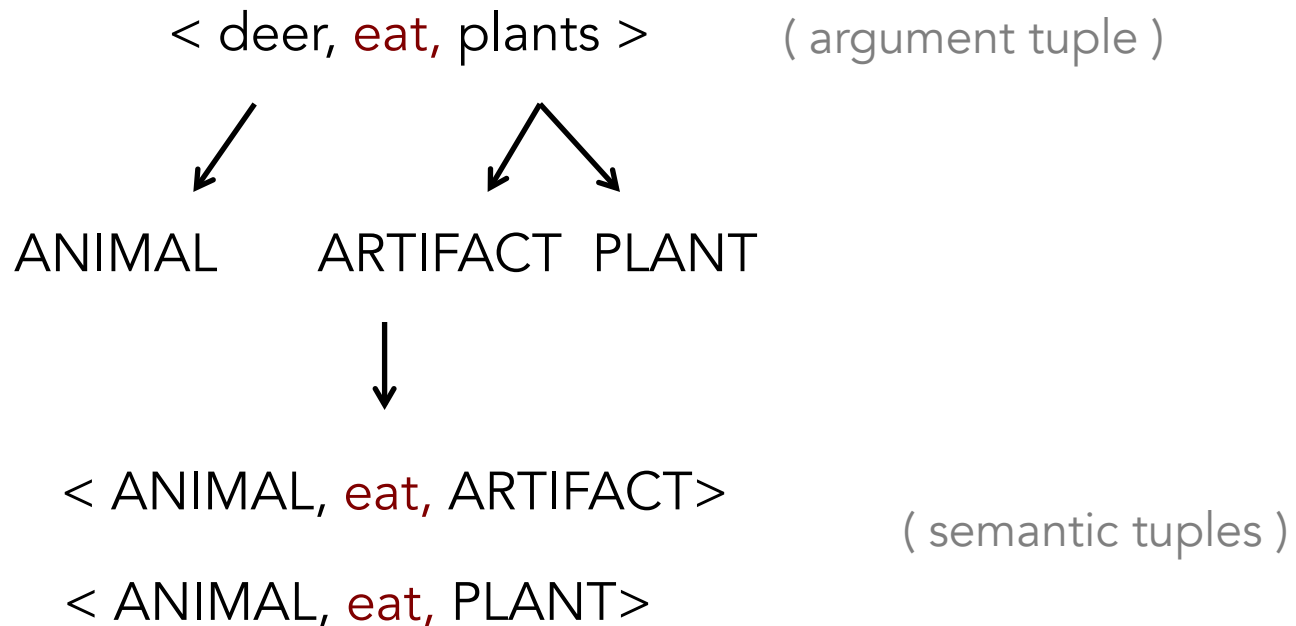


(subject) "The deer would eat your plants" (object)



< deer, eat, plants > (argument tuple)

Step 2. Obtaining semantic roles from a knowledge base



Step 3. Generating verb frames

<i>argument tuples</i>	<i>semantic tuples</i>
deer, eat, plants	< ANIMAL, eat, PLANT > < ANIMAL, eat, ARTIFACT >
birds, eat, seeds	< ANIMAL, eat, PLANT > < ANIMAL, eat, PERSON >
ants, eat, sugar	< ANIMAL, eat, FOOD >

Disambiguate and identify the *intended meaning* of arguments

< deer, eat, plants >



< ANIMAL, eat, ARTIFACT >

< ANIMAL, eat, PLANT >

Applying the Expectation-Maximization (EM) algorithm to disambiguate

- EM algorithm - Dempster et al., 1977
- Expectation (E-step)
- Maximization (M-step)

A flexible and extensible disambiguation module

$$p(v_f|v) = \alpha(p(v_f|v)) + \beta(\underset{\substack{\uparrow \\ \text{subject}}}{p(s|v)} \cdot \underset{\substack{\uparrow \\ \text{object}}}{p(o|v)} \cdot \prod_{po \in \text{adverbial}} \underset{\substack{\uparrow \\ \text{adverbial}}}{p(po|v)})$$

In this study, we set $\alpha = 1$, $\beta = 0$

Procedure GenerateVerbFrames(*ArgumentTuples*)

(1a) *FrameCandis* = {}

(1b) *allFrames* = {getVerbFrames(*t*) for each tuple *t* in *ArgumentTuples*}

(1c) $p(\text{frame} | v) = 1 / |\text{allFrames}|$

while not convergent

 count(*frame* | *v*) = 0 for *frame* in *allFrames*

 for each tuple *t* in *ArgumentTuples*

(2a) *tupleFrame* = getVerbFrames(*t*)

(2b) *nf* = sum(prob(*frame*) for all *frame* in *tupleFrames*)

 for each *frame* in *tupleFrames*

(2c) count(*frame*) += prob(*frame*) / *nf*

 for all *frame* in *allFrames*

(3) $p(\text{frame} | v) = \text{count}(\text{frame} | v) / |\text{ArgumentTuples}|$

 for each tuple *t* in *ArgumentTuples*

(4) append argmax(*frame*) to *FrameCandis*

(5) *RankedFrames* = Sort(*frames* in *FrameCandis* in decreasing order of frequency)

return top *K* ranking *RankedFrames*

EM Example

<i>argument tuples</i>	<i>semantic tuples</i>	
deer, eat, plants (a_1)	< ANIMAL, eat, ARTIFACT >	s_1
	< ANIMAL, eat, PLANT >	s_2
birds, eat, seeds (a_2)	< ANIMAL, eat, PLANT >	s_2
	< ANIMAL, eat, PERSON >	s_3
ants, eat, sugar (a_3)	< ANIMAL, eat, FOOD >	s_4

1st E-step: Initialize the probability of each candidate *uniformly*

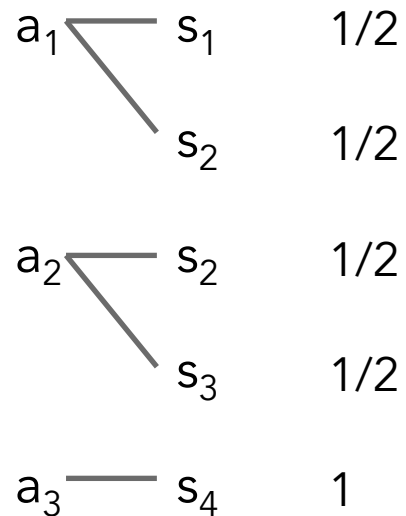
a_1 — s_1 $1/4 \rightarrow 1/2$
 \ s_2 $1/4 \rightarrow 1/2$

a_2 — s_2 $1/4 \rightarrow 1/2$
 \ s_3 $1/4 \rightarrow 1/2$

a_3 — s_4 $1/4 \rightarrow 1$

→ normalize

1st M-step:



$$s_1 : 1/2 / 3 = 1/6$$

$$s_2 : (1/2 + 1/2) / 3 = 1/3$$

$$s_3 : 1/2 / 3 = 1/6$$

$$s_4 : 1 / 3 = 1/3$$

2nd E-step:

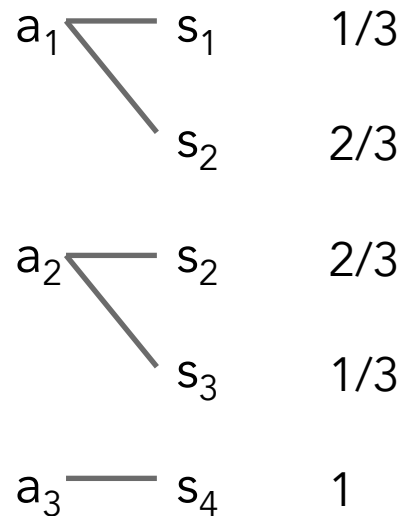
a_1 — s_1 $1/6 \rightarrow 1/3$
 \ s_2 $1/3 \rightarrow 2/3$

a_2 — s_2 $1/3 \rightarrow 2/3$
 \ s_3 $1/6 \rightarrow 1/3$

a_3 — s_4 $1/3 \rightarrow 1$

→ normalize

2nd M-step:



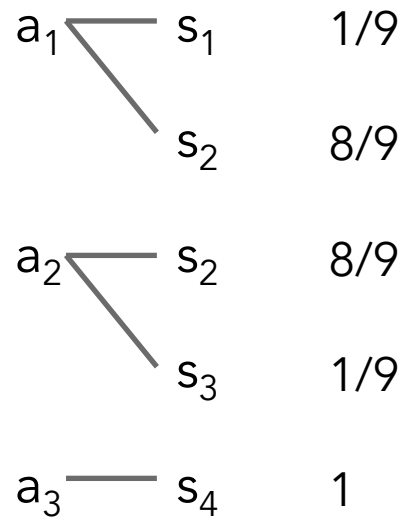
$$s_1 : 1/3 / 3 = 1/9$$

$$s_2 : (2/3 + 2/3) / 3 = 4/9$$

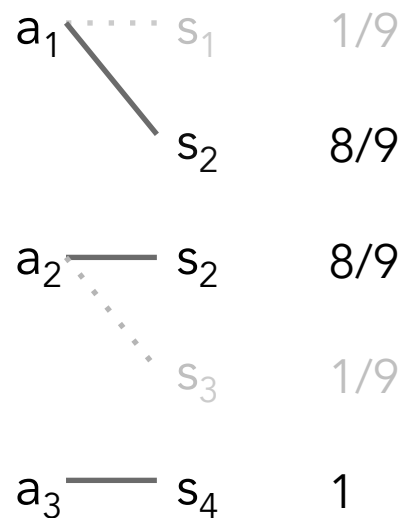
$$s_3 : 1/3 / 3 = 1/9$$

$$s_4 : 1 / 3 = 1/3$$

Until convergent...



Assign a semantic tuple for each argument tuple



Finally, counts and outputs the ranked semantic tuples

ANIMAL eat PLANT	s_2	2
ANIMAL eat FOOD	s_4	1
< ANIMAL, eat, ARTIFACT >	s_1	0
< ANIMAL, eat, PERSON >	s_3	0

Verb frame generation methods compared

- Most Frequent Sense (**MFS**)
 - always chooses the major sense
 - e.g., < deer, eat, plants > → < ANIMAL, **eat**, ARTIFACT >
- Expectation-Maximization: (**EM**)
 - considers probabilities of semantic categories

The extracted frames for the average language learners should be:

1. Valid
2. Cover common usages of the verb

Evaluate *FrameFinder's* performance using

1. Precision
2. Weighted recall
3. F-measure

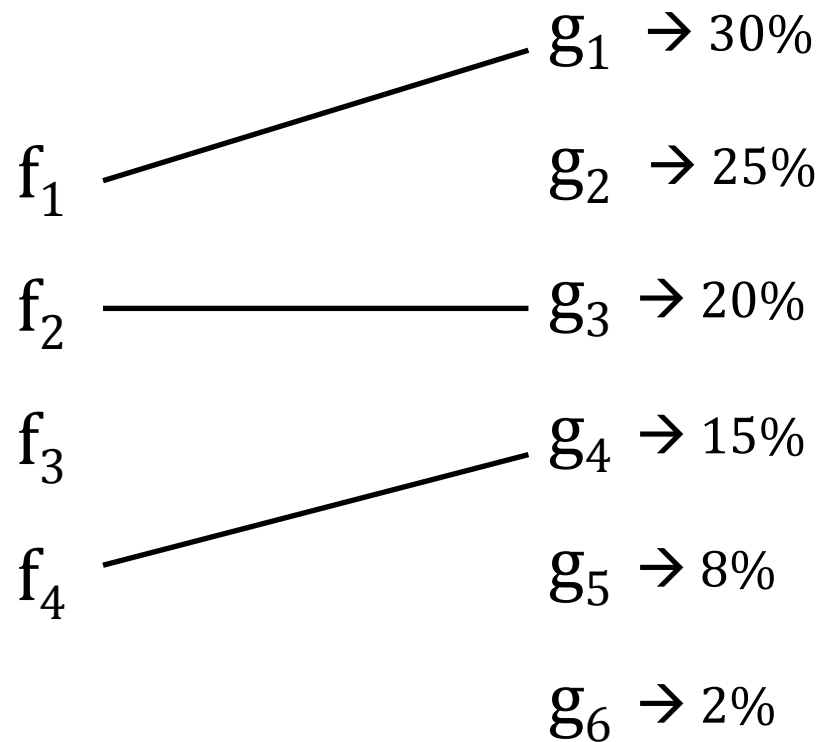
$$\textit{Precision}(v) = \frac{|F(v) \cap \textit{Gold}(v)|}{|F(v)|}$$

$$\textit{wtRecall}(v) = \sum_{f \in F(v) \cap \textit{Gold}(v)} \frac{\textit{freq}(f)}{|\textit{Gold}(v)|}$$

$$\textit{F - measure} = \frac{2PR}{P + R}$$

where $P = \textit{Precision}$, $R = \textit{wtRecall}$

For example, given a verb v



$$\textit{Precision} = 0.6$$

$$\begin{aligned}\textit{wtRecall} &= 0.3 + 0.2 + 0.15 \\ &= 0.65\end{aligned}$$

$$\textit{F-measure} = 0.62$$

Resources and tools we used

- British National Corpus (BNC)
 - 4,693,767 sentences
- Stanford Parser
 - subject, object, adverbial- p.22 Table 5
- WordNet
 - 25 supersenses - p.22 Table 4
- Corpus Pattern Analysis (CPA)
 - 812 verbs
 - 3,100 verb patterns

Mapping CPA labels to WordNet supersenses

e.g.,

HUMAN | INSTITUTION **abandon** PLAN | ACTIVITY

PERSON **abandon** COGNITION

PERSON **abandon** ACT

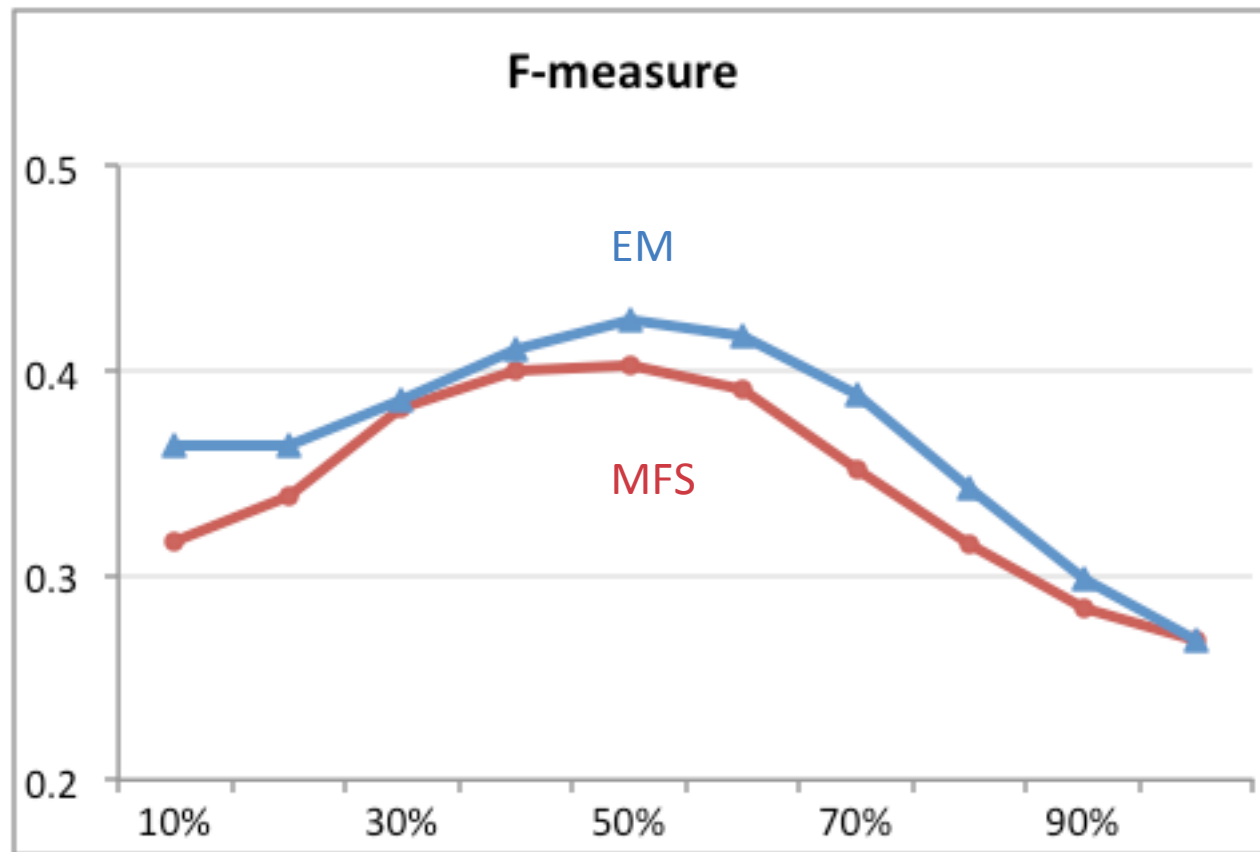
GROUP **abandon** COGNITION

GROUP **abandon** ACT

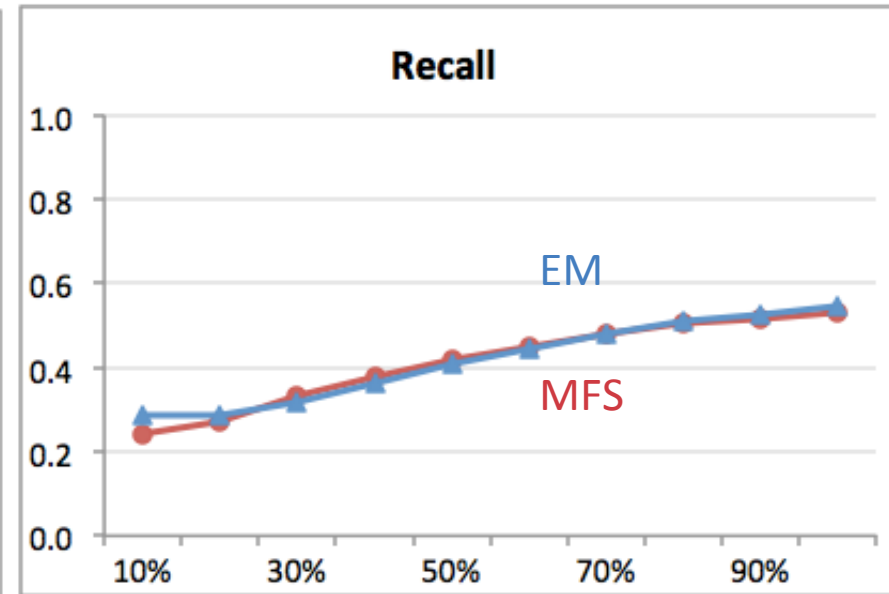
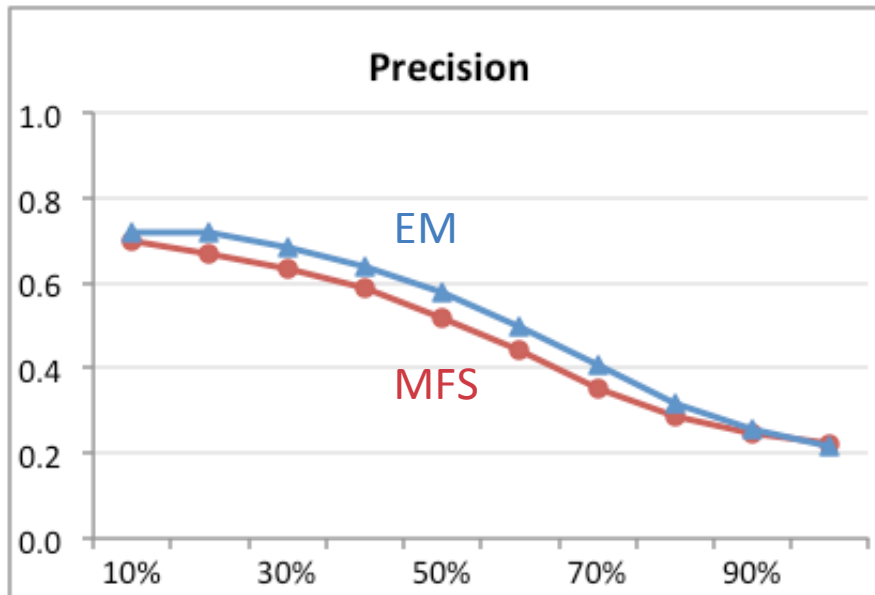
Calculate the coverage

Verb frame		Coverage	Accumulated coverage
PERSON	abandon COGNITION	28.3%	28.3%
40% →	PERSON abandon ACT	20.2%	48.5%
	PERSON abandon PERSON	18.18%	66.7%
	⋮	⋮	⋮
	⋮	⋮	⋮
	⋮	⋮	⋮

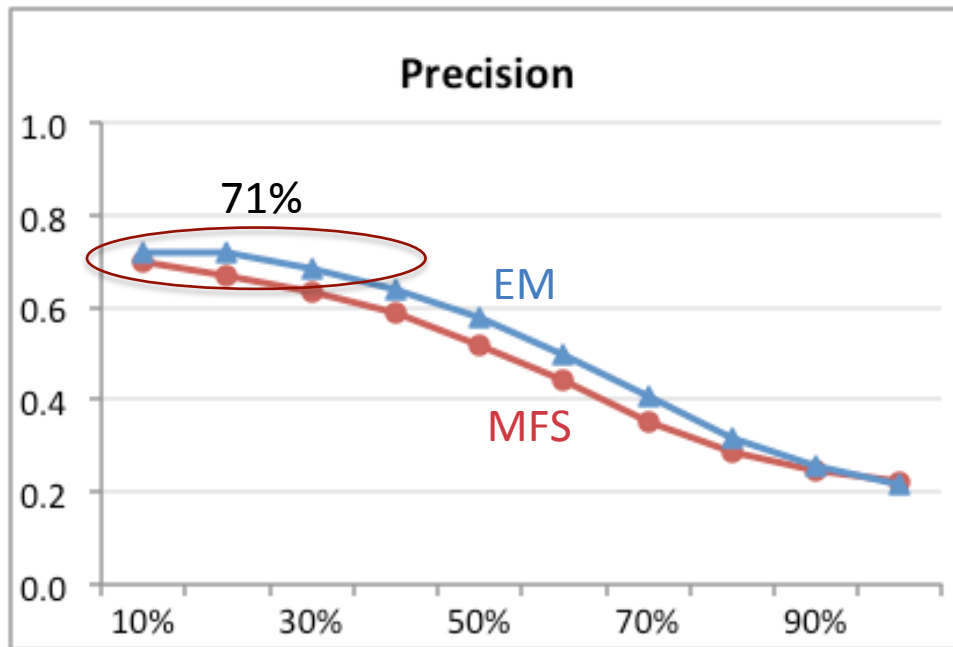
Overall performance of *FrameFinder*



Precision and Recall of *FrameFinder*

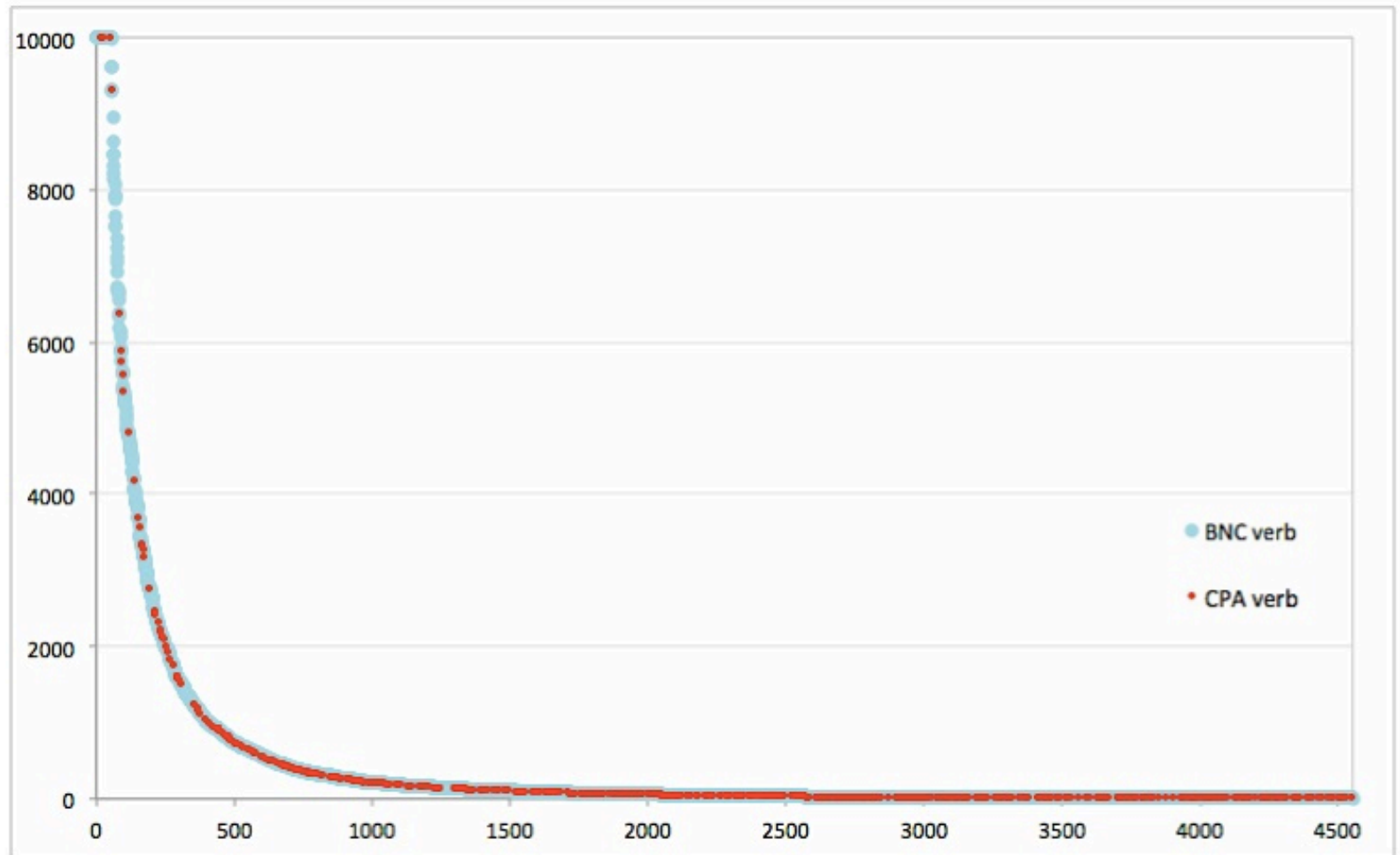


Our goal is to provide verb frames to facilitate language learning



We need high precision rate
in high frequency verbs

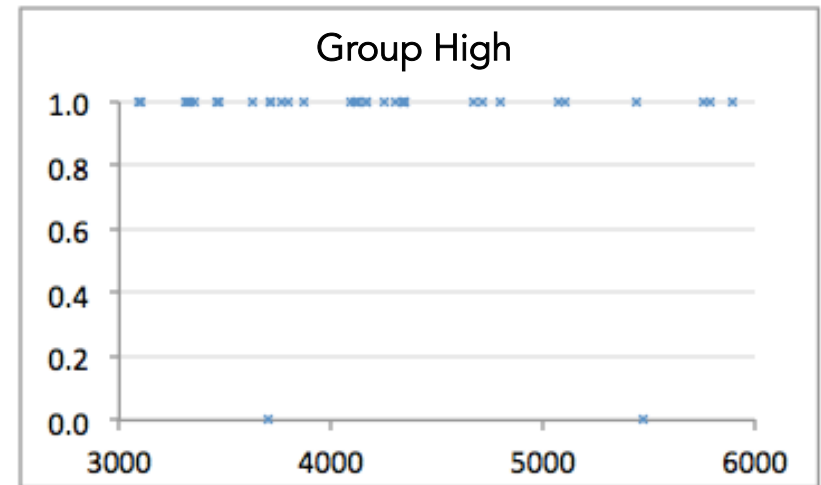
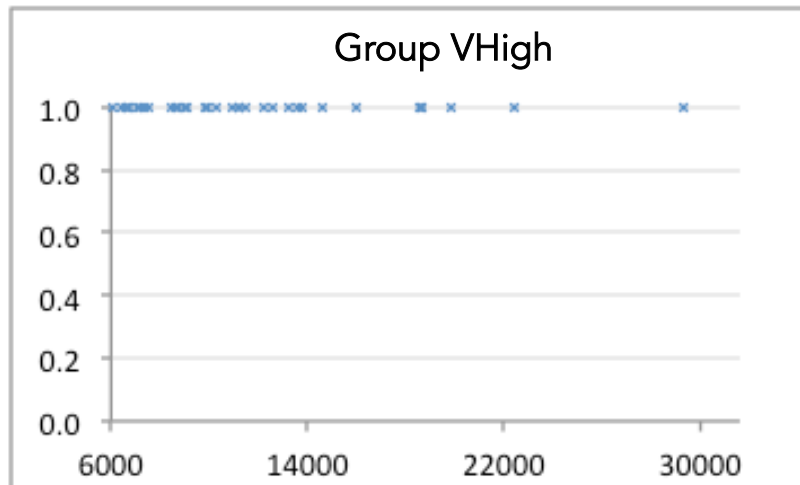
Number of occurrences of verbs in BNC and CPA



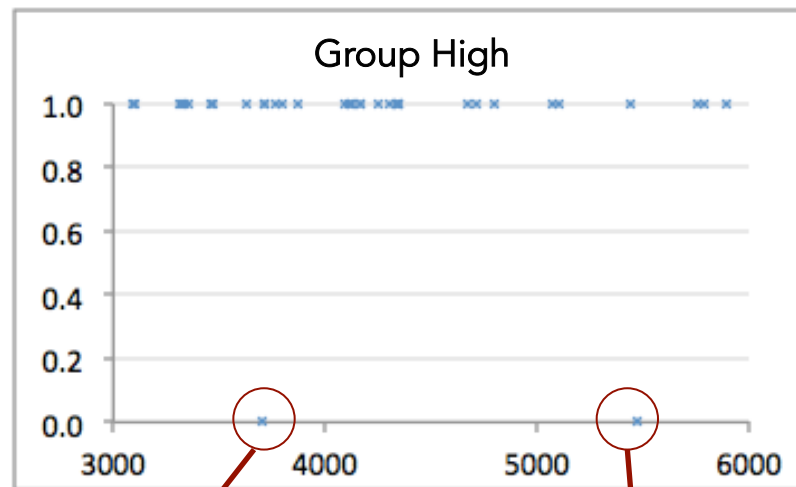
Five grouping criteria and the number of verbs in CPA

Group	Verb count criterion	# of verbs in CPA	Verb samples
VHigh	> 6000	26 → 3%	talk, say, tell
High	3000 - 6000	50 → 6%	propose, accuse, drink
Mid	500 - 3000	85 → 10%	abolish, dispose, pray
Low	150 - 500	77 → 9%	irritate, disregard, overflow
VLow	< 150	575 → 72%	petrify, abase, apostrophize

The precision of each verb in VHigh and High



Failure cases in Group High



recall

glance

PERSON recall

PERSON glance at PERSON

CPA internal labeling inconsistency

He glanced at his young colleague.



He glanced at his young colleague.

CPA (gold standard)

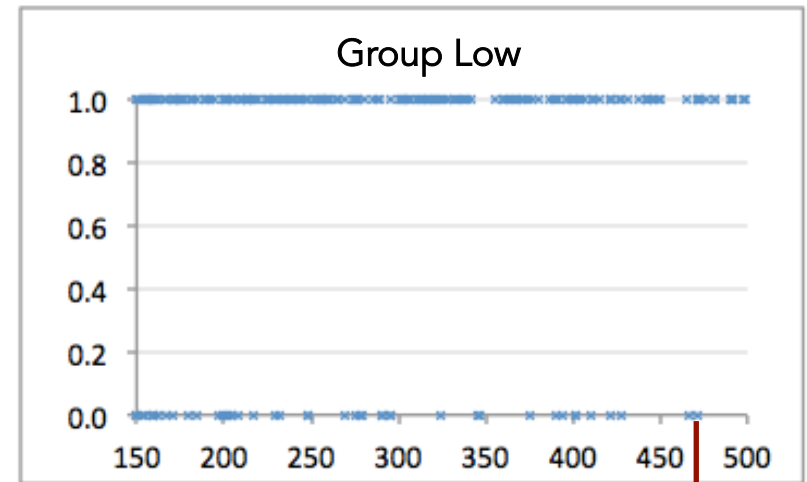
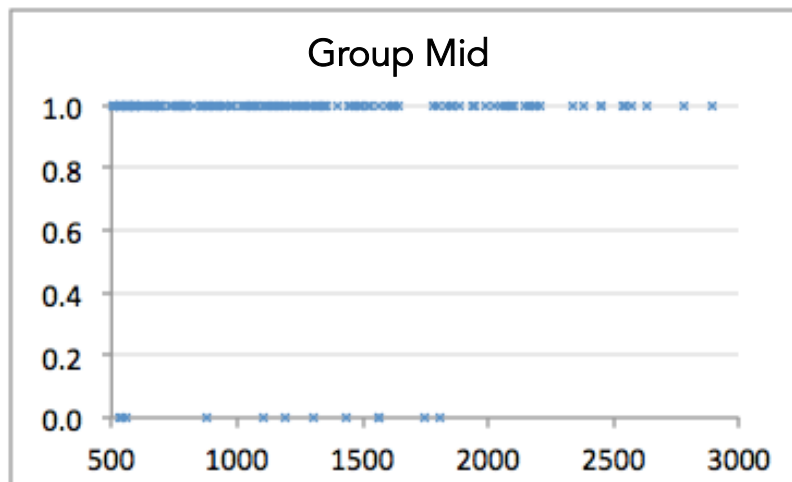
FrameFinder

PERSON glance [NO OBJ] Direction
Adv.

PERSPN glance at PERSON

PERSON yell [NO OBJ] at PERSON

The precision of each verb in **Mid** and **Low**



anger

CPA-WordNet mapping inconsistency

The announcement is bound to **anger** the fundamentalists.



The announcement is bound to **anger** the fundamentalists.

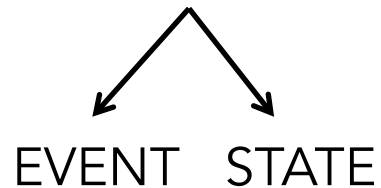


FrameFinder: COMMUNICATION **anger** PERSON

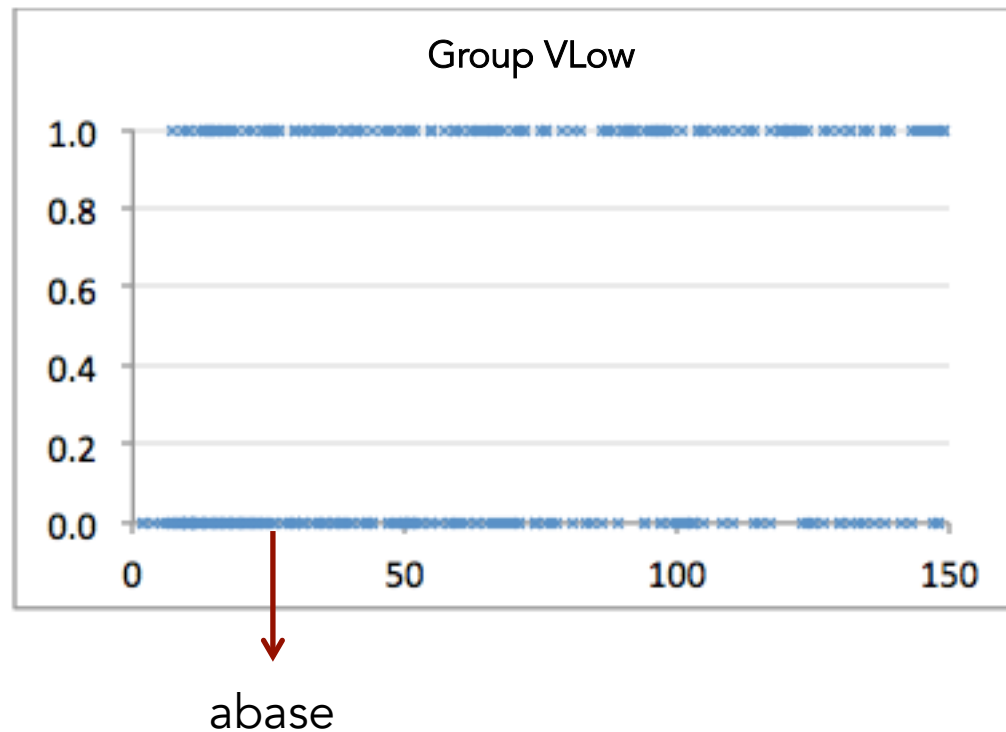
CPA-WordNet mapping inconsistency

FrameFinder: COMMUNICATION anger PERSON

CPA: [Event] anger PERSON



The precision of each verb in VLow

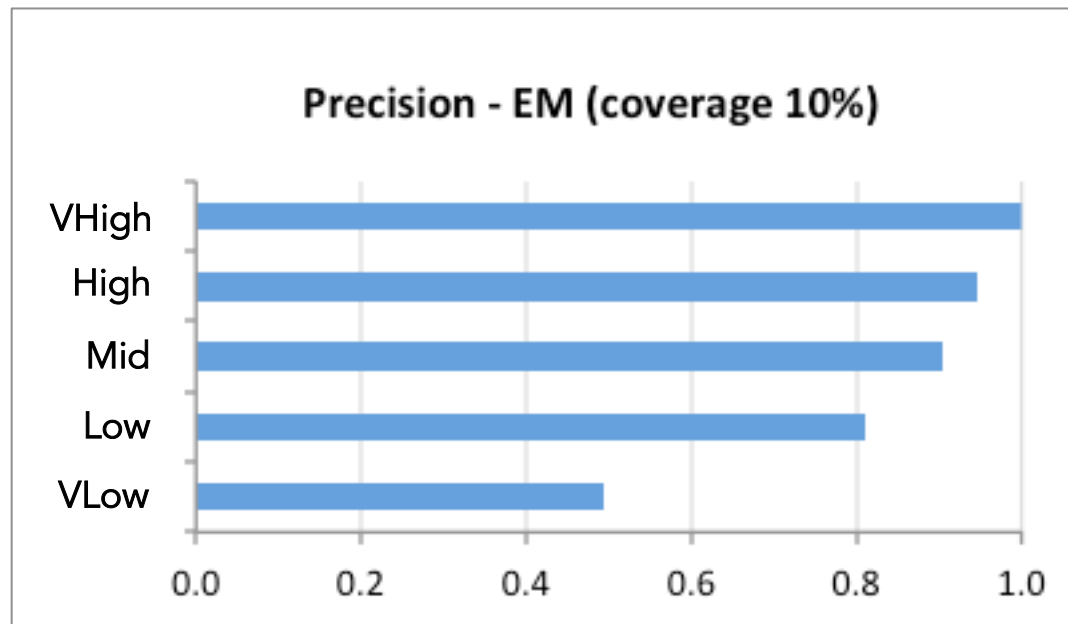


Optional adverbial

FrameFider: PERSON abase PERSON in STATE
PERSON abase PERSON before PERSON

CPA: PERSON abase PERSON

The precision of *FrameFinder* in each groups



Summary

- We have introduced a method combining corpus statistics and knowledge base to disambiguate and automatically acquire verb frames
- Evaluation results show that the method is able to extract verb frames with reasonable precision and recall of the most frequent usages for language learning

Future work

- Existing very large corpora can be exploited
 - UKWaC: 2 billion words
 - ClueWeb: 70 billion words
- More grammatical relations should be utilized
 - Passive voice
 - Clause
- The generated verb frames could be applied in other fields
 - Verb clustering
 - Grammatical error detection/correction
 - Learning or teaching the difference of verb usage

Thank you!