



Overview

- **Goal:** Evaluate semantic accuracy of non-native speaker (NNS) responses to picture description task (PDT)
 - Compare to gold standard (GS) of native speaker (NS) responses

Past Approach

1. Dependency parse NNS responses & GS
2. Use custom rules to extract and lemmatize *verb-subject-object* triple for each response
3. Attempt to match NNS triple to GS triples

Past Limitations

1. GS is small
2. NNS responses show more variation than NS responses
3. Matching exact triples is restrictive (no partial matching)
 - $kick(boy, ball) \neq kick(boy, football)$

Upshot: low coverage (50.8%)

Current Approach

Generalize methods by:

1. Representing responses as lists of dependencies
2. Scoring NNS response representation according to how closely it resembles GS representation
 - partial matching: **subj** **boy** **kick** + **obj** **ball** **kick** vs. **subj** **boy** **kick** + **obj** **football** **kick**

Data

Picture Description Task (PDT)

- 10 items (2 photos, 8 drawings) depicting transitive events
- PDT elicits natural productions but constrains form & content

Participants

- 39 NNSs, intermediate/advanced English; 390 sentences
 - Arabic, Chinese, Japanese, Korean, Spanish, Kurdish, Polish, Portugese
- 14 NSs; 140 sentences



Example Response (L1)

The man killing the beard. (Arabic)
A man is shutting a bird. (Chinese)
A man is shooting a bird. (English)
The man shouted the bird. (Spanish)

Generalizing the Methods

In the sections below, we explain the system parameter settings. The first two are closely related to generalizing the methods to overcome a limited GS, handle a wider range of sentence types (beyond transitives), and better reflect similarity to the GS.

- Response Representation: By moving to a “bag of dependencies” approach, we loosen the strict evaluation from *covered/not covered*; partial dependencies further loosen matching.
- Response Scoring: Averaging response term scores or calculating cosine distances allows for gradable rather than binary response scoring.

Response Representation

Responses are dependency parsed and treated as a list of *terms*, which are dependencies in one of the formats below (l, d, h = *label, dependent, head*; x = *placeholder*):

- **ldh**: subj_boy_kick
- **xdh**: x_boy_kick
- **lhx**: subj_x_kick
- **ldx**: subj_boy_x
- **xdx**: x_boy_x

Response Scoring

Scoring responses involves:

- Weighting terms (dependencies)
- Scoring responses: comparing weighted NNS terms with weighted GS terms

- **Frequency Average (FA):**
 - Weight: NNS terms assigned GS term frequencies
 - Response score: average of NNS term scores
- **Tf-idf Average (TA):**
 - Weight: NNS terms assigned tf-idf scores based on GS frequencies
 - Response score: average of NNS term scores
- **Frequency Cosine (FC):**
 - Weight: NNS & GS term frequencies are calculated
 - Response score: cosine distance between NNS & GS term scores
- **Tf-idf Cosine (TC):**
 - Weighting: NNS & GS tf-idf values are calculated
 - Response score: cosine distance between NNS & GS term scores

Reference Corpus

TA and TC require a reference corpus for deriving tf-idf scores. We experimented with two:

- **Brown Corpus (B)**
- **Wall Street Journal Corpus (W)**

NNS Source

We experiment with two forms of the NNS responses:

- **NNSO**: Original, uncorrected form
- **NNSLM**: Language Model autocorrected form

Future Directions: Clustering Items By Features

We used hierarchical clustering to explore for patterns among the items and parameters.

Set-up: cluster PDT items using features from response (e.g., type/token counts for terms) & features from system performance (i.e., average error score for parameter setting).

Goal: new PDT items could be placed into known clusters via response features & optimal parameter settings for that cluster could be applied automatically

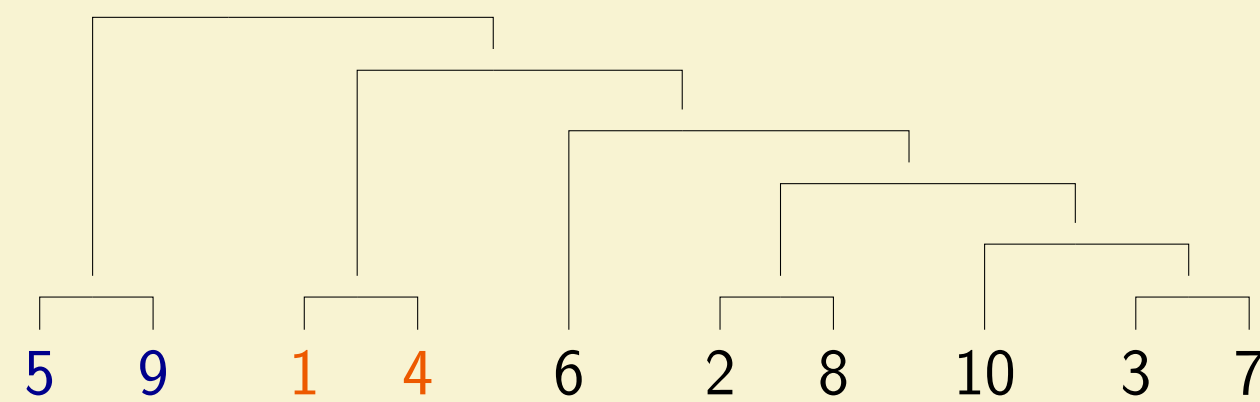


Figure 1: PDT items clustered by type and token counts of all NS, NNSO and NNSLM responses.

System Evaluation

1. Manually annotate responses; unacceptable responses are “errors”
2. Use each combination of parameters to produce a scored, ranked list of responses (Table 1)
 - Good parameter settings should rank good responses near GS and errors far from GS
3. Evaluate and rank parameter settings by (mean) average precision ((M)AP) (Table 3)
 - Also evaluate settings by non-normalized error score, which better illustrates differences in difficulty of PDT items (used in Figure 2)
4. Evaluate individual parameter values by MAP (Table 2)

<i>R</i>	<i>S</i>	Sentence	<i>E</i>	<i>V</i>
1	1.000	she is hurting.	1	1.5
	1.000	man mull bird	1	1.5
3	0.996	the man is hurting duck.	1	3.0
4	0.990	he is hurting the bird.	1	3.0
11	0.865	the man is trying to hurt a bird	1	11.0
12	0.856	a man hunted a bird.	0	0.0
17	0.775	the bird not shot dead.	1	17.0
18	0.706	he shot at the bird	0	0.0
19	0.669	a bird is shot by a un	1	19.0
20	0.646	the old man shooting the birds	0	0.0
37	0.086	the old man shot a bird.	0	0.0
38	0.084	a old man shot a bird.	0	0.0
39	0.058	a man shot a bird	0	0.0
Total Raw Score (not normalized)			17	169
Average Precision				0.75084

Table 1: Excerpt of rankings for Item 10 from the best system setting (TC.B.NNSLM.l dh) based on average precision scores. *R*: rank; *S*: sentence score; *E*: error; *V*: rank value.

Results

	Approach	Term Form	Ref. Corp. (TA/TC)	NNS Source
0.51577	TC	x dh	0.51810 Brown	0.51534 NNSLM
0.50780	FC	l dh	0.51677 WSJ	0.50798 NNSO
0.50755	TA	l x h	0.51350	
0.49464	FA	x d x	0.49901	
		l d x	0.49352	

Table 2: Approaches and parameters ranked by mean average precision for all 10 PDT items.

- Best approach: **TC**
 - **TC** > **FC**, **TA** > **FA**:
 - tf-idf weighting > frequency weighting
 - **TC&FC** > **TA&FA**:
 - cosine distance > weight averaging
- Term form: **x dh, l dh, l x h** > **x d x, l d x**
 - Importance of heads (**h**): with short transitive responses, verbs are salient (subj/obj head)
- Reference corpus: **Brown** > **WSJ**
 - Content & style of responses more like **Brown**
- NNS source: **NNSLM** > **NNSO**
 - More errors in NNSLM forms, inflating MAP values: use non-normalized scores? (see paper)

Rank	MAP	Settings
1	0.5534	TC.B.NNSLM.l x h
2	0.5445	TA.B.NNSLM.l x h
3	0.5435	TC.W.NNSLM.l x h
4	0.5422	TC.B.NNSLM.x d h
5	0.5368	TC.B.NNSLM.l d h
56	0.4816	TA.B.NNSO.x d x
57	0.4796	FA.na.NNSLM.l d x
58	0.4769	FC.na.NNSO.l x h
59	0.4721	TA.W.NNSO.x d x
60	0.4530	FA.na.NNSO.l x h

Table 3: Based on Mean Average Precision, the five best and five worst settings across all 10 PDT items.

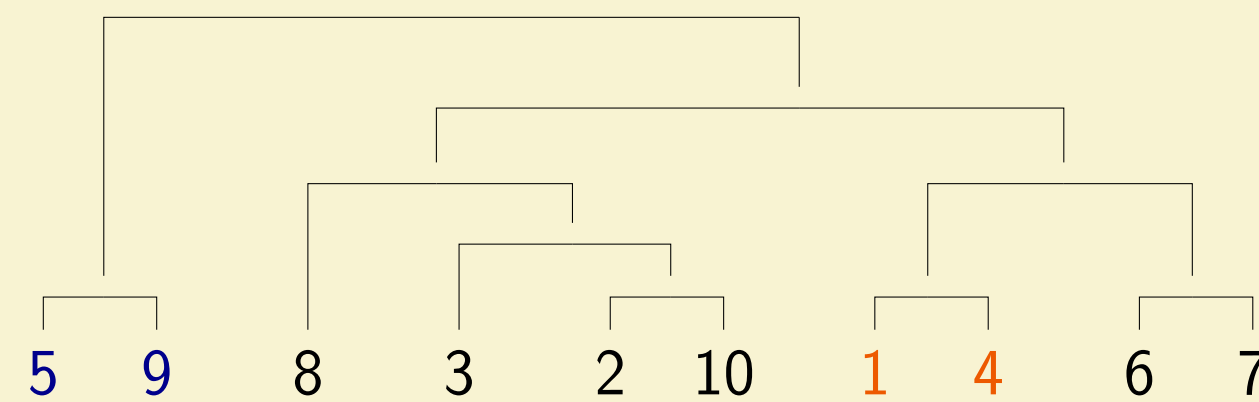


Figure 2: PDT items clustered by parameter performance.