Annotating Picture Description Task Responses for Content Analysis



Levi King & Markus Dickinson Indiana University

Overview & Background

- Introducing the Semantic Analysis of Image-based Learner Sentences (SAILS)

 Corpus
- ▶ 13,533 picture description task (PDT) responses from native speakers (NS) and non-native speakers (NNS), each annotated for five binary features
- ▶ **Goal:** Evaluate content of NNS sentences
- ► Compare to gold standard (GS) of NS sentences
- ▶ Needs: Adequate data, appropriately constrained
- ▶ Large set of PDT responses from NS and NNS participants
- ▶ Varied task prompts and participant demographics allowing for study of variability
- Annotation allowing for content analysis

Picture Description Task

- ▶ PDT elicits natural productions but constrains form & content
- ▶ 60 **items**: 30 images *x* 2 prompts
- ▶ 30 image
- simple vector graphics
- ▶ 10 transitive, 10 intransitive, 10 ditransitive actions
- ▶ 2 prompt versions:
- ▶ targeted: What is <the subject> doing?
- untargeted: What is happening?

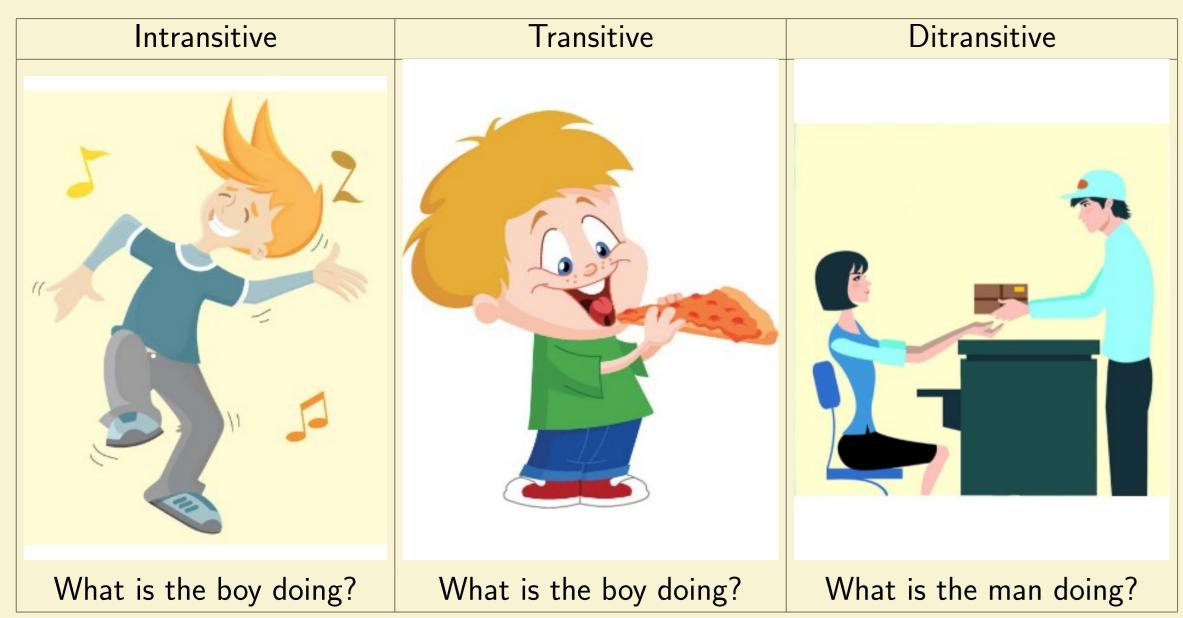


Table 1: Example PDT images with their **targeted** questions. In the **untargeted** form, the question for each is *What is happening?*

- ▶ PDT Instructions:
- ▶ Focus on the main action
- Respond in a complete sentence
- Task administered as online survey (SurveyMonkey.com)
- Multiple versions
- ► Most participants completed 30 items
- ► Roughly equal number of targeted & untargeted responses collected per image
- NNSs provide one response per item
- NSs asked to provide two non-identical responses per item
- ► Intended to increase variety of NS responses for a more robust GS

Participants

- 499 total participants
- ▶ 141 NNSs
- ▶ Recruited from intermediate & advanced English as a Second Language writing courses at Indiana University
- L1s: 125 Chinese (90%), 4 Korean, 3 Burmese, 2 Hindi; 1 each: Arabic, Indonesian, German, Gujarati, Spanish, Thai, Vietnamese
- ▶ 358 NSs
- ▶ 29 Familiar Native Speakers (FNSs)
- ► Relatives or friends of researchers
- ▶ Responses are assumed to be high quality
- ▶ 329 Crowdsourced Native Speakers (CNSs)
- ► Responses purchased via SurveyMonkey
- Responses parenased via surveymentage
 Responses are assumed to be lower quality

Responses

The SAILS Corpus contains a total of 13,533 PDT responses.

	Response Counts								
Group	First	Second	Total						
NNS	4290	0	4290						
NS (all)	4634	4609	9243						
FNS	642	641	1283						
CNS	3992	3968	7960						
Total	8924	4609	13,533						

Table 2: First and second response counts for the SAILS Corpus participant groups. Familiar (FNS) and crowdsourced (CNS) are subgroups of NS. NNS participants are not asked to provide a second response.

In order to examine the level of variation among responses, type to token ratios (TTRs) were calculated on the response level. Capitalization and final punctuation were ignored. We can see that variation increases with item complexity (intransitives < transitives < ditransitives) and that untargeted responses vary more than targeted responses.

	Targ	eted	Untar	geted			
Set	NS	NNS	NS	NNS			
Intransitives	0.628	0.381	0.782	0.492			
Transitives	0.752	0.655	0.859	0.779			
Ditransitives	0.835	0.817	0.942	0.936			

Table 3: NS and NNS type-to-token ratios (TTR) for complete responses (not words), for the full corpus.

TTRs were also calculated to compare the variability among NSs' first responses versus their second responses. As the TTRs for second responses are considerably higher than those for first responses, asking for two non-identical responses appears to effectively increase the variety of NS responses available for use in a GS.

	Targ	eted	Untargeted				
Set	R1	R2	R1	R2			
Intransitives	0.343	0.819	0.549	0.939			
Transitives	0.509	0.895	0.682	0.926			
Ditransitives	0.641	0.948	0.864	0.955			

Table 4: TTRs for complete responses, separated by first responses (R1) and second responses (R2). The ratios here are calculated from all NS responses; NNS responses are not included.

Annotation

Two annotators:

- ▶ NSs (US English), both with language teaching experience (child & adult learners).
- ▶ Annotator 1 (A1) annotated the complete corpus, Annotator 2 (A2) annotated development set & test set, each containing 1 intransitive, 1 transitive, 1 ditransitive.

Initial scheme (accurate + native-like > accurate + not native-like > not accurate) proved problematic; evolved to five binary features related to accuracy & native-likeness.

- 1. Core Event: Does the response capture the core event depicted in the image?
- 2. **Verifiability**: Does the response contain only information that is true and verifiable based on the image? Inferences allowed only when necessary; e.g., familial relationships between persons in image.
- 3. **Answerhood**: Does the response make a clear attempt to answer the question? This generally requires a progressive verb. For targeted items, the subject of the question or an appropriate pronoun must be used as the subject of the response.
- 4. **Interpretability**: Does the response evoke a clear mental image (even if different from the actual item image)? Any required verb arguments must be present and unambiguous.
- 5. **Grammaticality**: Is the response free from errors of spelling and grammar?

What is the boy doing?	С	V	Α		G	What is happening?	С	V	Α		G
eating food.	0	1	1	1	1	Child is eating pizza.	1	1	1	1	0
eatting.	0	1	1	1	0	Tommy is eating pizza.	1	0	1	1	1
The child is about to eat pizza.	1	1	0	1	1	The boy's eating his favorite food.	0	0	1	0	1
He may get fat eating pizza.	1	0	0	1	1	Pizza is this boy's favorite food.	0	0	0	0	1

Table 5: Targeted and untargeted sample responses shown with adjudicated annotations for the development set transitive item (see Table 1).

PLACEHOLDER One

PLACEHOLDER Two

 ${\tt BEA13,\,5\,\,June\,\,2018;\,\,New\,\,Orleans,\,\,LA}$