

Annotating Picture Description Task Responses for Content Analysis

Levi King & Markus Dickinson Indiana University

Overview

Semantic Analysis of Image-based Learner Sentences (SAILS) Corpus

- ► 13,533 picture description task (PDT) responses
- Both native (NS) & non-native speakers (NNS)
- Annotated for five binary features
- **Goal:** Evaluate content of NNS sentences
- Compare to gold standard (GS) of NS sentences
- ► **Need:** Adequate data, appropriately constrained
- Large set of PDT responses
- Varied task prompts & participant demographics
- Annotation for content analysis

Picture Description Task

PDT elicits natural productions but constrains form & content

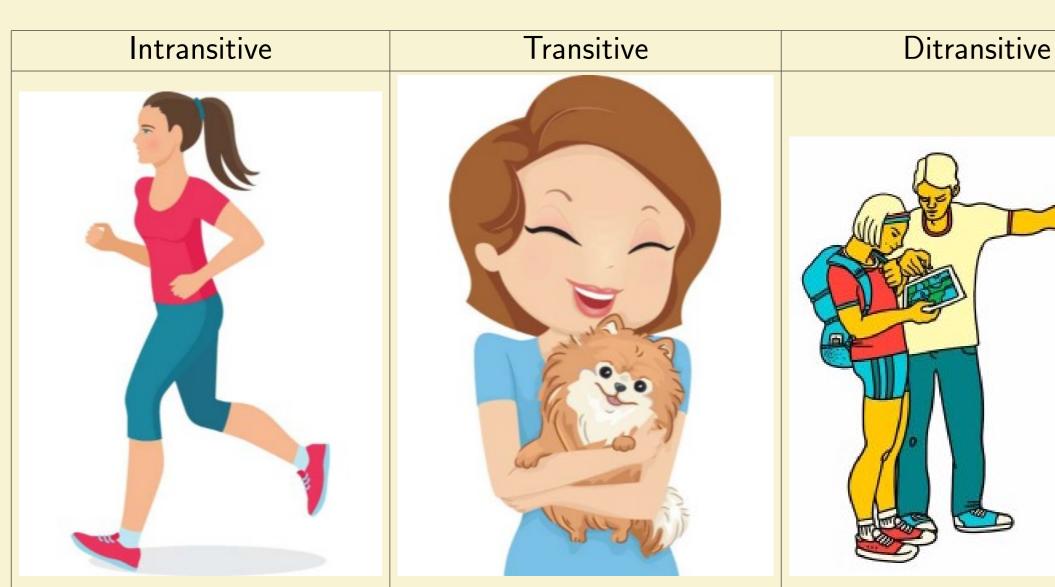
▶ 60 **items**: 30 images x 2 prompts

2 prompts

- Simple vector graphics
- ▶ 10 intransitive, 10 trans, 10 ditrans

30 images

- ► **Targeted**: What is < the subject
- Untargeted: What is happening?



What is the woman doing? What is the woman doing? What is the man
 Table 1: Example PDT images with their targeted questions.

Administered as online survey (SurveyMonkey.com)

PDT Instructions

- Focus on the main action
- Respond in a complete sentence

Multiple versions

- Most participants completed 30 items
- Roughly equal number of targeted & untargeted responses
- NNSs provide one response per item
- ► NSs provide two non-identical responses per item (more robust GS)

Participants

499 total participants

- ▶ 141 NNSs: students in intermediate & advanced ESL writing courses at IU
- L1s: 125 Chinese (90%), 4 Korean, 3 Burmese, 2 Hindi; 1 each: Arabic, Indon German, Gujarati, Spanish, Thai, Vietnamese

▶ 358 NSs

- 29 Familiar Native Speakers (FNSs)
- Relatives or friends of researchers (assumedly higher quality)
- ► 329 Crowdsourced Native Speakers (CNSs)
- Responses purchased via SurveyMonkey (assumedly lower quality)

	Responses							
us			Response Coun	nts				
		Group	Response First Second					
		NNS NS (all)	4290 (4634 4609	0 4290 9 9243				
		FNS	642 642	1 1283				
		CNS Total	3992 3968 8924 4609					
	Table 2: Firs	t & second resp	onse counts for S	AILS Corpus part	icipant groups			
		Туре-	-Token Ratios	(TTRs)				
			Targeted	Untargeted				
		Set Intransitives	NS NNS 0.628 0.381	NS NNS 0.782 0.492				
		Transitives	0.752 0.655	0.859 0.779				
	Table 3:	Ditransitives TTRs for <i>com</i>	0.835 0.817	0.942 0.936 not words), for fu	II corpus			
				,,,,				
t> doing? ~2	 Capitalization & fina Variation increases w 		gnored					
g?	 Item complexity (transitives $< ditr$	ransitives)				
/e	Less targeting (ta	argeted < untargeted	geted)					
	Type-Token Ratios (TTRs): first vs. second responses (NSs only)							
		Set	Targeted R1 R2	Untargeted R1 R2				
		Intransitives	0.343 0.819	0.549 0.939				
		Transitives Ditransitives	0.5090.8950.6410.948	0.682 0.926 0.864 0.955				
	Table 4: TTRs for	complete respo	nses, separated b	y first (R1) & sec	ond responses (R2)			
	► TTRs for R2s consid	derably higher th	nan for R1s					
doing?	\Rightarrow Asking for two re	sponses increase	es variety of langu	uage available for	use in GS			
	Annotation Scheme							
	Initial scheme: accur	rate + native-lik	ke > accurate + k	not native-like >	not accurate)			
	Final scheme: five bi				,			
	1. Core Event (C): [Does response c	apture the core e	vent depicted in in	mage?			
	2. Verifiability (V):	Does response c	ontain only true a	& verifiable info, l	based on image?			
	Inferences allowed	d only when nec	essary; e.g., fami	lial relationships o	f persons in image			
	3. Answerhood (A):			empt to answer th	e question?			
	Generally requires				*			
					t be response subject			
	 Interpretability (I Any required vert 							
	5. Grammaticality (ımar?			
onesian,		· ·						
	Annotators							
	Two annotators:							
	► NSs (US English), b	oth with langua	ge teaching expe	rience (child & ad	lult learners).			
	► Annotator 1 (A1): c	complete corpus						
	► Annotator 2 (A2): c	levelopment & t	test sets, each wit	th 1 intransitive, 1	trans, 1 ditrans			

- age
- ubject PDT)?

Annotation Results



 Table 5:
 Sample responses from development transitive item, with adjudicated annotations

Set	Total	A1Yes	A2Yes	AvgYes	Chance	Agree	Kappa
Intransitive	2155	0.863	0.855	0.859	0.758	0.978	0.910
Transitive	2155	0.780	0.774	0.777	0.653	0.949	0.853
Ditransitive	2155	0.812	0.786	0.799	0.678	0.924	0.764
Targeted	3390	0.829	0.818	0.824	0.709	0.949	0.823
Untargeted	3075	0.806	0.790	0.798	0.678	0.952	0.872
Core Event	1293	0.733	0.717	0.725	0.601	0.923	0.808
Verifiability	1293	0.845	0.817	0.831	0.719	0.968	0.884
Answerhood	1293	0.834	0.831	0.833	0.721	0.982	0.936
Interpretability	1293	0.818	0.787	0.802	0.682	0.919	0.744
Grammaticality	1293	0.861	0.872	0.866	0.768	0.960	0.827
	Intransitive Transitive Ditransitive Ditransitive Targeted Untargeted Untargeted Core Event Verifiability Answerhood Interpretability Grammaticality	Intransitive2155Transitive2155Ditransitive2155Ditransitive2155Targeted3390Untargeted3075Core Event1293Verifiability1293Answerhood1293Interpretability1293Grammaticality1293	Intransitive21550.863Transitive21550.780Ditransitive21550.812Targeted33900.829Untargeted30750.806Core Event12930.733Verifiability12930.845Answerhood12930.834Interpretability12930.818Grammaticality12930.861	Intransitive21550.8630.855Transitive21550.7800.774Ditransitive21550.8120.786Targeted33900.8290.818Untargeted30750.8060.790Core Event12930.7330.717Verifiability12930.8450.817Answerhood12930.8180.787Grammaticality12930.8610.872	Intransitive21550.8630.8550.859Transitive21550.7800.7740.777Ditransitive21550.8120.7860.799Targeted33900.8290.8180.824Untargeted30750.8060.7900.798Core Event12930.7330.7170.725Verifiability12930.8450.8170.831Answerhood12930.8180.7870.802Grammaticality12930.8610.8720.866	Intransitive21550.8630.8550.8590.758Transitive21550.7800.7740.7770.653Ditransitive21550.8120.7860.7990.678Targeted33900.8290.8180.8240.709Untargeted30750.8060.7900.7980.678Core Event12930.7330.7170.7250.601Verifiability12930.8340.8310.8330.721Interpretability12930.8180.7870.8020.682Grammaticality12930.8610.8720.8660.768	Intransitive21550.8630.8550.8590.7580.978Transitive21550.7800.7740.7770.6530.949Ditransitive21550.8120.7860.7990.6780.924Targeted33900.8290.8180.8240.7090.949Untargeted30750.8060.7900.7980.6780.952Core Event12930.7330.7170.7250.6010.923Verifiability12930.8450.8170.8310.7190.968Answerhood12930.8180.7870.8020.6820.919Grammaticality12930.8610.8720.8660.7680.960

Table 6: Agreement scores broken down by different properties of test set

- Cohen's kappa needed as measure of inter-annotator agreement

- Guidelines less complicated for untargeted items
- **Feature:** Answerhood has highest kappa, interpretability has lowest Matches annotator reporting of easiest & hardest features to annotate

Accessing the SAILS Corpus

Download the entire annotated SAILS Corpus, PDTs, & annotation guidelines at: https://github.com/sailscorpus/sails

SAILS corpus can be used for:

- Language testing & ICALL
- Possibilities for expansion from other researchers:
- New participants, items, approaches for processing

Annotation Examples

What is the boy doing? (Targeted)			Α		G	
eating pizza			1	1	1	
eating food.			1	1	1	
eatting.	0	1	1	1	0	
The child is eating pizza.	1	1	0	1	1	
He may get fat eating pizza.			0	1	1	
The boy is hungry.	0	1	0	0	1	
Pizza is this boy's favorite food.			0	0	1	
What is happening? (Untargeted)			A		G	
The kid's eating pizza		1	1	1	1	
Child is eating pizza.		1	1	1	0	
Tommy is eating pizza.		0	1	1	1	
The boy's eating his favorite food.		0	1	0	1	
A youngster anticipates the taste of pizza		1	0	1	1	
Pepperoni pizza makes the boy smile		0	0	1	1	
He sure is happy.		1	0	1	1	

Inter-Annotator Agreement

Observations from Table 6

Average yes rates (AvgYes) show all features skew toward yes annotations

Cohen's kappas well above conventional 0.67 threshold for meaningful agreement

 \Rightarrow Annotation scheme can be implemented reliably by following guidelines

► Verb Type: Agreement decreases with item complexity (intransitive > trans > ditrans)

Prompt: Agreement slightly higher for untargeted than targeted items

Question answering, dialog systems, pragmatic modeling, visual references