

# Semantic Analysis of Image-Based Learner Sentences

Levi King  
Indiana University

August 6, 2021

# Motivation

Most intelligent computer-assisted language learning (ICALL) applications (*Rosetta Stone*, *Duolingo*, etc.) rely on outdated, ineffective methods:

- ▶ rote memorization & grammatical error detection; menu-based vs. free input;
- ▶ “*engineering first*”: no second language acquisition, pedagogy;

SLA research → communicative & task-based learning

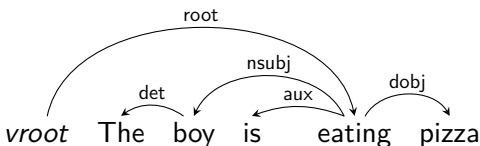
*How can we bridge this gap?*

- ▶ My vision: open source app; transparent; pipeline of existing tools;
- ▶ teachers create new games/stories by adding visual prompts and crowdsourcing native speaker (NS) responses;
- ▶ use NS model to evaluate non-native speaker (NNS) responses

# Research Questions

- RQ1. Are the picture description task (PDT) responses of L2 English learners sufficiently similar to those of NSs to allow automatic evaluation based on a collection of NS responses?
- RQ2. For PDT responses, what are appropriate representations for the purpose of providing meaning-oriented feedback or evaluation?
- RQ3. What kinds of NLP tools are appropriate here?
- RQ4. How do “bag-of-words” and “bag-of-dependencies” approaches compare in terms of performance?
- RQ5. Can the accuracy of the system be improved with information from semantic tools (e.g., BERT)?
- RQ6. What is the annotation scheme for this task and can the system perform within the range of human performance?

## Step 1: Dependency parse:



### Get dependencies:

root(eating, *vroot*)

det(the, boy)

nsubj(boy, eating)

aux(is, eating)

dobj(pizza, eating)

### Step 2: Lemmatize:

→ root(**eat**, *vroot*)

→ det(the, boy)

→ nsubj(boy, **eat**)

→ aux(**be**, **eat**)

→ dobj(pizza, **eat**)

# System

## Step 3: tf-idf (term frequency-inverse document frequency)

NS model: [He is eating pizza. The boy is eating pizza.]

NNS 1: He is eating food.      NNS 2: He is eating pizza.

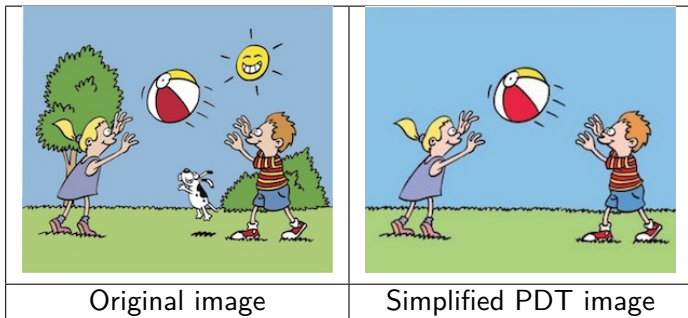
NS $\cup$ NNSs	NS model		NNS 1		NNS 2	
	tf	tf-idf	tf	tf-idf	tf	tf-idf
aux(be,eat)	2	.04	1	.02	1	.02
det(the,boy)	1	.04	-	0	-	0
dobj(food,eat)	-	0	1	.06	-	0
dobj(pizza,eat)	2	.16	-	0	1	.08
nsubj(boy,eat)	1	.08	-	0	-	0
nsubj(he,eat)	1	.04	1	.02	1	.02
root(eat,vroot)	2	.02	1	.01	1	.01

Response scores: cosine(NS model tf-idf vector, NNS tf-idf vector)

NNS 1: 0.139;    NNS 2: 0.886     $\rightarrow$     *NNS 2 is closest to the model.*

# Data collection



PDT with very simple images only:



Intended to focus participants' attention on the main action

# Data collection




Two PDT prompt versions:

Targeted	Untargeted
	
<i>What is <b>the baby</b> doing?</i>	<i>What is happening?</i>

Intended for exploring the specificity needed for my approach

# Data collection

3 verb types:


10 <b>intransitive</b> items	10 <b>transitive</b> items	10 <b>ditransitive</b> items
		
What is the girl doing?	What is the boy doing?	What is the girl doing?

Intended for exploring whether my approach can generalize to a range of sentence types



# Data collection

The pilot study *rake* problem; 100% of NS used the verb *rake*:

	NNS Responses
	The gardener is <i>cleaning</i> the street.
	a man <i>removing</i> the tree leafs.
	The man is <i>sweeping</i> the floor.
	A man is <i>gathering</i> lots of leafs.

- ▶ NNS responses without *rake* are penalized;
- ▶ I address this by asking NSs for two non-identical responses.

# Main study: Data collection

499 participants, 13,533 responses:

- ▶ 141 NNSs (ELIP at IU), 4,290 responses;
  - ▶ 125 Mandarin, 4 Korean, 3 Burmese, 2 Hindi; 1 each: Arabic, Indonesian, German, Gujarati, Spanish, Thai, Vietnamese;
- ▶ 358 NSs, 9,243 responses:
  - ▶ 329 crowdsourced, purchased via SurveyMonkey;
    - ▶ 7,960 responses;
  - ▶ 29 familiar, unpaid colleagues;
    - ▶ 1,283 responses;


# Annotation features

5 binary features:

- ▶ CORE EVENT: Does response capture main action?
- ▶ ANSWERHOOD: Does response directly answer prompt?
- ▶ GRAMMATICALITY: Is response free from grammar problems?
- ▶ INTERPRETABILITY: Does response evoke a clear mental image?
- ▶ VERIFIABILITY: Is all response info supported by image?

# Annotation features

Core event, **A**nswerhood, **G**rammaticality, **I**nterpretability, **V**erifiability

					
<i>What is the boy doing?</i>	C	A	G	I	V
He is eating food.	0	1	1	1	1
he eating pizza.	1	1	0	1	1
The boy is smiling pizza.	0	1	0	0	0
He may get fat eating.	0	0	1	1	0

Inter-rater reliability (Cohen's kappa): 0.744 (**I**) – 0.936 (**A**)

# Evaluating performance

To evaluate system performance, I need **benchmark rankings** for the NNS test set.

- ▶ Mean average precision (**MAP**) to see how system rankings predict **individual features**;
  - ▶ Also need MAP from benchmark rankings (*upper bound*);
- ▶ **Spearman** rank correlation: Compare system rankings with benchmark rankings to see how system predicts overall quality;

How do we get benchmark rankings from 5 binary annotations?

- ▶ Determine feature weights and apply to annotations to obtain benchmark holistic scores and then rankings.
- ▶ Annotators performed a preference task for pairs of responses.
- ▶ Feature weights were derived according to how frequently each feature is “yes” among preferred responses.

# Benchmark rankings

*Weighted annotation score (WAS); weighted annotation ranking (WAR)*

<i>What is happening?</i>	C	A	G	I	V	WAS	WAR
The boy is eating pizza	.365	.093	.055	.224	.263	10	1
Child is eating pizza	.365	.093	0	.224	.263	.945	2
Tommy is eating pizza	.365	.093	.055	.224	0	.737	3
The boy's eating his favorite food	0	.093	.055	0	0	.513	4
Pizza is this boy's favorite food	0	0	.055	0	0	.055	5

Agreement for two annotators on a sample of 300 pairs:

Chance Agree	Observed Agree	Cohen's Kappa
.621	.883 (265/300)	.692

# SBERT for comparison

I also use SBERT for comparing my system's performance.

- ▶ State-of-the-art sentence embedding for semantic textual similarity.
- ▶ Replaces dependency parser + lemmatizer + tf-idf cosine pipeline.
- ▶ Provides distance between NNS response and NS model; rankable.
- ▶ Not explainable; Internal representations are not suitable for informing a feedback module.

# System configuration

Optimizing means finding the best system settings:

- ▶ **Transitivity:** intransitive, transitive, ditransitive;
- ▶ **Targeting:** targeted, untargeted;
- ▶ **Familiarity:** familiar, crowdsourced;
- ▶ **Primacy:**
  - ▶ primary: NS model contains only 1st responses;
  - ▶ mixed: NS model: 1st & 2nd responses (50-50);
- ▶ **Term Representation:**
  - ▶ ldh: label-dependent-head; i.e., labeled dependencies;
  - ▶ xdh: dependent-head; i.e., unlabeled dependencies;
  - ▶ xdx: dependent only; cf. *bag of words*;
  - ▶ Does not apply to SBERT (operates on plain text);

A **system configuration** combines one setting from each.



# Sampling data

## **NNS test sets:**

- ▶ All experiments rank the same randomly sampled NNS test sets;
- ▶ 70 targeted, 70 untargeted per PDT item (max available for NNS data);

## **NS models:**

- ▶ 14-response models (max available for familiar data);
- ▶ 50-response models (max available for crowdsourced data);

# Sampling data: Complexity

Standardized type-to-token ratio (STTR)  
for response samples. Tokens here are  
*dependencies*.

	n14		n50	n70
	Fam	Crd	Crd	NNS
Intrans	.558	.525	.535	.391
Trans	.569	.580	.581	.517
Ditrans	.598	.640	.637	.606
Target	.545	.535	.545	.481
Untarg	.610	.633	.621	.528
Primary	N/A	.517	.523	.505
Mixed	.576	.652	.645	N/A
xdx	.364	.424	.421	.364
x dh	.658	.661	.660	.572
l dh	.665	.664	.671	.578
Total	.576	.583	.584	.505

Complexity often correlates with  
parameter settings in terms of  
system performance.

Within each parameter block,  
complexity increases as we move  
down the rows. E.g.:

$\text{Intrans} < \text{Trans} < \text{Ditrans}$

In some settings (e.g.,  
Intrans), Crowd complexity is  
closer to NNS than is Familiar;  
other settings vice versa (e.g.,  
Ditrans).

## Predicting features: CORE EVENT MAP

	Crowd NS model = 14					Crowd NS model = 50				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	<b>.859</b>	.856	.854	.865	.835	<b>.855</b>	.854	.852	.865	.831
Tran	<b>.737</b>	.735	.728	.742	.703	<b>.736</b>	.733	.725	.742	.701
Ditr	<b>.665</b>	.661	.664	.660	.634	.657	.656	<b>.661</b>	.660	.629
Targ	<b>.739</b>	.738	.732	.735	.708	<b>.737</b>	.735	.729	.735	.704
Untg	<b>.768</b>	.763	.765	.777	.740	.762	.759	<b>.763</b>	.777	.736
Prim	<b>.754</b>	.752	.747	.756	.723	<b>.750</b>	.748	.745	.756	.719
Mix	<b>.753</b>	.749	.750	.756	.725	<b>.749</b>	.746	.746	.756	.721
Total	<b>.735</b>	.751	.748	.756	.724	<b>.750</b>	.747	.746	.756	.720

- ▶ In all cases, ldh + 14NS is best (slightly);
- ▶ xdx becomes more competitive for larger model (50NS);
  - ▶ ditrans, untarg: *most complex*—i.e., highest STTRs;
  - ▶ In general: ldh STTR > x dh STTR > x dx STTR

## Predicting features: CORE EVENT MAP

	Familiar NS model = 14					Crowd NS model = 14				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.859	.859	<b>.865</b>	.865	.838	<b>.857</b>	.852	.848	.865	.833
Tran	<b>.740</b>	.737	.726	.742	.703	<b>.738</b>	.735	.728	.742	.702
Ditr	.651	.648	<b>.660</b>	.660	.625	.663	.659	<b>.673</b>	.660	.641
Targ	<b>.733</b>	.732	.732	.735	.707	<b>.739</b>	.736	.733	.735	.709
Untg	.767	.764	<b>.769</b>	.777	.737	.767	.761	<b>.767</b>	.777	.742
Total	.750	.748	<b>.751</b>	.756	.722	<b>.753</b>	.749	.750	.756	.725

- ▶ \*mixed only (due to sparse familiar data);
- ▶ Totals: crowdsourced outperforms familiar (slightly);
- ▶ crowdsourced works best with ldh;
- ▶ familiar works best with xdx;

## Predicting features: MAP Results

For all 5 features, my system outperforms SBERT.

ANSWERHOOD, in *all* cases:

- ▶ `xdx > xdh > ldh`;
- ▶ Model size makes no difference;
- ▶ `familiar > crowdsourced`;

GRAMMATICALITY, in *most* cases:

- ▶ `xdx > xdh > ldh`;
- ▶ `familiar 14NS > crowd 14NS > crowd 50NS`;

Predicting ANSWERHOOD or GRAMMATICALITY is relatively simple; requires only small model and bag-of-words representation.

# Predicting features: MAP Results

## INTERPRETABILITY:

- ▶ 14NS crowd > 14NS familiar > 50NS crowd;

## VERIFIABILITY:

- ▶ 14NS crowd > 50NS crowd > 14NS familiar;
- ▶ Model size effect is most pronounced with untargeted & mixed;
  - ▶ Unconstrained settings; larger models have more noise;

## For both INTERPRETABILITY & VERIFIABILITY:

- ▶ intransitives & ditransitives work best with xdx;
- ▶ transitives work best with ldh;
  - ▶ Why? Transitive responses are relatively homogenous;  
Annotators relatively strict;

# Predicting quality

Holistic quality experiments use one set of 360 Spearman correlations:

targeting (2)  $\times$  primacy (2)  $\times$  term rep (3)  $\times$  items (30) = 360.  
(Familiar vs. Crowd handled separately due to sparse data.)

Each experiment focuses on one variable, e.g., targeting:

Divide 360 into 180 targeted scores and 180 untargeted scores;  
compare mean, median, etc.

SBERT uses plain text (no term rep), thus only 120 total.

(SBERT always wins over system.)

# Predicting quality: Transitivity

Spearman rank correlations: System vs. WAR (benchmark)

		intrans		trans		ditrans	
		Sys	SBERT	Sys	SBERT	Sys	SBERT
	count	120	40	120	40	120	40
14NS	mean	<b>.439</b>	.497	.314	.563	.267	.400
	median	<b>.416</b>	.479	.304	.555	.276	.444
50NS	mean	<b>.423</b>	.516	.345	.566	.278	.446
	median	<b>.426</b>	.517	.331	.561	.286	.471

- ▶ SBERT, regardless of model size: trans > intrans > ditrans;
- ▶ System, regardless of model size: intrans > trans > ditrans;
- ▶ More complex items (TTR) work best with larger models;
  - ▶ trans & ditrans: 50NS model is best;
  - ▶ intrans: 14NS gives best mean, 50NS gives best median;



# Predicting quality: Results

## Targeting:

- ▶ targeted > untargeted
- ▶ 50NS models > 14NS models
  - ▶ Model size effect is most pronounced for targeted

## Familiarity (14NS models only):

- ▶ System: No discernible difference for familiar vs crowdsourced
- ▶ SBERT: familiar > crowdsourced
  - ▶ NNS STTR < familiar STTR < crowdsourced STTR

# Predicting quality: Results

## Privacy:

- ▶ System: 14NS: primary  $<$  mixed (slight difference)
- ▶ System: 50NS: primary  $\approx$  mixed
- ▶ SBERT: primary  $>$  mixed
- ▶ System & SBERT: 50NS  $>$  14NS
  - ▶ System: model size effect is greatest for primary

## Term representation:

- ▶ SBERT: 50NS  $>$  14NS
- ▶ System: for 1dh & xdh: 50NS  $>$  14NS;
  - ▶ Model size effect is greater for 1dh
- ▶ System: for xdx: NS14  $>$  NS50 (very slight)

# Summary

- ▶ Collected 13,533 PDT responses from 499 participants;
- ▶ Annotated for 5 features, focused on content;
- ▶ Established feature weights and benchmark rankings;
  - ▶ Features and weights are reliable;
- ▶ Developed explainable semantic textual similarity system based on dependencies and tf-idf;
- ▶ Uncovered some exploitable patterns for predicting features and holistic quality;

## Future work

With more data, I would:

- ▶ Explore results for broader range of L1s;
- ▶ Compare results across L2 English proficiency levels;
- ▶ Further map relationship between complexity and optimal settings;
- ▶ For a given PDT item, try clustering responses into multiple models;
  - ▶ Route NNS responses to most appropriate model;

## Backup slides

(The following slides are all backup for Q&A.)


# Research Questions (full)

- RQ1. Are the responses of L2 English learners sufficiently similar to those of NSs to allow automatic evaluation based on a collection of NS responses? In other words, do learners demonstrate significant overlap with native-like usage in a picture description task (PDT) setting?
- RQ2. In the constrained communicative environment of a PDT, what are appropriate response and model representations for the purpose of providing meaning-oriented feedback or evaluation? In other words, which linguistic components are crucial and which are superfluous?
- RQ3. What kinds of existing NLP tools and language resources can be integrated to form a content analysis system for open response language learning tasks?

# Research Questions (full)

- RQ4. How do “bag-of-words” and “bag-of-dependencies” approaches compare in terms of performance? Is a bag-of-words approach alone adequate for our needs?
- RQ5. Can the accuracy of the system be improved by the inclusion of semantic information from tools like semantic role labelers, WordNet, or word and sentence embeddings?
- RQ6. What is the annotation scheme for this task and can the system perform within the range of human performance? Relatedly, what does it mean for a response to be *appropriate* and how can this be captured with annotation?

## Pilot study: Data

	Response (L1)
	He is droning his wife pitcher. (Ar)
	The artist is drawing a pretty women. (Ch)
	The artist is painting a portrait of a lady. (En)
	The painter is painting a woman's paint. (Sp)

**Figure:** Example item from the pilot study showing responses from native speakers of Arabic (Ar), Chinese (Ch), English (En) and Spanish (Sp).

- ▶ 10 (transitive) PDT items  $\times$  53 participants = 530 responses;
  - ▶ 14 NSs (grad students), 39 NNSs (ESL students);
- ▶ Annotation: *Given the prompt, would the response be acceptable to most English speakers? Acceptable/unacceptable*
  - ▶ 1 annotator (me)



# Pilot study: Processing

First approach: **Rule-based** triple extraction and matching

Dependency parser  $\rightarrow$  lemmatizer  $\rightarrow V(S,O)$  extraction rules;

Compare NNS  $V(S,O)$  & NS  $V(S,O)$  list  $\rightarrow$  covered / not covered;

- ▶ Dependency-based
  - ▶ Captures aspects of form and meaning;
  - ▶ Subjects, objects, verbs clearly labeled;
- ▶  $V(S,O)$  extraction
  - ▶ Decision tree based on dependency indexing & labels, POS;
  - ▶ Custom for my transitive PDT, not generalizable, not robust;
  - ▶  $\approx 92\%$  accurate,  $\approx 8\%$  extraction errors;
- ▶ Overall accuracy: 58.9%
  - ▶ I.e., *Acceptable* covered, *unacceptable* not covered;

## Pilot study: Processing

Second approach: **Semantic similarity** scoring

Dependency parser  $\rightarrow$  lemmatizer  $\rightarrow$  term frequency-inverse document frequency (tf-idf; “term” = lemmatized dependency);

NNS response score = cosine distance of NS and NNS tf-idf scores;

- ▶ tf-idf: Score dependencies according to importance;
- ▶ Vectorize & Score
  - ▶ Get *sorted union set* of NS and NNS dependencies;
  - ▶ NNS vector: Replace deps with their **NNS** tf-idf scores;
  - ▶ NS vector: Replace deps with their **NS** tf-idf scores;
  - ▶ Response score = *cosine distance* for NNS & NS vectors;
- ▶ Rank by scores & calculate Mean Average Precision (MAP);
  - ▶ MAP *acceptable* responses:  $\approx 51\%$
- ▶ Process is more robust & generalizable;
- ▶ Dataset (especially NS models) and annotation are weak;

# System configuration

All parameters or variables and their settings:

Transitivity	Targeting	Familiarity	Primacy	Term Rep.
intransitive	targeted	familiar	primary	ldh
transitive	untargeted	crowdsourced	mixed	x dh
ditransitive				x dx

A **system configuration** combines one setting from each column.

If particular settings correlate highly with item characteristics (intransitive / transitive / ditransitive; response complexity), I can optimize the system for new items.

## Sampling data: Response length

	n=14		n=50	n=70
	Fam	Crowd	Crowd	NNS
Intrans	5.5	4.9	4.9	4.9
Trans	6.9	6.3	6.2	6.7
Ditrans	7.8	7.2	7.2	8.3
Target	6.5	5.4	5.4	6.3
Untarg	6.9	6.8	6.8	6.9
primary	N/A	5.7	5.8	6.6
mixed	6.7	6.5	6.4	N/A
Total	6.7	6.1	6.1	6.6

**Table:** Comparing average response length (in words) for the samples used throughout this chapter as NS models and NNS test sets, in total and by parameter setting.

# Annotation features

First iteration: **accuracy (A)** & **native-likeness (NL)**

- ▶ **2:** +A, +NL > **1:** +A, -NL > **0:** -A, -NL
- ▶ Not operationalizable: e.g., response is accurate w.r.t. prompt but adds unverifiable details; is this still *accurate*?
- ▶ Not *reliable*, not *valid*;

This was scrapped and I settled on the 5 binary features.

## Annotation features

Inter-rater reliability for two annotators and 10% of the dataset:  
yes annotations for Annotator 1 (note skewedness), expected  
chance agreement (*Chance*), actual observed agreement  
(*Observed*) and Cohen's kappa (*Kappa*)

Set	A1Yes	Chance	Observed	Kappa
CORE EVENT	.733	.601	.923	.808
ANSWERHOOD	.834	.721	.982	.936
GRAMMATICALITY	.861	.768	.960	.827
INTERPRETABILITY	.818	.682	.919	.744
VERIFIABILITY	.845	.719	.968	.884
Intransitive	.863	.758	.978	.910
Transitive	.780	.653	.949	.853
Ditransitive	.812	.678	.924	.764

# Weighting features

Raters perform holistic preference test (blind to annotations)

<i>What is the boy doing?</i>	Pref?	Core	Ansr	Gram	Intrp	Verif
He is eating food.	yes	0	1	1	1	1
He may get fat eating.	no	0	0	1	1	0
He is hungry.	no	0	0	1	0	1
the boy is eating pizza	yes	1	1	1	1	1
The child is about to eat pizza.	yes	1	0	1	1	1
he eating.	no	0	1	0	1	1
Totals preferred responses		2	2	3	3	3
Totals dispreferred responses		0	1	2	2	2
Net preferred (pref - dispref)		2	1	1	1	1
Feature weight		.333	.167	.167	.167	.167
*Real feature weight		.365	.093	.055	.224	.263

# Mean Average Precision

*Average precision* represents the area under the precision-recall curve.

*Mean average precision* is an average over multiple average precisions (here it's from multiple PDT items or datasets).

This is a simplification, but for our purposes here, we can think of MAP as a measure of how well a ranking separates "yes" and "no" annotations.

Bad	Okay	Good
-----	------	------

<b>yes</b>	<b>yes</b>	<b>yes</b>
no	<b>yes</b>	<b>yes</b>
no	no	<b>yes</b>
<b>yes</b>	<b>yes</b>	no
<b>yes</b>	no	no

AP →	.51	.64	1.0
------	-----	-----	-----

$$\text{MAP} = (0.51 + 0.64 + 1.00)/3 \approx 0.72$$



## Predicting features: ANSWERHOOD MAP

	Crowd NS model = 14					Crowd NS model = 50				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.868	.871	<b>.878</b>	.881	.869	.866	.868	<b>.874</b>	.881	.868
Tran	.816	.819	<b>.846</b>	.845	.838	.818	.823	<b>.851</b>	.845	.838
Ditr	.824	.826	<b>.841</b>	.837	.833	.821	.822	<b>.840</b>	.837	.833
Targ	.787	.788	<b>.810</b>	.817	.799	.787	.789	<b>.811</b>	.817	.798
Untg	.885	.890	<b>.900</b>	.892	.894	.883	.886	<b>.899</b>	.892	.895
Prim	.837	.840	<b>.854</b>	.854	.845	.837	.840	<b>.854</b>	.854	.846
Mix	.835	.838	<b>.857</b>	.854	.848	.833	.835	<b>.856</b>	.854	.847
Total	.836	.839	<b>.855</b>	.854	.847	.835	.838	<b>.855</b>	.854	.846

- ▶  $x dx > x dh > ldh$ ;
- ▶ Model size makes little difference;

Predicting ANSWERHOOD is relatively simple; requires only small model and bag-of-words representation.

## Predicting features: ANSWERHOOD MAP

	Familiar NS model = 14					Crowd NS model = 14				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.868	.871	<b>.882</b>	.881	.868	.869	.873	<b>.878</b>	.881	.870
Tran	.824	.826	<b>.852</b>	.845	.840	.817	.818	<b>.847</b>	.845	.840
Ditr	.820	.822	<b>.846</b>	.837	.832	.820	.822	<b>.845</b>	.837	.835
Targ	.786	.787	<b>.815</b>	.817	.798	.785	.787	<b>.813</b>	.817	.802
Untg	.889	.892	<b>.904</b>	.892	.896	.885	.889	<b>.900</b>	.892	.894
Total	.837	.840	<b>.860</b>	.854	.847	.835	.838	<b>.857</b>	.854	.848

- ▶  $x dx > x dh > ldh$ ;
- ▶ familiar > crowdsourced;

Predicting ANSWERHOOD is relatively simple; requires only small model and bag-of-words representation.

## Predicting features: GRAMMATICALITY MAP

	Crowd NS model = 14					Crowd NS model = 50				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.868	.870	<b>.872</b>	.887	.866	.863	.864	<b>.866</b>	.887	.864
Tran	.753	.756	<b>.757</b>	.781	.757	.758	.760	<b>.761</b>	.781	.757
Ditr	.682	.685	<b>.700</b>	.695	.694	.679	.685	<b>.697</b>	.695	.693
Targ	.777	.778	<b>.784</b>	.800	.782	.776	.776	<b>.783</b>	.800	.781
Untg	.758	.763	<b>.769</b>	.776	.762	.757	.762	<b>.766</b>	.776	.761
Prim	.769	.773	<b>.776</b>	.788	.770	.768	.770	<b>.774</b>	.788	.770
Mix	.766	.768	<b>.776</b>	.788	.774	.765	.768	<b>.775</b>	.788	.772
Total	.768	.770	<b>.776</b>	.788	.772	.767	.769	<b>.775</b>	.788	.771

In *most* cases:

- ▶  $x dx > x dh > ldh$ ;

Predicting GRAMMATICALITY is relatively simple; requires only small model and bag-of-words representation.

## Predicting features: GRAMMATICALITY MAP

	Familiar NS model = 14					Crowd NS model = 14				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.863	.864	<b>.873</b>	.887	.863	.868	.869	<b>.874</b>	.887	.869
Tran	.760	.759	<b>.762</b>	.781	.760	.752	.754	<b>.757</b>	.781	.758
Ditr	.678	.685	<b>.698</b>	.695	.698	.678	.680	<b>.699</b>	.695	.696
Targ	.776	.776	<b>.787</b>	.800	.783	.776	.777	<b>.786</b>	.800	.786
Untg	.757	.762	<b>.768</b>	.776	.764	.756	.759	<b>.767</b>	.776	.763
Total	.767	.769	<b>.778</b>	.788	.773	.766	.768	<b>.776</b>	.788	.774

In *most* cases:

- ▶  $x dx > x dh > ldh$ ;
- ▶ familiar 14NS > crowd 14NS > crowd 50NS;

Predicting GRAMMATICALITY is relatively simple; requires only small model and bag-of-words representation.

## Predicting features: INTERPRETABILITY MAP

	Crowd NS model = 14					Crowd NS model = 50				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.932	.931	<b>.933</b>	.930	.922	.928	.927	<b>.933</b>	.930	.923
Tran	<b>.823</b>	.821	.811	.803	.806	<b>.821</b>	.816	.812	.803	.804
Ditr	.789	.784	<b>.794</b>	.721	.777	.786	.782	<b>.792</b>	.721	.772
Targ	.835	.832	<b>.836</b>	.804	.828	.833	.829	<b>.834</b>	.804	.826
Untg	<b>.862</b>	.858	.856	.833	.842	<b>.857</b>	.855	<b>.857</b>	.833	.840
Prim	<b>.847</b>	.845	.846	.818	.837	.845	.842	<b>.846</b>	.818	.833
Mix	<b>.849</b>	.846	.846	.818	.833	.844	.841	<b>.845</b>	.818	.833
Total	<b>.848</b>	.845	.846	.818	.835	<b>.845</b>	.842	<b>.845</b>	.818	.833

- ▶ 14NS crowdsourced > 50NS crowdsourced;
- ▶ intransitives & ditransitives work best with xdx;
- ▶ transitives work best with ldh;
  - ▶ Why? Transitive responses are relatively homogenous;  
Annotators relatively strict;

## Predicting features: INTERPRETABILITY MAP

	Familiar NS model = 14					Crowd NS model = 14				
	ldh	xdh	xdx	WAR	SBERT	ldh	xdh	xdx	WAR	SBERT
Intr	.930	.930	<b>.934</b>	.930	.923	<b>.933</b>	.931	.932	.930	.922
Tran	<b>.822</b>	.819	.811	.803	.805	<b>.826</b>	.824	.811	.803	.805
Ditr	.787	.786	<b>.796</b>	.721	.782	.788	.783	<b>.795</b>	.721	.772
Targ	.835	.833	<b>.836</b>	.804	.830	<b>.835</b>	.832	<b>.835</b>	.804	.825
Untg	<b>.858</b>	.857	<b>.858</b>	.833	.843	<b>.863</b>	.859	.857	.833	.841
Total	<b>.847</b>	.845	<b>.847</b>	.818	.837	<b>.849</b>	.846	.846	.818	.833

- ▶ 14NS crowdsourced > 14NS familiar;
- ▶ intransitives & ditransitives work best with xdx;
- ▶ transitives work best with ldh;
  - ▶ Why? Transitive responses are relatively homogenous;  
Annotators relatively strict;

## Predicting features: VERIFIABILITY MAP

	Crowd NS model = 14					Crowd NS model = 50				
	ldh	xdh	xdx	WAR	SBERT	ldh	xdh	xdx	WAR	SBERT
Intr	.852	.852	<b>.853</b>	.866	.840	.849	.849	<b>.851</b>	.866	.836
Tran	<b>.809</b>	.808	.803	.798	.787	<b>.807</b>	.806	.803	.798	.785
Ditr	.814	.812	<b>.815</b>	.780	.798	.811	.809	<b>.812</b>	.780	.796
Targ	<b>.825</b>	.824	.825	.815	.812	<b>.825</b>	.824	.823	.815	.810
Untg	<b>.825</b>	.824	.822	.815	.805	<b>.820</b>	.819	.820	.815	.802
Prim	<b>.826</b>	.824	.823	.815	.808	<b>.824</b>	.823	.822	.815	.806
Mix	<b>.825</b>	.824	.824	.815	.808	<b>.821</b>	.821	.821	.815	.805
Total	<b>.825</b>	.824	.824	.815	.808	<b>.823</b>	.822	.822	.815	.806

- ▶ 14NS crowd > 50NS crowd;
- ▶ intransitives & ditransitives work best with xdx;
- ▶ transitives work best with ldh;
  - ▶ Why? Transitive responses are relatively homogenous;  
Annotators relatively strict;

## Predicting features: VERIFIABILITY MAP

	Familiar NS model = 14					Crowd NS model = 14				
	ldh	x dh	x dx	WAR	SBERT	ldh	x dh	x dx	WAR	SBERT
Intr	.847	.847	<b>.852</b>	.866	.836	.852	.852	<b>.854</b>	.866	.843
Tran	<b>.808</b>	.807	.803	.798	.787	<b>.807</b>	.807	.802	.798	.786
Ditr	.811	.811	<b>.812</b>	.780	.802	.815	.812	<b>.817</b>	.780	.796
Targ	.821	.821	<b>.822</b>	.815	.814	.824	.824	<b>.826</b>	.815	.811
Untg	<b>.824</b>	.822	.823	.815	.803	<b>.825</b>	.824	.823	.815	.806
Total	.822	.822	<b>.823</b>	.815	.808	<b>.825</b>	.824	.824	.815	.808

- ▶ 14NS crowd > 14NS familiar;
- ▶ intransitives & ditransitives work best with xdx;
- ▶ transitives work best with ldh;
  - ▶ Why? Transitive responses are relatively homogenous;  
Annotators relatively strict;



# Predicting quality: Targeting

Spearman rank correlations: System vs. WAR (benchmark)

		targeted		untargeted	
		System	SBERT	System	SBERT
	count	180	60	180	60
14NS	mean	<b>.380</b>	.530	.300	.444
	median	<b>.369</b>	.545	.314	.472
50NS	mean	<b>.393</b>	.550	<i>.305</i>	.469
	median	<b>.389</b>	.564	<i>.323</i>	.496

- ▶ targeted > untargeted
- ▶ 50NS models > 14NS models
  - ▶ Model size effect is most pronounced for targeted

## Predicting quality: Familiarity

Spearman rank correlations: System vs. WAR (benchmark)

		familiar		crowdsourced	
		System	SBERT	System	SBERT
	count	180	60	180	60
14NS	mean	.338	.499	<b>.339</b>	.481
	median	<b>.329</b>	.513	.326	.500

- ▶ System: No discernible difference for familiar vs crowdsourced
- ▶ SBERT: familiar > crowdsourced
  - ▶ NNS STTR < familiar STTR < crowdsourced STTR

## Predicting quality: Primacy

Spearman rank correlations: System vs. WAR (benchmark)

		primary		mixed	
		System	SBERT	System	SBERT
	count	180	60	180	60
14NS	mean	.339	.493	<b>.340</b>	.481
	median	.326	.517	<b>.334</b>	.500
50NS	mean	<b>.354</b>	.514	.344	.505
	median	.345	.532	<b>.350</b>	.518

- ▶ System: 14NS: primary < mixed (slight difference)
- ▶ System: 50NS: primary  $\approx$  mixed
- ▶ SBERT: primary > mixed
- ▶ System & SBERT: 50NS > 14NS
  - ▶ System: model size effect is greatest for primary

# Predicting quality: Term Representation

Spearman rank correlations: System vs. WAR (benchmark)

		System			SBERT
		ldh	xdh	xdx	(text)
	count	120	120	120	40
14NS	mean	.333	.336	<b>.351</b>	.487
	median	.318	<b>.344</b>	.330	.507
50NS	mean	<b>.350</b>	.349	.348	.509
	median	.364	<b>.374</b>	.331	.523

- ▶ SBERT: 50NS > 14NS
- ▶ System: for ldh & xdh: 50NS > 14NS;
  - ▶ Model size effect is greater for ldh
- ▶ System: for xdx: NS14 > NS50 (very slight)

## Predicting quality: Term Normalization

Response A	Response B	Non-norm -alized weight	Normal -ized weight
The girl is singing	The girl in the cute purple dress is singing a song		
det(the, girl)	det(the, girl)	.143 (2/14)	.175
nsubj(girl, sing)	nsubj(girl, sing)	.143	.175
	det(the, dress)	.071	.050
	<b>amod(cute, dress)</b>	<b>.071 (1/14)</b>	<b>.050</b>
	<b>amod(purple, dress)</b>	<b>.071</b>	<b>.050</b>
	prep_in(dress, girl)	.071	.050
aux(be, sing)	aux(be, sing)	.143	.175
root(sing, ROOT)	root(sing, ROOT)	.143	.175
	det(a, song)	.071	.050
	dobj(song, sing)	.071	.050
4	10	1.0 (14/14)	1.0

A 2-response toy NS model. Normalizing for response length so each *response* (not *dependency*) in model carries equal weight reduces weight of some extraneous dependencies, but performance suffers overall.