# Sprout: Crowd-Powered Task Design for Crowdsourcing

Jonathan Bragg University of Washington (Stanford University) Mausam

Indian Institute of Technology

Delhi

Daniel S. Weld University of Washington

PAUL G. ALLEN SCHOOL

JTER SCIENCE & ENGINEERING





### Crowdsourcing Task Platforms



2

#### Is this a car?



#### Is this a car?



#### Is this a real (not cartoon) car?



#### Is this a real (not cartoon) car that a person can drive on a street?

Yes

No

Low Inter-annotator

Agreement



6

Is this a single real (not cartoon) car that a person can drive on a street?





Is this a single real (not cartoon) car that a person can drive on a street? The whole car must be visible for it to count.



## Iterative Task Design Loop

#### Perform (explore) the task

- Task definition changes
- Concept evolution [Kulesza et al., CHI '14]



Instructions Revise Training & Testing

Revise

#### Run (debug) the task

- Worker mistakes initially unknown
- Model mismatch / debugging



Low Inter-Annotator Agreement

## Task Design is a Major Problem

- Difficulty for Requesters; Importance for Work Quality
  - [Alonso & Mizzaro, Information Processing and Management '12; Papoutsaki et al. HCOMP '15; Liu et al., NAACL '16]
- Risk to Workers: Task Design Flaws & Unfair Rejections
  - [McInnis et al., CHI '16; Gadiraju et al., HT '17; Wu & Quinn, HCOMP '17]

## ...and it matters



[Freelancing in America: 2017. Upwork and Freelancers Union] <sup>11</sup>

## How can we design tools for task design?

#### Perform (explore) the task

- Task definition changes
- Concept evolution [Kulesza et al., CHI '14]



Instructions Revise Training & Testing

Revise

#### Run (debug) the task

- Worker mistakes initially unknown
- Model mismatch / debugging

Low Inter-

Annotator

Agreement



## Sprout: Debugging-Prioritized Exploration



## Sprout: Debugging-Prioritized Exploration



## Sprout: Debugging-Prioritized Exploration





## Sprout Compiles Test Q's into Training & Testing



Gated Instruction [Liu et al., NAACL '16]

Confusions viewed	
34	
Confusions	a)
No 7	Ves
✓is a train car	6
✓is multiple cars	
✓is a truck	
✓Is a train	0
✓Is a bumper car	0
has a car as the main subject	0
✓Is a military vehicle	0
✓Is the interior of a car	0
is an automobile	•
✓is a car on a ferry boat	•
349 444	
Related is car on a boat (1 in common)	
Finckulas more than one (1 in	
common)	
444	
is a street car	
is a trolley car	
is a subway train car	
is a toy car	0

Pre	evie	w					
				444			4
	Answ	er Reas	on				
1	7	is a c	ar on a fe	my boat			
1	9	includ	ies more t	than one			
1	yes	There photo	are seve graph.	nal clearly	visible ve	hicles in t	10
	-	549	10			-	

Ð

#### Instructions Your instructions for workers go here. Use twitter mention notation to reference items, for example, @18 refers to item 10 and will preview as . You may also use other types of Markdown to format your instructions, like This will be a list item . This will be anbolides νÖ ¥1 Write Preview С is this an image of a car? Select "yes" even if there are multiple cars (e.g., @444) Recommended test questions 549 because you mentioned 444 in the d instructions Test questions O yes 100 Submit (07:21 remaining)



#### Structured labeling implementation

## Requester User Study

- Goal: improve underspecified instructions prompts
- 2 domains from prior work (car images and travel websites)
- Baseline no-crowd interface: Structured Labeling [Kulesza et al. CHI 2014]
- Within-subjects (counterbalanced interface)
- 11 participants
  - varying crowdsourcing experience
  - graduate and undergraduate students
  - 5 male, 5 female, 1 other

yes						
552	526	533	176	238		
	multi	ple cars				\$
449	•					
+ Dra	ig here f	or new (	group			
maybe						
273						
+ 0.00	a hara f				_	
T Dra	g nere i	or new ş	jro	458		
no						
			-		_	
497	292	245	510	458	175	545
145				200	)	

## Sprout preferred by requesters

"Which interface did you prefer for creating instructions?"



"[c]ategorizing the inputs, showing me the cases where there was confusion, etc., made it SUPER easy to identify cases that needed clarification." (P1)

"The similarity metrics seem to be working great and the suggested items are great for testing the points I emphasized in the instructions." (P2)

20-29% of test questions from suggestions

### Sprout aids instruction comprehensiveness

- Longer instructions and more examples (mean=4.6 vs 2.6 on travel task)
- More resolved ambiguous categories (mean=4.0 vs. 2.8 on cars task)
  - 2 experts coded number of distinct categories

## Some requesters wanted to use both tools

- Structured labeling benefits
  - No bias (P2)
  - "Let me try doing the task myself" (P3)
- Tradeoff: "It sucks that you have to start from a completely blank slate. [SPROUT] gave you some more support." (P4)



## Future Work

- Expanded roles for workers & algorithms
  - More automated instructions: Infer requester decisions ("legal precedent")
  - Engaging, progressive tutorial
  - Interface design
  - Workflow design (task decomposition)
- More studies
  - Field deployments
- Beyond labeling tasks
  - Open-ended & creative tasks

## Thanks! https://jonathanbragg.com

## Sprout: Crowd-Powered Task Design for Crowdsourcing

Jonathan Bragg	Mausam	Dan Weld	NSF
		1 1\A/	ONR
	III Deim	U VV	WRF
(Stanford)			Google

IBM

Bloomberg

Government of India

![](_page_25_Picture_6.jpeg)

![](_page_25_Picture_7.jpeg)

![](_page_25_Picture_8.jpeg)