

# Human Factors in Crowdsourcing

Gianluca Demartini

University of Queensland, Australia

<http://gianlucademartini.net>

@eglu81

# Research Interests

- **Entity-centric Information Access (2005-now)**
  - Structured/Unstruct data (SIGIR 12), TRank (ISWC 13, WSemJ 16)
  - NER in Scientific Docs (WWW 14), Prepositions (CIKM 14)
  - IR Evaluation (ECIR 16 Best Paper Award, IRJ 2015, CIKM 17)
- **Hybrid Human-Machine Systems (2012-now)**
  - ZenCrowd (WWW 12, VLDBJ), CrowdQ (CIDR 13)
  - Human Memory based Systems (WWW 14, PVLDB)
  - Hybrid systems overview (COMNET, 2015)
- **Better Crowdsourcing Platforms (2013-now)**
  - **Platform Dynamics** (WWW 15)
  - Pick-a-Crowd (WWW 13), **Malicious Workers** (CHI 15)
  - Scale-up Crowdsourcing (HCOMP 14), Scheduling (WWW 16)
  - **Timeout** (HCOMP 16), **Environment** (UBICOMP 17)

Thanks to:



Project Duration 2017-2019. Funded under the H2020-ICT-14-2016 topic Big Data PPP: cross-sectorial and cross-lingual data integration and experimentation. Total cost: 2.9M EUR.

# FashionBrain: Understanding Europe's Fashion Data Universe



eXascale Infolab

## Project Objectives:

- Novel Shopping Experience: **Make Images Searchable**
  - Product search and recommendation
- Shift Traffic away from Web Search Engines to **Retailer's Mobile Apps**
  - By providing custom shopping experiences and advanced search tools
- Detect Influencers and **Predict Fashion Trends**
  - Time Series Analysis; Social Media data
- **Share Insights** with Cross Industry Partner Network
  - Data Integration infrastructure based on HDFS and column stores

[fashionbrain-project.eu](http://fashionbrain-project.eu)

# Crowdsourcing

- "Simply defined, crowdsourcing represents the act of a company or institution taking a function once performed by employees and outsourcing it to an undefined (and generally large) network of people in the form of an **open call**. This can take the form of peer-production (when the job is performed **collaboratively**), but is also often undertaken by sole **individuals**. The crucial prerequisite is the use of the open call format and the **large network of potential laborers**."

[Howe, 2006]

# Incentives in Crowdsourcing

- **Extrinsic motivation** if task is considered boring, dangerous, useless, socially undesirable, dislikable by the performer.
  - Paid Crowdsourcing
- **Intrinsic motivation** is driven by an interest or enjoyment in the task itself.
  - Fun (enjoyment) / Games with a purpose
  - Community (belonging, desire to help)
  - Citizen Science

# Paid Micro-Task Crowdsourcing

A Crowdsourcing Platform allows **requesters** to publish a crowdsourcing request (*batch*) composed of multiple tasks (*HITs*)

Programmatically Invoke the crowd with APIs or using a website

**Workers** in the crowd complete tasks and obtain a monetary reward

# Amazon MTurk



## Make Money by working on HITs

HITs - *Human Intelligence Tasks* - are individual tasks that you work on. [Find HITs now.](#)

### As a Mechanical Turk Worker you:

- Can work from home
- Choose your own work hours
- Get paid for doing good work



## Get Results from Mechanical Turk Workers

Ask workers to complete HITs - *Human Intelligence Tasks* - and get results using Mechanical Turk. [Register Now](#)

### As a Mechanical Turk Requester you:

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results



# MTurk is a Marketplace for HITs

## All HITs

1-10 of 3454 Results

Sort by:   

[Show all details](#) | [Hide all details](#)

1 2 3 4 5 [Next](#) [Last](#)

<b>Provide Information about a Product</b> Requester: <a href="#">requester</a>	<b>HIT Expiration Date:</b> May 23, 2015 (4 weeks 1 day) <b>Time Allotted:</b> 25 minutes	<b>Reward:</b> \$0.05 <b>HITs Available:</b> 11526	<a href="#">View a HIT in this group</a>
<b>Product Attribute Tagging - April 17th Please read the instructions</b> Requester: <a href="#">slee</a>	<b>HIT Expiration Date:</b> May 23, 2015 (4 weeks 2 days) <b>Time Allotted:</b> 60 minutes	<b>Reward:</b> \$0.03 <b>HITs Available:</b> 23887	<a href="#">View a HIT in this group</a>
<b>Inv B 2</b> Requester: <a href="#">rohzi0d</a>	<b>HIT Expiration Date:</b> May 22, 2015 (4 weeks 1 day) <b>Time Allotted:</b> 48 minutes	<b>Reward:</b> \$0.00 <b>HITs Available:</b> 19740	<a href="#">View a HIT in this group</a>
<b>Geo Result Relevance-Tue Apr 21 10:40:14 PDT 2015</b> Requester: <a href="#">Amazon Requester Inc.</a>	<b>HIT Expiration Date:</b> May 22, 2015 (4 weeks 1 day) <b>Time Allotted:</b> 60 minutes	<b>Reward:</b> \$0.00 <b>HITs Available:</b> 10734	<a href="#">View a HIT in this group</a>
<b>Type the text from the images, carefully. Productivity and bonuses guaranteed.</b> Requester: <a href="#">CopyText Inc.</a>	<b>HIT Expiration Date:</b> Apr 30, 2015 (6 days 23 hours) <b>Time Allotted:</b> 10 minutes	<b>Reward:</b> \$0.01 <b>HITs Available:</b> 10590	<a href="#">View a HIT in this group</a>
<b>Transcribe up to 25 Seconds of Media to Text - Earn up to \$0.12 per HIT!</b> Requester: <a href="#">Crowdsurf Support</a>	<b>HIT Expiration Date:</b> Apr 21, 2016 (51 weeks 6 days) <b>Time Allotted:</b> 15 minutes	<b>Reward:</b> \$0.08 <b>HITs Available:</b> 6702	<a href="#">View a HIT in this group</a>
<b>Fun and Fast Fashion Tagging</b> Requester: <a href="#">gavin</a>	<b>HIT Expiration Date:</b> Apr 28, 2015 (5 days 11 hours) <b>Time Allotted:</b> 60 minutes	<b>Reward:</b> \$0.02 <b>HITs Available:</b> 6460	<a href="#">View a HIT in this group</a>
<b>Geo Result Relevance-Wed Apr 08 14:30:08 PDT 2015</b> Requester: <a href="#">Amazon Requester Inc.</a>	<b>HIT Expiration Date:</b> May 10, 2015 (2 weeks 2 days) <b>Time Allotted:</b> 60 minutes	<b>Reward:</b> \$0.00 <b>HITs Available:</b> 6182	<a href="#">View a HIT in this group</a>
<b>Transcribe up to 25 Seconds of General Content to Text - Earn up to \$0.14 per HIT!</b> Requester: <a href="#">Crowdsurf Support</a>	<b>HIT Expiration Date:</b> Apr 21, 2016 (51 weeks 6 days) <b>Time Allotted:</b> 15 minutes	<b>Reward:</b> \$0.09 <b>HITs Available:</b> 6043	<a href="#">View a HIT in this group</a>
<b>Whac-a-mole by Gaze (hard mode)! Play a 1min eye tracking game in the web browser! 0416</b> Requester: <a href="#">px</a>	<b>HIT Expiration Date:</b> Apr 23, 2015 (8 hours 40 minutes) <b>Time Allotted:</b> 60 minutes	<b>Reward:</b> \$0.10 <b>HITs Available:</b> 4682	<a href="#">View a HIT in this group</a>

1 2 3 4 5 [Next](#) [Last](#)



You must **accept** this HIT before working on it.

#### Data Collection Instructions!

Find the postal address for this Australian company.

- Search on Google, the company's website, YellowPages or Facebook to find the correct postal address for the company below.
- Enter the **full Australian postal address** for the business.
- You may use the research links provided to help.
- **Do not enter incomplete or incorrect details!**

<b>Company name:</b>	Stellar Electrical And Solar Systems
<b>Location:</b>	Australia
<b>Company website:</b>	
<b>Company YellowPages:</b>	
<b>Company Facebook:</b>	
<b>Google search:</b>	<a href="https://www.google.com.au/search?q=%22Stellar Electrical And Solar Systems%22+Australia+postal+address">https://www.google.com.au/search?q=%22Stellar Electrical And Solar Systems%22+Australia+postal+address</a>

**Australian Street Address (ONLY this field is required if complete):**

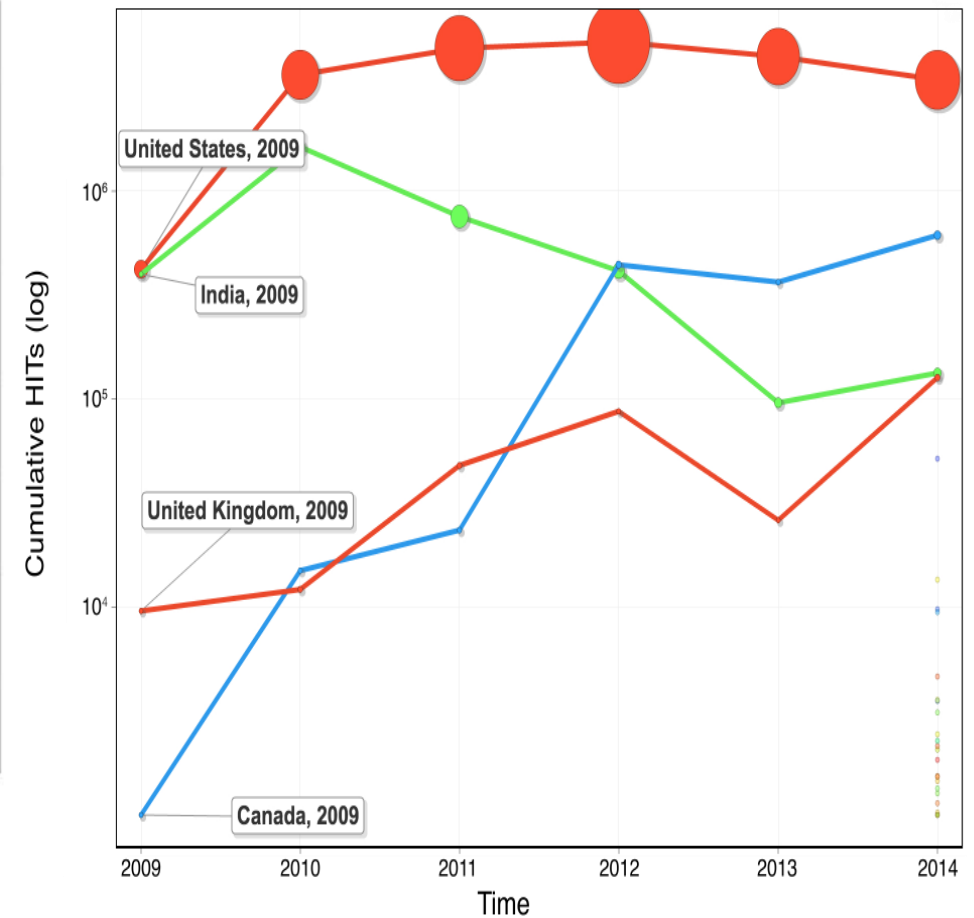
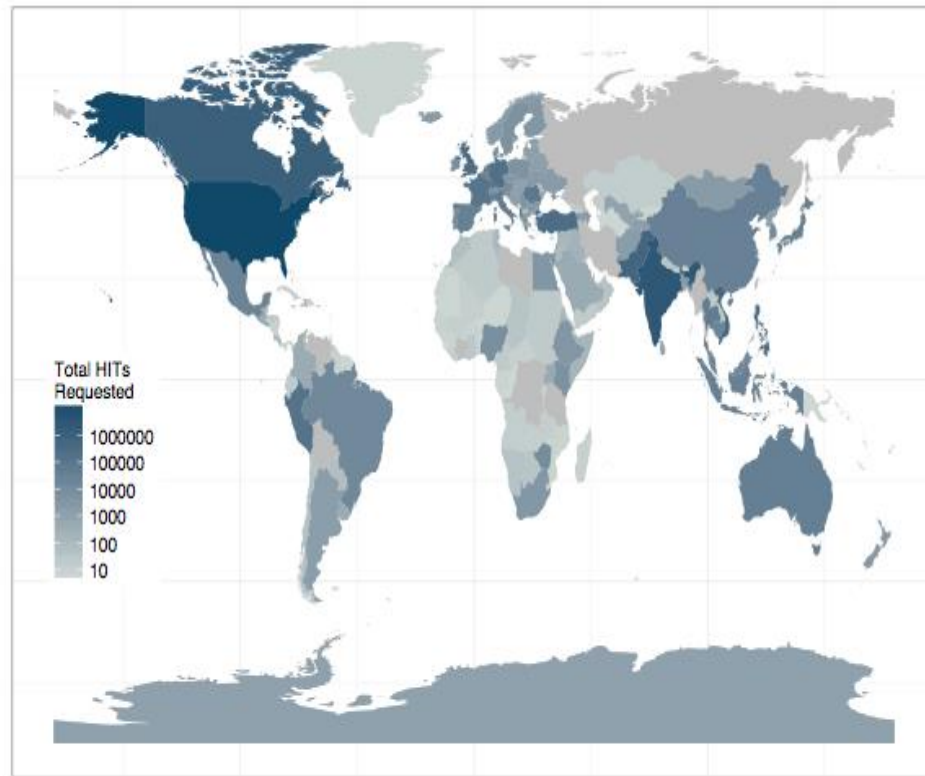
Start typing Australian Street Address...

# MTurk is a Marketplace for HITs

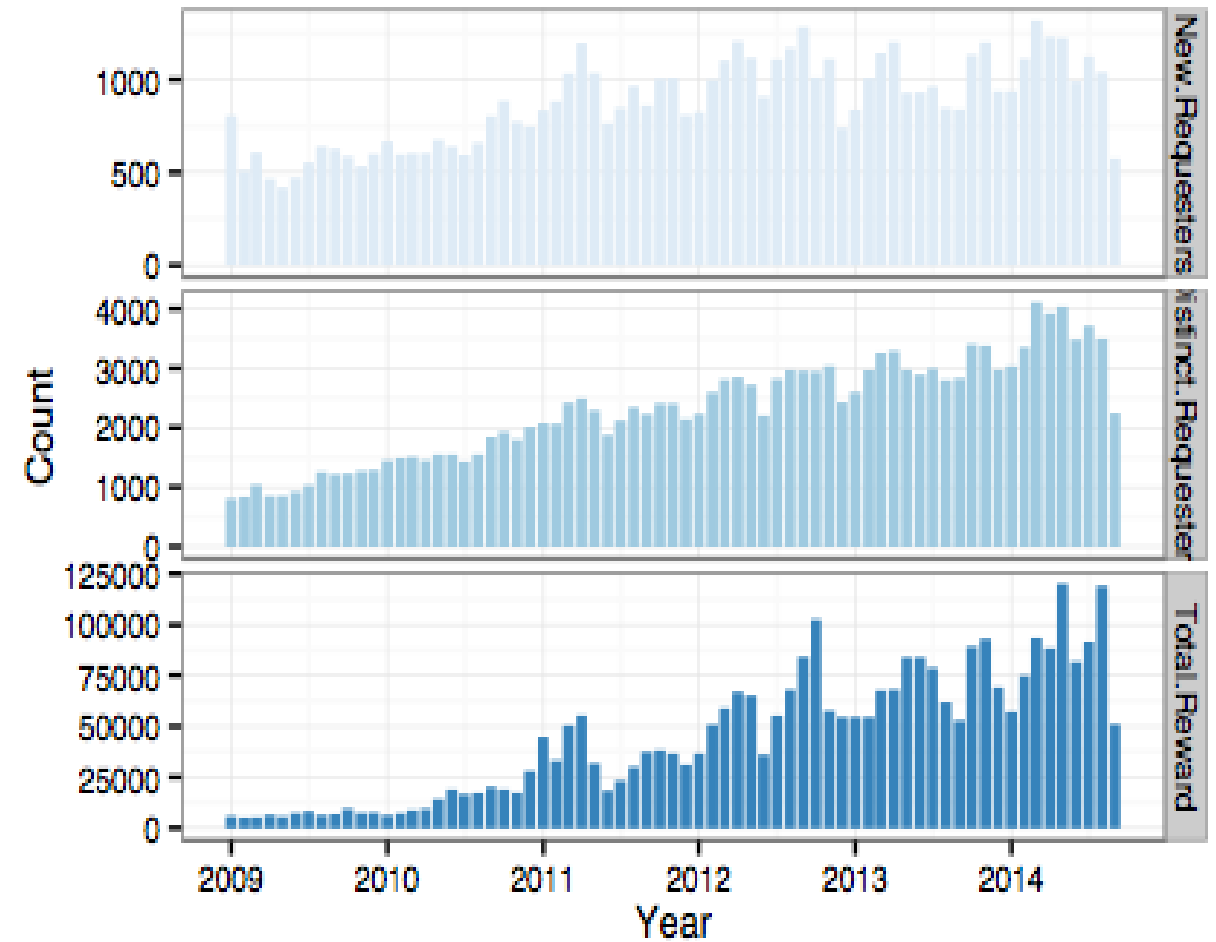
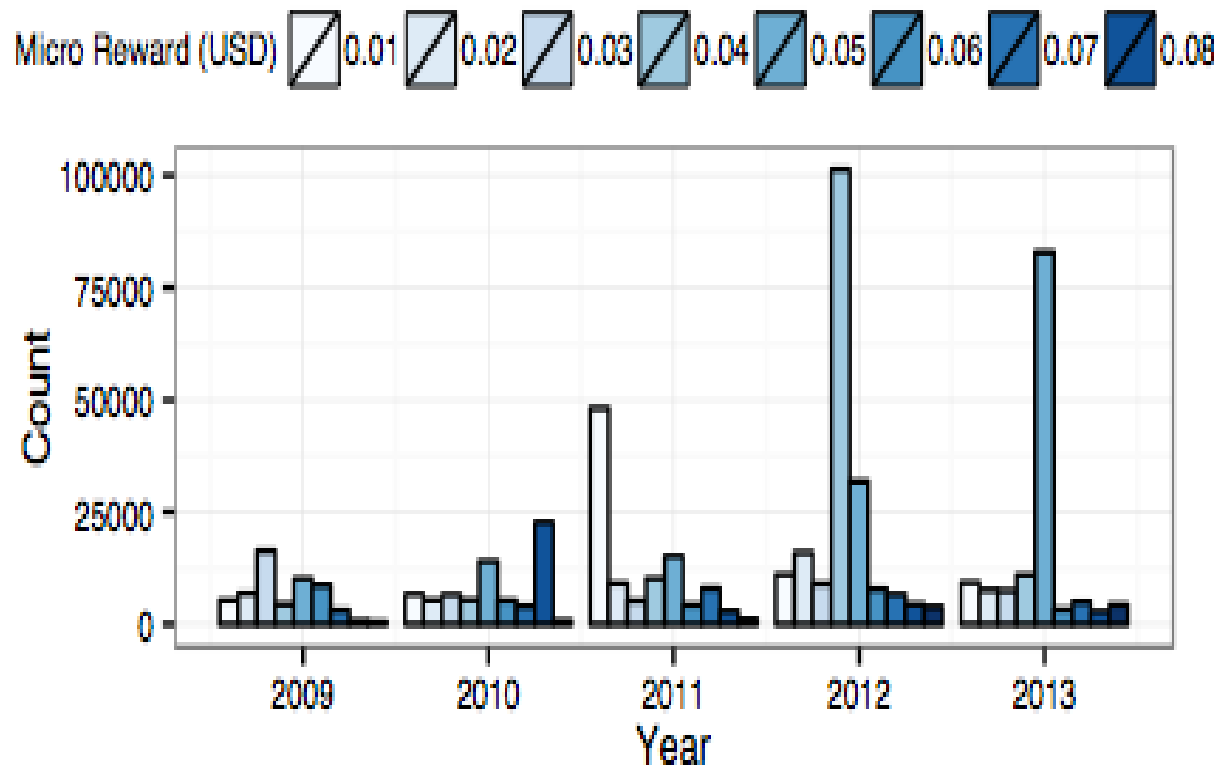
Top-1000 Requesters, report for April 16, 2016 to May 16, 2016

Requester name	hits	reward
Speechpad	23857	\$172,994.63
Percy Liang	883	\$7,320.48
Princeton Vision	51187	\$5,762.44
Stanford GSB Behavioral Lab	3749	\$2,110.70
Chris Callison-Burch	8157	\$2,064.29
RC.org Mechanical Turk	6591	\$2,011.33
VacationrentalAPI	399	\$1,373.50
Med Expertise	869	\$1,303.50
Bluejay Labs	13613	\$1,288.59
YL Testing	1051	\$1,236.83

# Requested Workers



# Reward Distribution

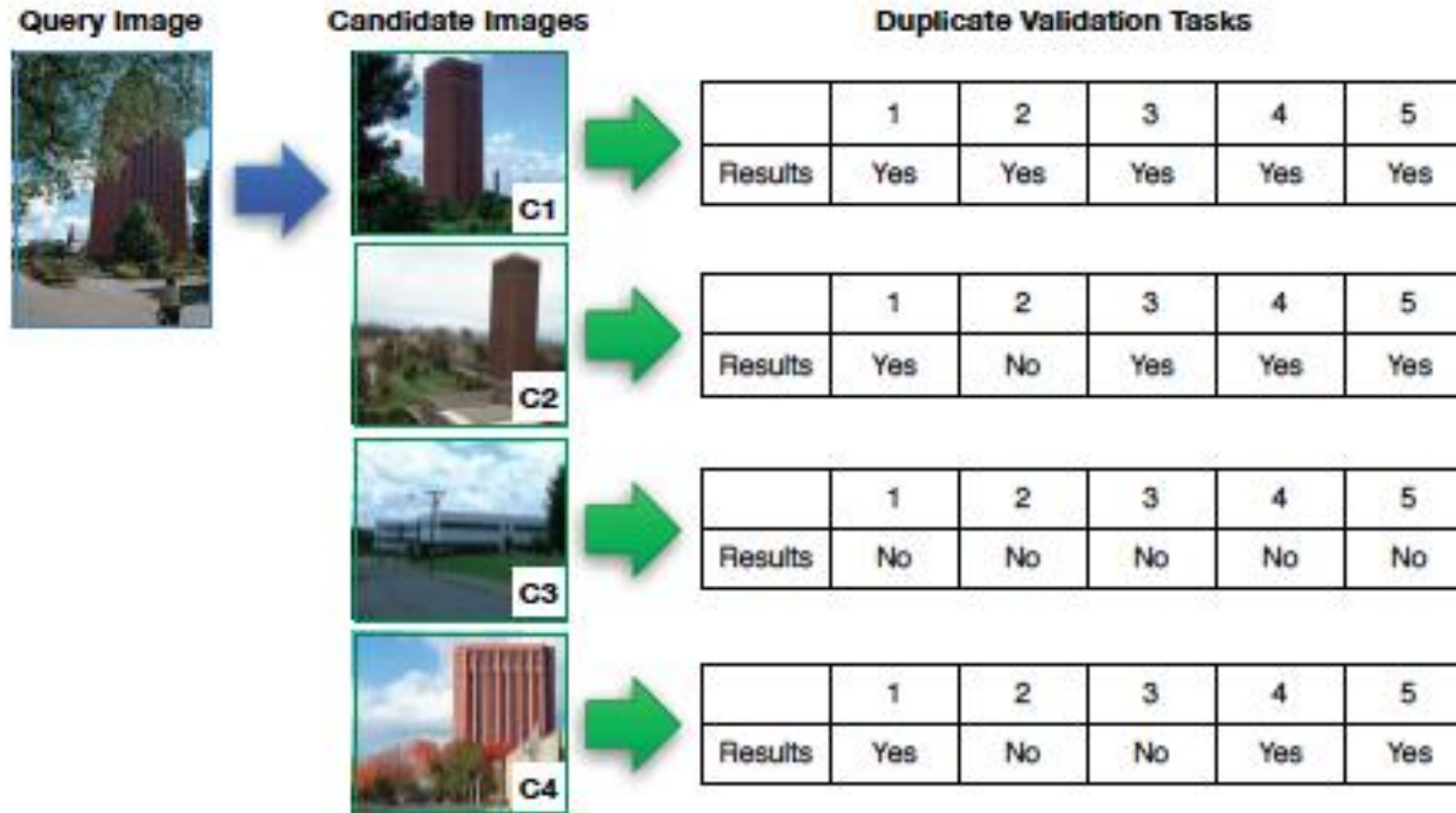


# Hybrid Human-Machine Systems

- Use Machines to scale over large amounts of data
- Keep humans in the loop
  - By means of Crowdsourcing
  - To make sure the quality of the data processing is good
- Crowd for Pre-processing vs Post-processing

G Demartini. Hybrid human-machine information systems: Challenges and opportunities. In: **Computer Networks**, 90, 5-13. 2015

# Hybrid Image Search



Yan, Kumar, Ganesan, CrowdSearch: Exploiting Crowds for Accurate Real-time Image Search on Mobile Phones, Mobisys 2010.

# Human Computation 101 - Summary

- Crowdsourcing is growing in popularity
- It is used both in industry and academia
- For a number of applications across disciplines
- Open questions:
  - How to make sure we get quality results back from a crowdsourcing platforms? (**Effectiveness**)
  - Can we optimize the cost and execution in paid micro-task crowdsourcing? (**Efficiency**)

# Human Factors - Outline

- The effect of limiting **task time** (HCOMP 2016)
- Understanding **malicious behaviors** in paid crowdsourcing (CHI 2015)
- The **modus operandi** of crowd workers (UBICOMP 2017)



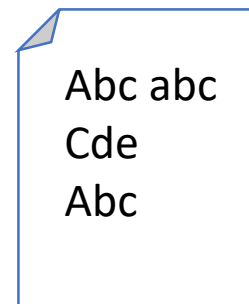
# The Unexpected Benefits of Limiting the Time to Judge

Eddy Maddalena, Marco Basaldella, Dario De Nart, Dante Degl'Innocenti, Stefano Mizzaro, and Gianluca Demartini. Crowdsourcing Relevance Assessments: The Unexpected Benefits of Limiting the Time to Judge. In: **The 4th AAI Conference on Human Computation and Crowdsourcing (HCOMP 2016)**. Austin, Texas, October 2016.

# Crowdsourcing Relevance Judgements

- Task:            Given a Query, Document pair  
                      Is the doc  
                      highly relevant, relevant, partially relevant, not relevant?
- Ask multiple workers
- Aggregate answers to obtain a relevance label

Query: jaguar



Abc abc  
Cde  
Abc

- Highly relevant
- Relevant
- Partially relevant
- Not relevant

## Our Research Question

**Can we limit the time to judge  
to reduce the cost (\$\$) of  
creating IR test collections?**

## Hypothesis

Yes, but with quality loss

# Our Experimental Setup

- **E1 Unbound time** (i.e., the standard approach)
  - 5 judgements per doc, 8 documents, 5 topics, 2 crowds = 400 workers
- **E2 Document shown for a predefined amount of time**
  - **30, 15, 7, 3 seconds.** Each worker to judge 8 docs
- **E3 Same timeout** for all 8 documents (**15 or 30 sec**)
- **E4 Fixed budget:** comparison between
  - more quick judgements
  - few slow judgements

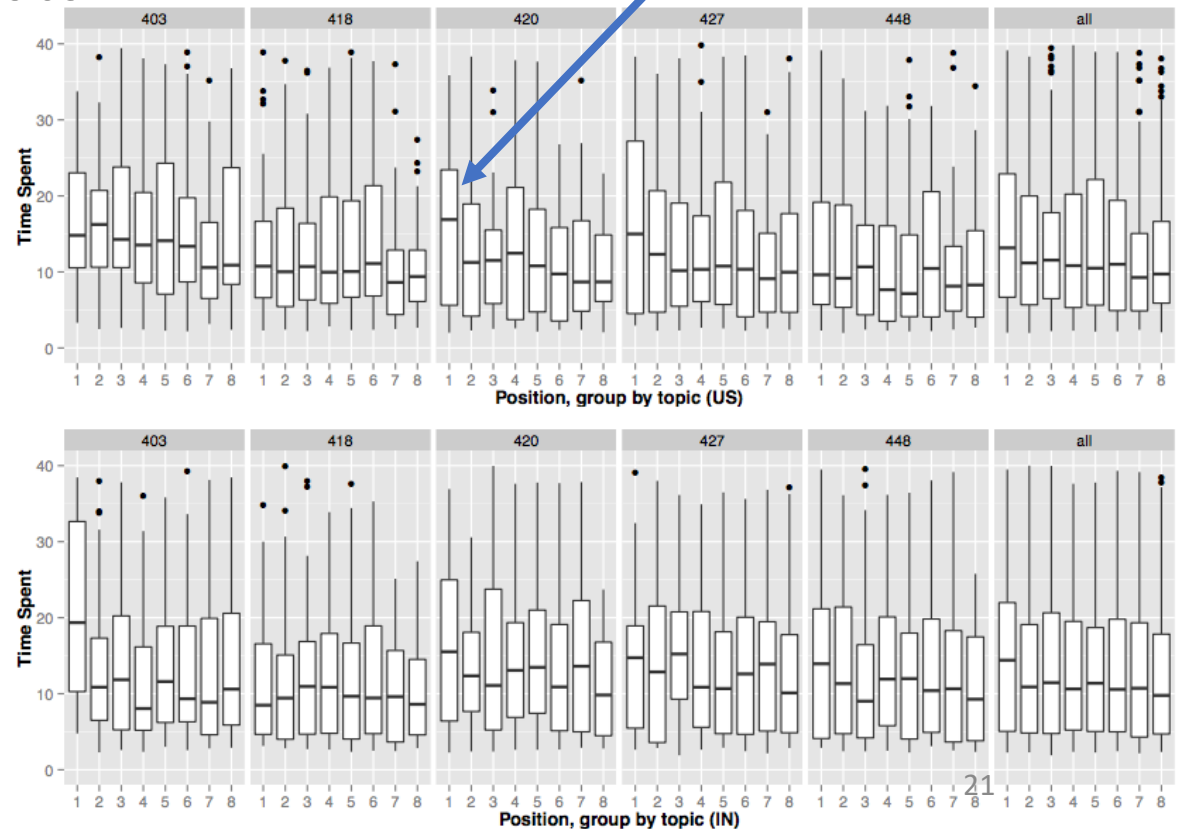
# E1: We Have All the Time in the World

- RQ: **How much time** do crowd workers take to judge the relevance of a document **if no time constrain** is set?
  - 5 workers to judge a permutation of 8 docs



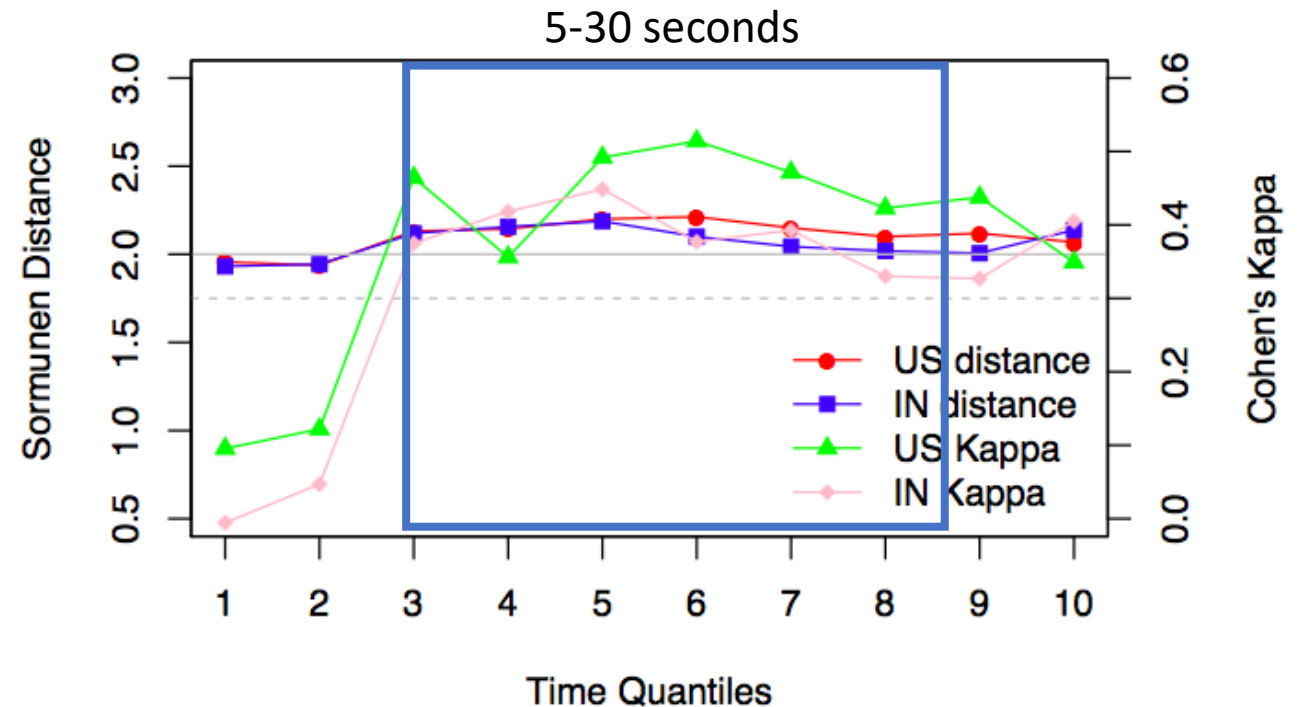
Median:  
13 sec  
Mean  
24-25  
sec

First doc takes longer (learning)



# E1: We Have All the Time in the World

- No correlation of time with
  - Doc length
  - Doc readability
  - Topic
  - Relevance level
- Time vs Quality



	0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
U.S.	2.0	3.2	<b>5.1</b>	7.6	10	13	17	<b>23</b>	32	51	580
IN	1.9	3.4	<b>4.5</b>	7.0	9.9	13	17	<b>22</b>	31	46	630

## E2: Faster! Faster! Sorry, Too Late

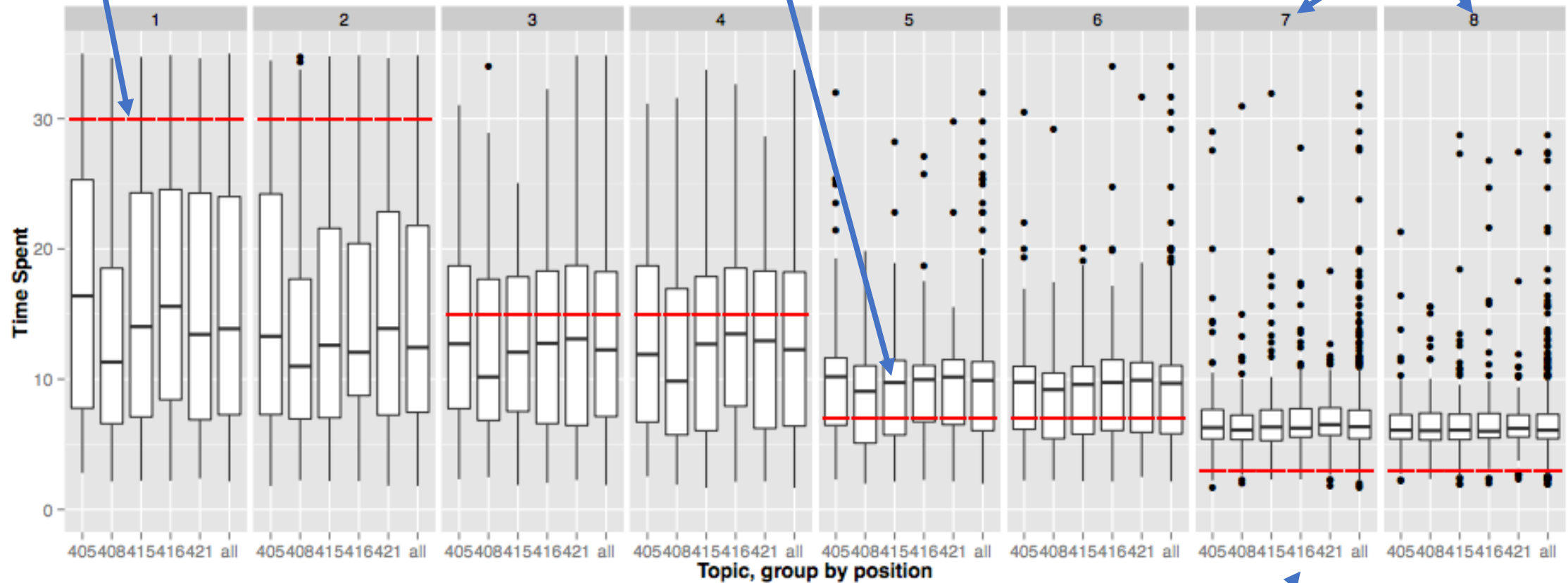
- Understand which is the **minimum amount of time required** to perform relevance judgments
- (max) timeouts: 30, 15, 7, 3 seconds
- Each worker to judge 8 docs, 2 for each timeout (one long, one short)
- Looking at Quality measures:
  - 3 and 7 secs are not enough
  - 15 slightly better than 30 (learning bias for position 1-2?)

# E2: Faster! Faster! Sorry, Too Late

Time when document disappears

Time when judgement is made

Position of the document judged (1-8)



Variance across topics



# E3: Selecting the Best Timeout

- We repeated E1 using 15 and 30 sec timeouts
- 15 seconds timeouts yield consistently better quality judgements
  - Than 30 seconds timeouts
  - Than no timeouts (E1 quality values)

## Our Research Question

**Can we limit the time to judge  
to reduce the cost (\$\$) of  
creating IR test collections?**

~~Hypothesis~~

**Yes, and it improves the quality!**

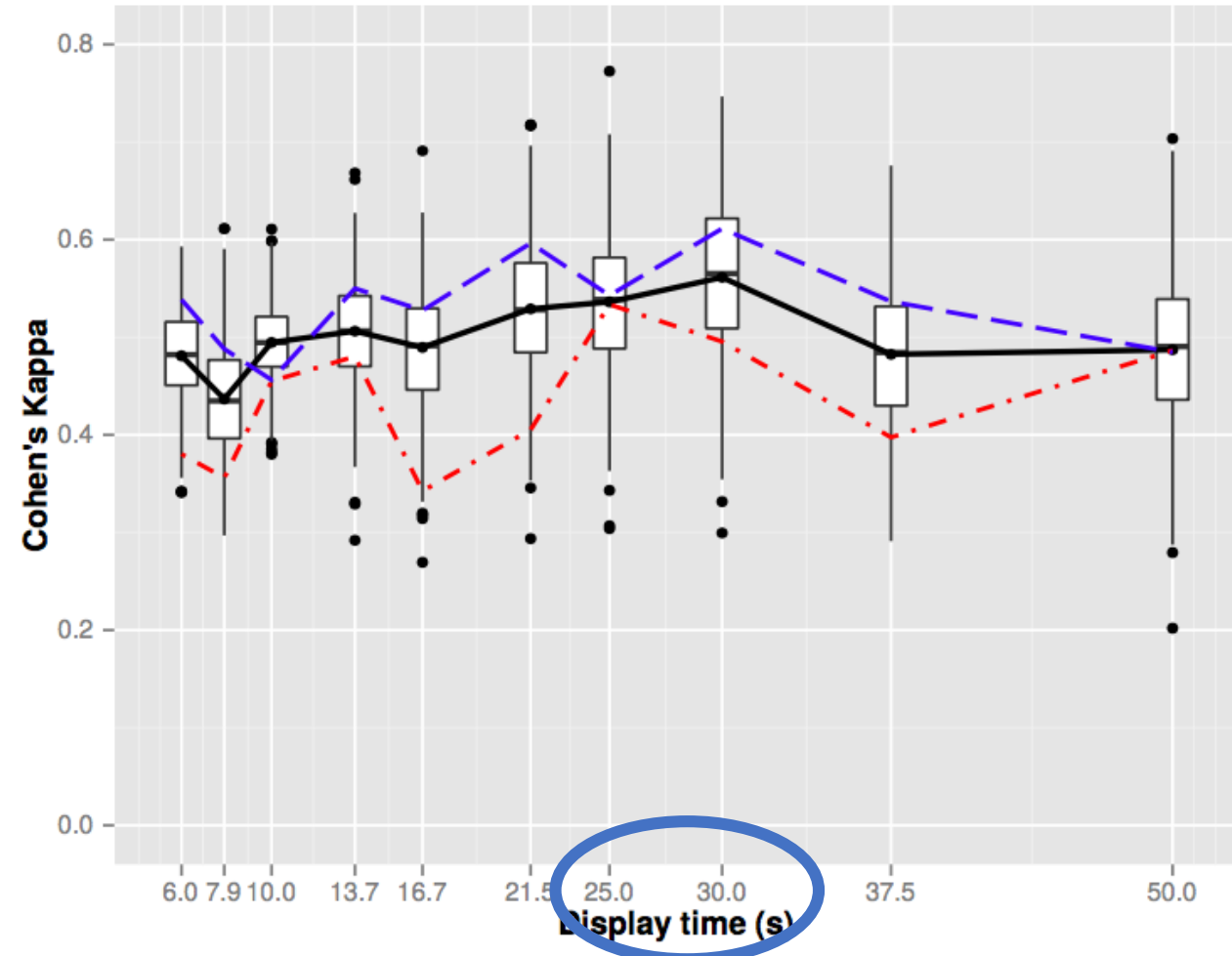
~~Yes, but with quality loss~~

# E4: Many Fast&Furious or a Few Laid-Back?

- **Fixed budget:**
  - small timeout, more workers
  - Long timeout, less workers
- We compared 10 combinations with the same budget

Timeslot(sec)	6	7.9	10	13.7	16.7	21.5	25	30	37.5	50
Assignments	25	19	15	11	9	7	6	5	4	3

- **Highest quality at 25-30 sec**



# Findings

- The **first** couple of judgments done by a worker are of **lower quality**
- Judgements that take **more than 30** seconds are of **lower quality**
- **Time-outs** in relevance judgements HITs can **increase quality**
- The **best timeout** to be used lies in the interval of **25-30 seconds** and does not depend on topic, document, or crowd.

# Discussion

- Crowdsourcing Relevance Judgements for IR Evaluation can be **expensive to scale**
- **Limiting the time** to judge can **control the cost**
- But can also **increase the quality!**
  - By inducing workers to look at the document for a predefined amount of time
- Why? (Hypotheses)
  - With a balance between boredom and stress -> “in the flow”
  - System I and System II thinking

# Understanding Malicious Behaviors

Ujwal Gadiraju, Ricardo Kawase, Stefan Dietze, and Gianluca Demartini. Understanding Malicious Behaviour in Crowdsourcing Platforms: The Case of Online Surveys. In: **Proceedings of the ACM Special Interest Group on Computer Human Interaction (CHI 2015)**. Seoul, South Korea, April 2015

# Quality Control in Paid Crowdsourcing

- Diverse pool of crowd workers
  - Wide range of behavior
  - Various motivations
- Typically adopted solution to prevent/flag malicious activity :  
**Gold-Standard Questions**

# Research Questions

RQ1: Do untrustworthy workers adopt different **methods to complete tasks**, and exhibit different kinds of behavior?

RQ2: Can **behavioral patterns** of malicious workers in the crowd be identified and quantified?



# Design

- CrowdFlower Platform to deploy survey
- Survey questions
  - Demographics
  - Educational & general background
- 34 Questions in total
  - Open-ended
  - Multiple Choice
  - Likert-type
- Responses from 1000 crowd workers
  - Monetary Compensation per worker : 0.2 USD

# RQ1 - Behavioral Patterns

Ineligible  
Workers (IW)

Instruction: Please attempt this microtask ONLY IF you have successfully completed 5 microtasks previously.

Response: *'this is my first task'*

Fast Deceivers  
(FD)

eg: Copy-pasting same text in response to multiple questions, entering gibberish, etc.

Response: *'What's your task?' , 'adasd' , 'fgfgf gsd ljlkj'*

Rule Breakers  
(RB)

Instruction: Identify 5 keywords that represent this task (separated by commas).

Response: *'survey, tasks, history' , 'previous task yellow'*

Smart  
Deceivers (SD)

Instruction: Identify 5 keywords that represent this task (separated by commas).

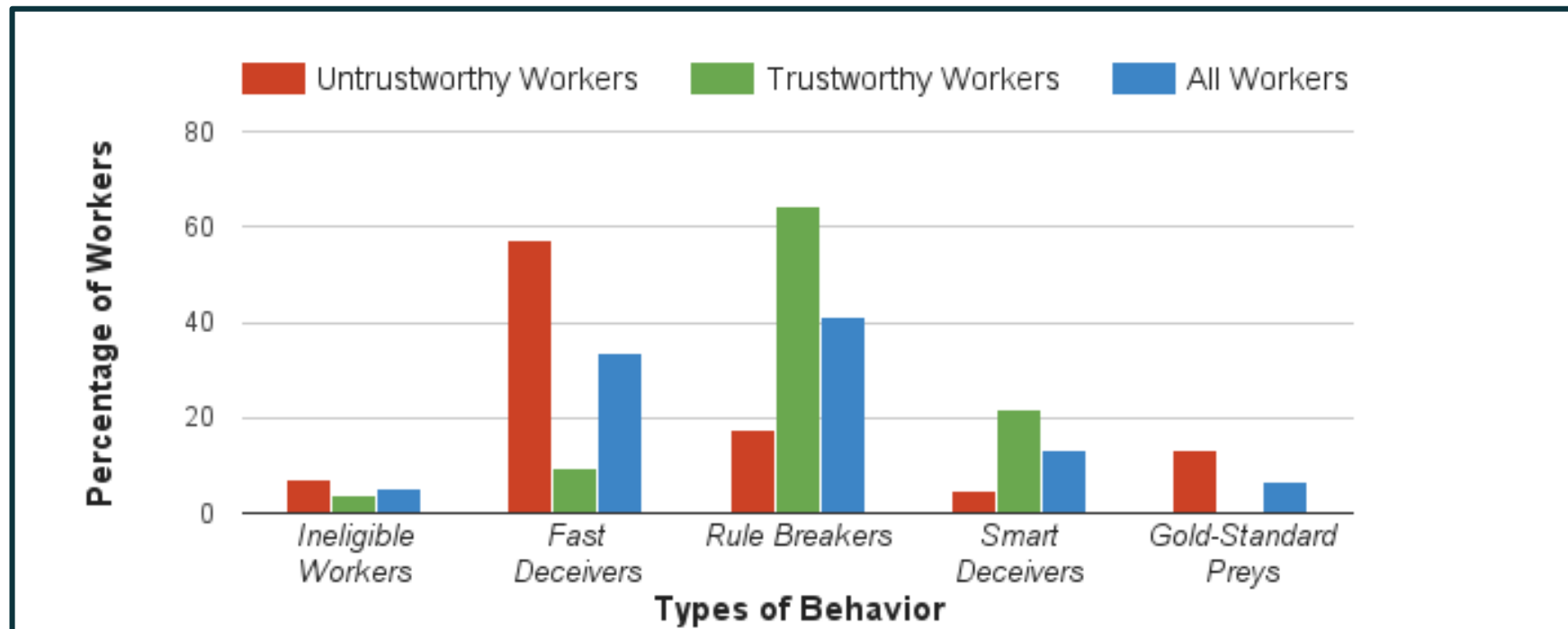
Response: *'one, two, three, four, five'*

Gold Standard  
Preys (GSP)

These workers abide by the instructions and provide valid responses, but stumble at the gold-standard questions!

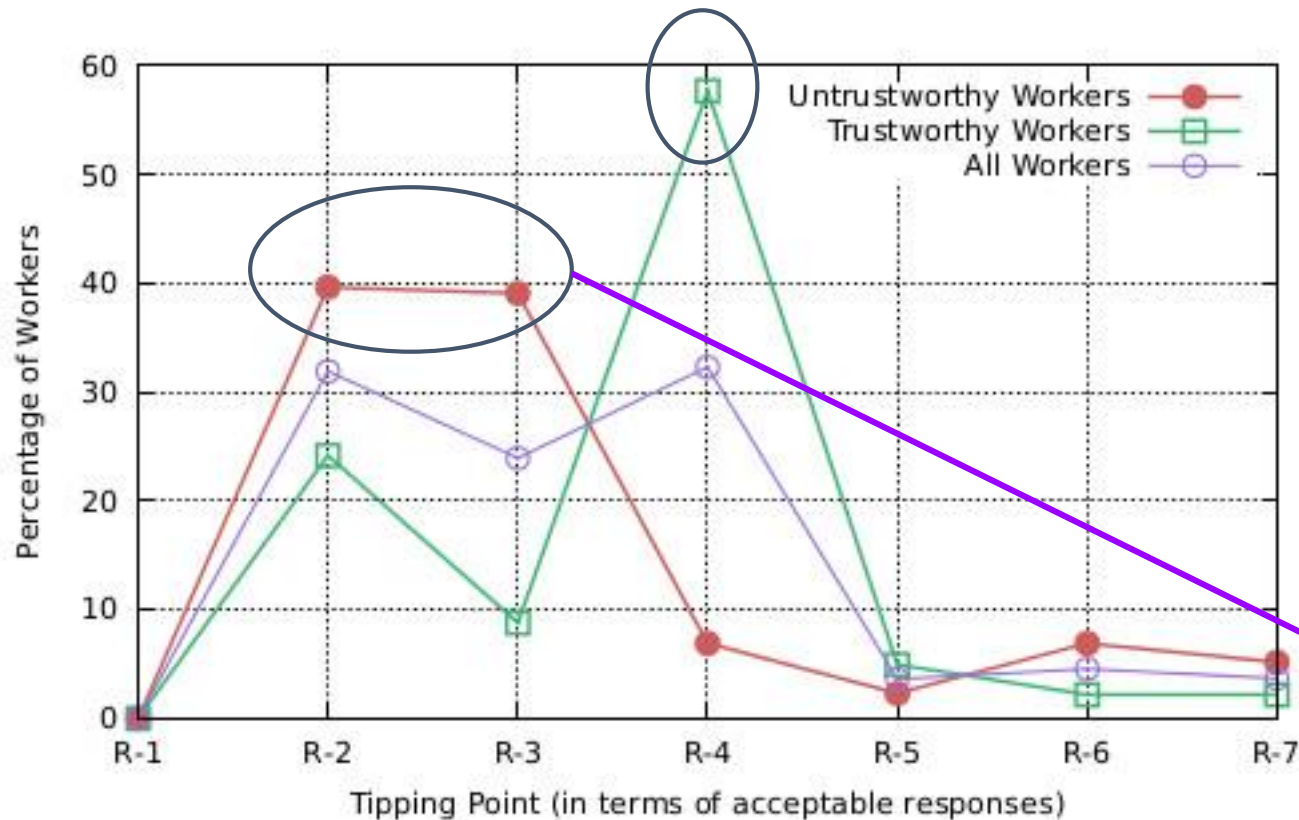
# RQ2 - Distribution of Low-quality Workers

- passed the gold-standard: **Trustworthy workers (TW)**
- failed to pass the gold-standard: **Untrustworthy workers (UW)**



# Tipping Point

- “the first point at which a worker begins to exhibit malicious behavior after having provided an acceptable response”



**Table 1. Relationship between the Maliciousness and Tipping Point of untrustworthy and trustworthy workers (percentage of workers having tipping point @ R).**

Maliciousness	UW	TW
$0 < M \leq 0.2$	40.9% @ R-7 31.8% @ R-6	28.5% @ R-7 28.5% @ R-5
$0.2 < M \leq 0.4$	43.47% @ R-3 21.73% @ R-6	30% @ R-5 30% @ R-3
$0.4 < M \leq 0.6$	66.19% @ R-3 25.35% @ R-2	88% @ R-4 5.1% @ R-3
$0.6 < M \leq 0.8$	71.05% @ R-2 28.95% @ R-3	60% @ R-3 40% @ R-2
$0.8 < M \leq 1$	100% @ R-2	100% @ R-2

# Findings

- Identified different types of malicious behavior exhibited by crowd workers.
- Measuring ‘maliciousness’ of workers to quantify their **behavioral traits**, and ‘**tipping point**’ to further understand worker behavior.
- This understanding helps requesters in effective task design, ensures adequate utilization of the crowdsourcing platform(s).
- Guidelines for efficient design of Surveys by limiting malicious activity.
  - Pre-screening (ineligible)
  - Validators (fast deceivers, rule breaker)
  - Psychometric approaches (smart deceivers)

# Modus Operandi of Crowd Workers

Ujwal Gadiraju, Alessandro Checco, Neha Gupta, and Gianluca Demartini. Modus Operandi of Crowd Workers: The Invisible Role of Microtask Work Environments. In: Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies (IMWUT) presented at The ACM International Joint Conference on Pervasive and Ubiquitous Computing (**UBICOMP 2017**). Maui, Hawaii, September 2017.

# Context

- Crowd workers are embedded in diverse work environments
- Work environment: hardware/software at disposal
- Usually requesters provide an undifferentiated task to all workers
- How do task UI elements and work environments interact?

# Studies

- **Study I** - Survey on 100 people with questions about experience and problems related to UI
  - Problems with input (text areas, checkboxes, radio buttons), multimedia (audio,video), links, colors, buttons
- **Study II** – Measured performances of task design variants
  - 43 synthetic variations x 3 tasks x 50 judgements x 2 countries = 12 900 resp
  - American workers were faster than Indian workers
  - American workers outperformed Indian workers in audio transcription tasks (coping well with poor quality audio as well)
  - Workers with faster devices (laptops were found to be faster than desktops) provided higher quality responses (more tags, more unique tags)



# Studies

- **Study III** – 1:1 interviews with workers who participated to study II
  - Different devices are used for different tasks
  - Internet speed and cost is a variable for task selection (e.g., multimedia content)

*“Sometimes the Internet fee is greater than the rewards I earn (due) to images, audios or videos in tasks.”*

– CrowdFlower Worker from India

- Language proficiency has great impact on accuracy
- ModOp: a tool to check for crowdsourcing task design problems

# Conclusions

- Paid micro-task crowdsourcing to build hybrid human-machine systems
- Human-in-the-loop systems means to consider human factors!
- Timeouts to increase efficiency and effectiveness of crowd work
  - Does it generalize to other task types?
- Malicious behaviors
  - Supervised worker type classification
- The effect of work environment on work efficiency and effectiveness
  - Build recommender systems / assign tasks based on complexity

<http://gianlucademartini.net>

@eglu81