# CRQA: Crowd-powered Real-time Automated Question Answering System
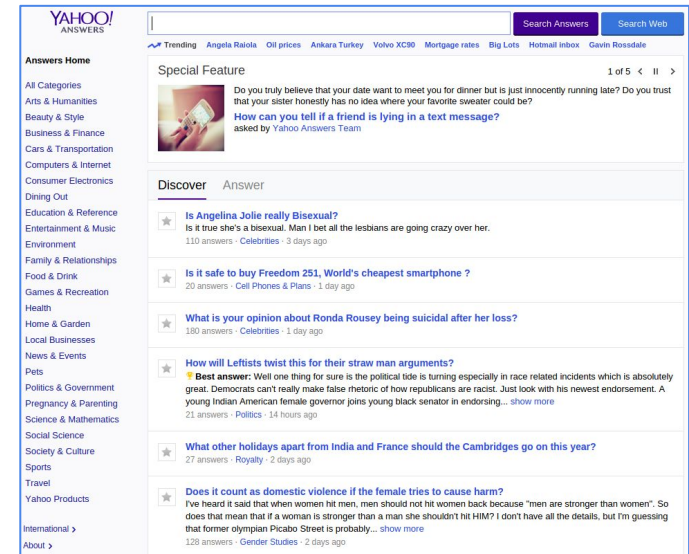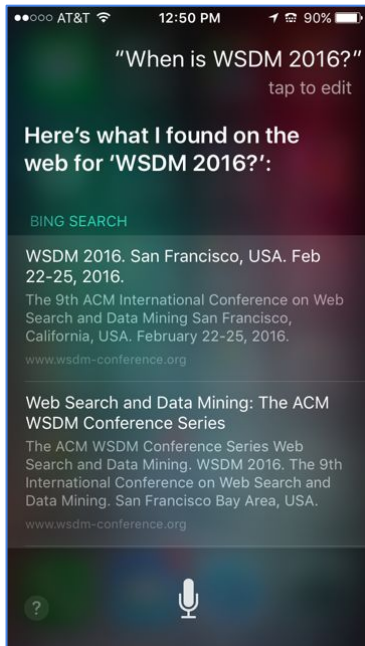
### Denis Savenkov
Emory University
dsavenk@emory.edu

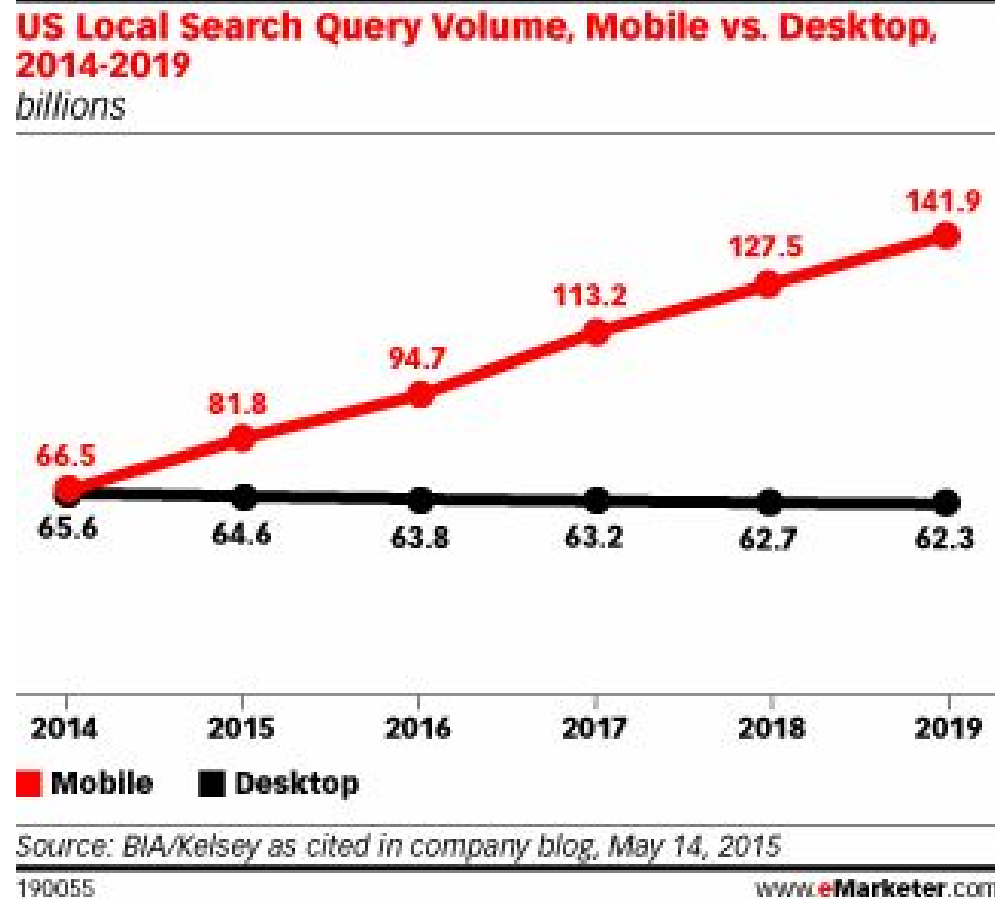### Eugene Agichtein
Emory University
eugene@mathcs.emory.edu

HCOMP, Austin, TX
October 31, 2016

# Volume of question search queries is growing[1]



[1] "Questions vs. Queries in Informational Search Tasks", Ryen W. White et al, WWW 2015

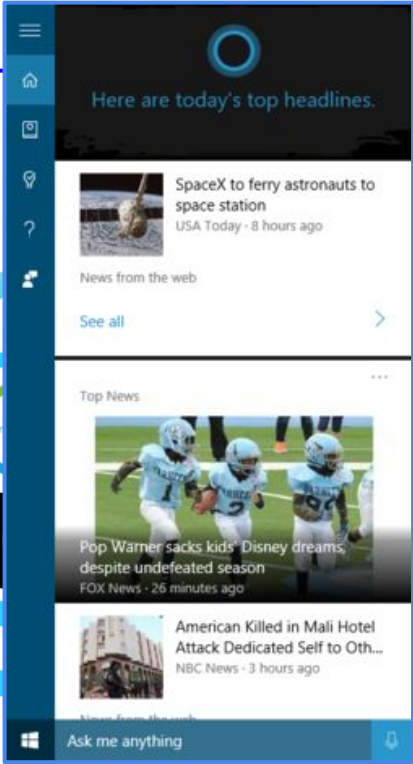# And more and more of this searches are happening on mobile

**US Local Search Query Volume, Mobile vs. Desktop, 2014-2019**

*billions*

Mobile: 66.5 (2014), 81.8 (2015), 94.7 (2016), 113.2 (2017), 127.5 (2018), 141.9 (2019)

Desktop: 65.6 (2014), 64.6 (2015), 63.8 (2016), 63.2 (2017), 62.7 (2018), 62.3 (2019)

■ Mobile   ■ Desktop

Source: BIA/Kelsey as cited in company blog, May 14, 2015

190055                                                    www.eMarketer.com

# Mobile Personal Assistants are popular

# Automatic Question Answering works relatively well for some questions



(AP Photo/Jeopardy Productions, Inc.)

# ... but not sufficiently well for many other questions

… when there is no answer, digging into "10 blue links" is even harder on mobile devices

It is important to improve question answering for complex user information needs

# Goal of **TREC LiveQA** shared task is to advance research into **answering real user questions in real time**

category: Pets > Cats

title: Why does my cats behavior change so diradticly at night?

body: So I have two litter mate cats, over one year old and have been neutered for more than six months. one of the cats has the same attitude all day and just gets more lovey at night. The other one is super cute and friendly during the day, and loves attention from me. But when night rolls around and I m about to get ready for bed, he avoids my touch. Walks away when I walk up and even one time when I picked him up, he hissed at me. It s not like he s mad at me for going to bed Bc he doesn t do this to my mother or sister. The other thing he does is when his brother is sleeping with me, he ll come in, wake him up, and then leave my room with his brother following him. Does anything describe this behavior? Why does this happen?

YAHOO! ANSWERS

Question Answering System

1 minute

24 hours

**Answers**

≤ 1000 chars

# LiveQA Evaluation Setup

Answers are pooled and judged by NIST assessors

- o **1: Bad** - contains no useful information
- o **2: Fair** - marginally useful information
- o **3: Good** - partially answers the question
- o **4: Excellent** - fully answers the question

# LiveQA 2015:

# Even the best system returns a fair or better answer only for ~50% of the questions!

|  | Avg score (0-3) | % questions with fair or better answer | % questions with excellent answer |
|---|---|---|---|
| Best system | 1.08 | 53.2 | 19.0 |

# The architecture of baseline automatic QA system



1. Search data sources
   a. CQA archives
      i. Yahoo! Answers
      ii. Answers.com
      iii. WikiHow
   b. Web search API
2. Extract candidates and their context
   a. Answers to retrieved questions
   b. Content blocks from regular web pages
3. Represent candidate answers with a set of features
4. Rank them using LambdaMART model
5. Return the top candidate as the answer

# Common Problem: Automatic systems often return an answer about the same topic, but irrelevant to the question

Pets > Rodents

## Scatter feeding my hamster?

I always leave 1/4th to 1/2 of her bowl filled at all times. Every night though i like to scatter food all around her cage to give her something to do,underneath toys,inside toys ect. She seems to enjoy seeking it all out but then immediately buries it...is she eating it later? Her bowl does lose alittte food from... show more

☆ Follow   ⁎ 2 answers

Throwback to when my friends hamster ate my hamster and then my friends hamster died because she forgot to feed it karma

# Incorporate crowdsourcing to assist an automatic real-time question answering system

Or: combine human insight and automatic QA with machine learning

# Existing research

✓ **"Direct answers for search queries in the long tail"** by M.Bernstein et al, 2012
  ○ Offline crowdsourcing of answers for long-tail search queries

✓ **"CrowdDB: answering queries with crowdsourcing"** by M.Franklin et al, 2011
  ○ Using crowd to perform complex operations in SQL queries

✓ **"Answering search queries with crowdsearcher"** by A.Bozzon et al, 2012
  ○ Answering queries using social media

✓ **"Dialog system using real-time crowdsourcing and twitter large-scale corpus"** by F. Bessho et al, 2012
  ○ Real-time crowdsourcing as a backup plan for dialog

✓ **"Chorus: A crowd-powered conversational assistant"** by W.Lasecki, 2013
  ○ Real-time chatbot powered by crowdsourcing

... and many other works

# Research Questions

- ○ **RQ1.** <u>Can crowdsourcing be used</u> to improve the performance of a near <u>real-time automatic question answering system</u>?

# Research Questions

- ○ **RQ1.** <u>Can crowdsourcing be used</u> to improve the performance of a near real-time automatic question answering system?

- ○ **RQ2**. <u>What kind of contributions</u> from crowd workers can help improve automatic question answering and what is the relative impact of different types of feedback to the overall question answering performance?

# Research Questions

- ○ **RQ1.** <u>Can crowdsourcing be used</u> to improve the performance of a near real-time automatic question answering system?

- ○ **RQ2**. <u>What kind of contributions</u> from crowd workers can help improve automatic question answering and what is the relative impact of different types of feedback to the overall question answering performance?

- ○ **RQ3**. What are the trade-offs in <u>performance, cost, and scalability</u> of using crowdsourcing for real-time question answering?

# CRQA: Integrating crowdsourcing with automatic QA system



1. After receiving a question, it is forwarded to the crowd
2. Can start working on the answer, if possible
3. When system ranks candidates, top-7 are pushed to workers for rating
4. Rated human and automatically generated answers are returned
5. System re-rank them based on all available information
6. Top candidate is returned as the answer

# We used the retainer model for real-time crowdsourcing

tasks

Our crowdsourcing UI

15 mins

$

labels

# UI for crowdsourcing answers and ratings

# Heuristic answer re-ranking (during TREC LiveQA)

# CRQA uses a learning-to-rank model to re-rank

Answer candidate    Answer candidate    Answer candidate    Answer candidate

**> sort answers -k crowd_rating**

**if** top candidate rating > 2.5 or no crowd generated candidates

**False**                                      **True**

**return longest crowd generated candidate**                        **return top candidate**

# CRQA uses a learning-to-rank model to re-rank

Answer candidate

Answer candidate

Answer candidate

Answer candidate

Answer re-ranking model features:
- answer source
- initial rank/score
- # crowd ratings
- min, median, mean, max crowd rating

final answer

- Offline crowdsourcing to get ground-truth labels

- Included Yahoo!Answers community response, crawled 2 days after challenge

- Trained GBRT model, 10-fold cross validation

# Evaluation

# Evaluation setup

**Methods compared**:

➢ Automatic QA

➢ CRQA (heuristic): re-ranking by crowdsourced score

➢ CRQA (LTR): re-ranking using a learning-to-rank model

➢ Yahoo! Answers (crawled 2 days later)

**Metrics**:

➢ avg-score: average answer score over all questions

➢ avg-prec: average answer score

➢ success@i+: fraction of questions with answer score $\geq i$

➢ precision@i+: fraction of answers with score $\geq i$

# Dataset

➢ 1,088 questions from LiveQA 2016 run
➢ Top 7 system and crowd-generated answers
➢ Answer quality labelling on a scale from 1 to 4
  - offline
  - also using crowdsourcing (different workers)

| | |
|---|---:|
| Number of questions received | 1,088 |
| Number of MTurk 15 minutes assignments completed | 889 |
| Average number of questions per assignment | 11.44 |
| Total cost per question | $0.81 |
| Avg number of answers provided by workers per question | 1.25 |
| Average number of ratings per answer | 6.25 |

# Main Results

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA: (heuristic) | 2.416 | 2.421 | 0.75 | 0.32 | 0.03 | 0.75 | 0.32 | 0.03 |
| CRQA (LTR) | **2.550** | **2.556** | **0.80** | **0.40** | 0.03 | **0.80** | 0.40 | 0.03 |
| Yahoo! Answers | 2.229 | 2.503 | 0.66 | 0.37 | **0.04** | 0.74 | **0.42** | **0.05** |

# Crowdsourcing improves performance of automatic QA system

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA: (heuristic) | 2.416 | 2.421 | 0.75 | 0.32 | 0.03 | 0.75 | 0.32 | 0.03 |
| CRQA (LTR) | **2.550** | **2.556** | **0.80** | **0.40** | 0.03 | **0.80** | 0.40 | 0.03 |
| Yahoo! Answers | 2.229 | 2.503 | 0.66 | 0.37 | **0.04** | 0.74 | **0.42** | **0.05** |

# Learning-to-rank model allows to more effectively combine all available signals and return a better answer

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA: (heuristic) | 2.416 | 2.421 | 0.75 | 0.32 | 0.03 | 0.75 | 0.32 | 0.03 |
| CRQA (LTR) | **2.550** | **2.556** | **0.80** | **0.40** | 0.03 | **0.80** | 0.40 | 0.03 |
| Yahoo! Answers | 2.229 | 2.503 | 0.66 | 0.37 | **0.04** | 0.74 | **0.42** | **0.05** |

# CRQA reaches the quality of community responses on Yahoo! Answers

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA: (heuristic) | 2.416 | 2.421 | 0.75 | 0.32 | 0.03 | 0.75 | 0.32 | 0.03 |
| CRQA (LTR) | **2.550** | **2.556** | **0.80** | **0.40** | 0.03 | **0.80** | 0.40 | 0.03 |
| Yahoo! Answers | 2.229 | 2.503 | 0.66 | 0.37 | **0.04** | 0.74 | **0.42** | **0.05** |

# ... and it has much better coverage

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA: (heuristic) | 2.416 | 2.421 | 0.75 | 0.32 | 0.03 | 0.75 | 0.32 | 0.03 |
| CRQA (LTR) | **2.550** | **2.556** | **0.80** | **0.40** | 0.03 | **0.80** | 0.40 | 0.03 |
| Yahoo! Answers | 2.229 | 2.503 | 0.66 | 0.37 | **0.04** | 0.74 | **0.42** | **0.05** |

# Both worker answers and ratings make an equal contribution to the answer quality improvements

| Method | avg-score | avg-prec | s@2+ | s@3+ | s@4+ | p@2+ | p@3+ | p@4+ |
|---|---|---|---|---|---|---|---|---|
| Automatic QA | 2.321 | 2.357 | 0.69 | 0.30 | 0.02 | 0.71 | 0.30 | 0.03 |
| CRQA (LTR) | 2.550 | 2.556 | 0.80 | 0.40 | 0.03 | 0.80 | 0.40 | 0.03 |
| no worker answers | 2.432 | 2.470 | 0.75 | 0.35 | 0.03 | 0.76 | 0.35 | 0.03 |
| no worker ratings | 2.459 | 2.463 | 0.76 | 0.35 | 0.03 | 0.76 | 0.36 | 0.03 |

# Crowdsourcing helps to improve empty and low quality answers

# Yahoo! Answers have both higher percentage of excellent and missing and low quality answers



Many questions on Yahoo! Answers are unanswered

Community experts provide an "excellent" answer more often than CRQA

Legend: CRQA (green), Yahoo! Answers (red)

Answer score

# Crowdsourced answers are especially good for general knowledge questions

## Is it bad not wanting to visit your family?

My mom and stepdad are planning to go to New York at the beginning of July, and I don't want to go because I really don't feel like it , they're only going to be in New York for 4 days . I'm 18 now , but I'd rather be home. We live in Georgia by the way. Is it bad or is it okay

☆ Follow    ✻ 3 answers

Is it bad not wanting to visit your family? It's nt bad. Just be honest with them. They may be upset but they should understand
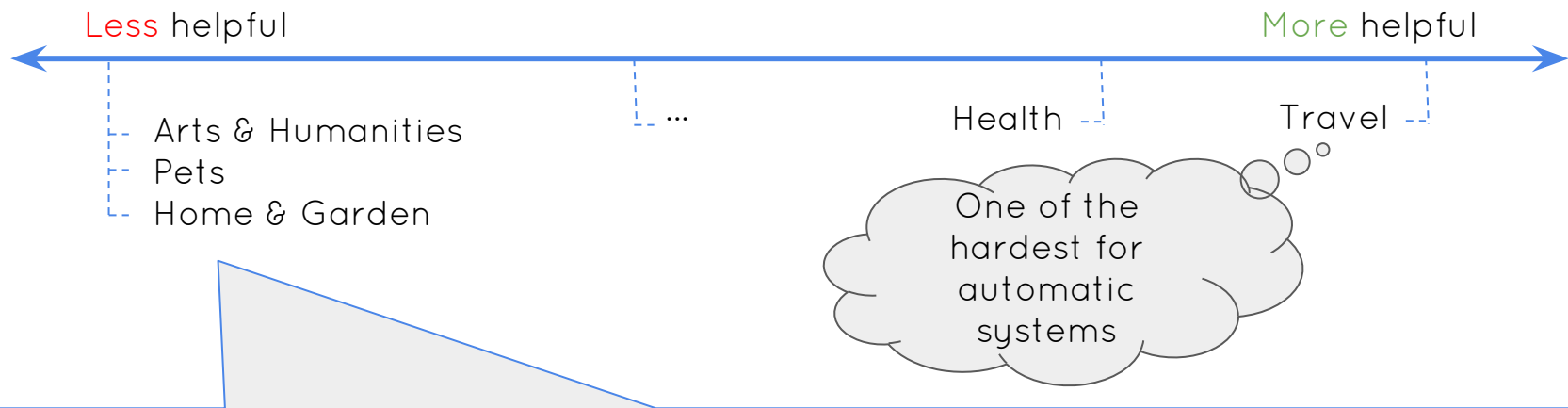
## What is a good remedy/medicine for stomach aches? Specifically ones caused by stress or anxiety.?

☆ Follow    ✻ 2 answers

Chamomile tea should help

36

# But less effective for questions which require domain expertise

Less helpful ← ———————————————————————————— → More helpful

Arts & Humanities
Pets
Home & Garden

...

Health

Travel

One of the hardest for automatic systems

---

Arts & Humanities > History

**What was pol pot's vision of an agrarian society?**

☆ Follow    ✻ 1 answer

---

Health > General Health Care > Injuries                                        Next >

**What is the average time for recovery of a complete Scapholunate (scaphoid) ligament tear in my wrist?**

I injured my wrist while at work. I was taken to the ER that night and diagnosed with a scaphoid dislocation. I saw a hand surgeon two days later, and had surgery the next week (all within about 10 days of the initial injury). The scope showed a complete tear of the ligament, at which time he performed the... show more
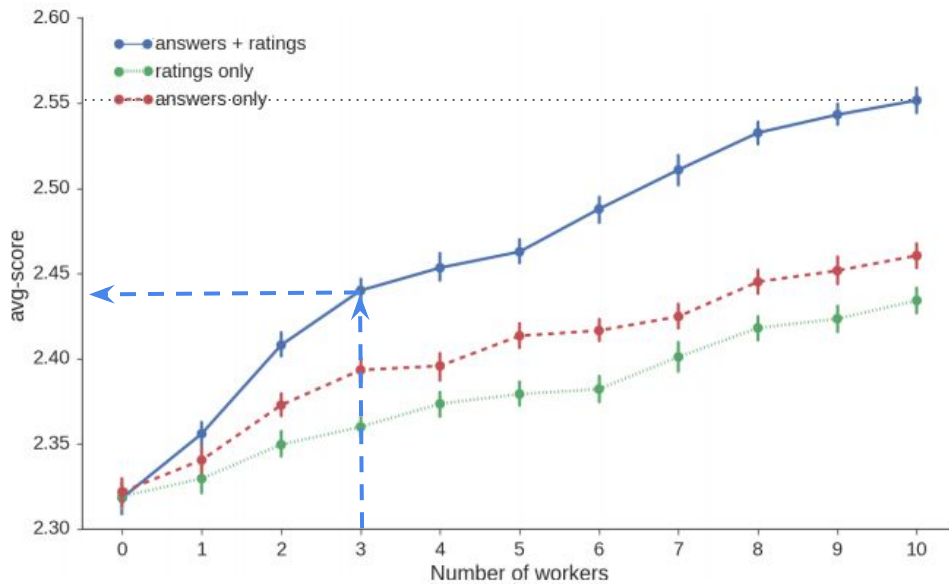
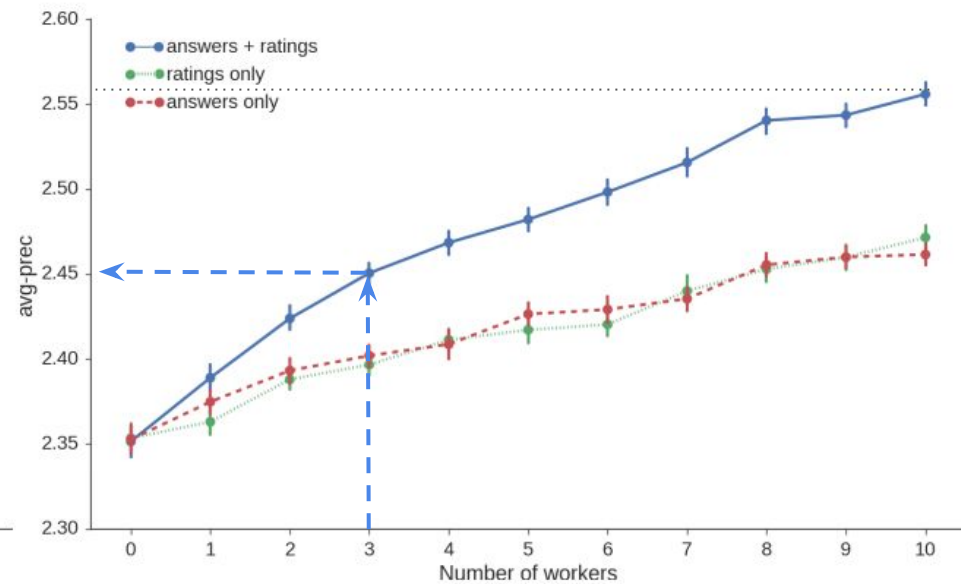☆ 1 following    ✻ 2 answers

Ok, but what about the costs?

$0.81 per question is a lot of money

# Half of the overall improvements can be achieved with only 3 workers per question (30% of cost)



(a) avg-score: Average score per question

(b) avg-prec: Average score per answer (ignoring non-answered questions)

# Limitations and Future work

➢ Limitations
   ○ Fixed and uniform load for the system over 24 hours
      - Need variable size pool of workers based on the current load
➢ Ideas
   ○ Allocate crowdsourcing resources based on expected performance of the automated system
   ○ Use other types of feedback:
      - Search query generation
      - Key phrases to look for the the answer
      - …
   ○ Online learning from crowd feedback
   ○ Cost optimization
      - Decide which feedback, in what amount and when to request

We conducted large scale experiments on real user questions, which showed:

- Crowdsourcing helps for real-time QA
  - Workers can contribute answers and rate candidates
  - Humans can immediately reject off-topic candidates
- Answers from our system are often even preferred to community answers
  - Which are collected 2 days after
  - With 20% of the questions were still unanswered by the community
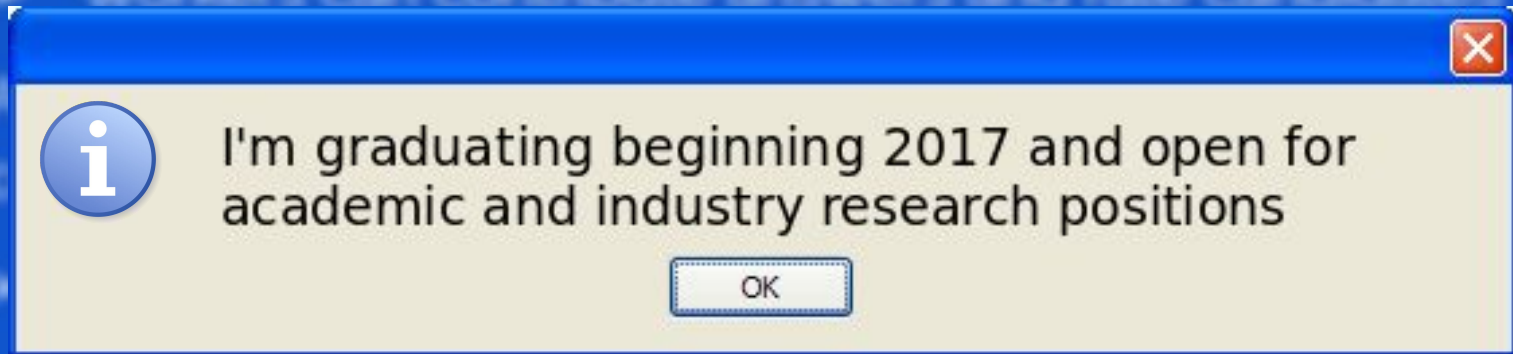
## Thank you!

We conducted large scale experiments on real user questions, which showed:

○ Crowdsourcing helps for real-time QA

  ► Workers can contribute answers and rate candidates

○ A

  ► With 20% of the questions still unanswered

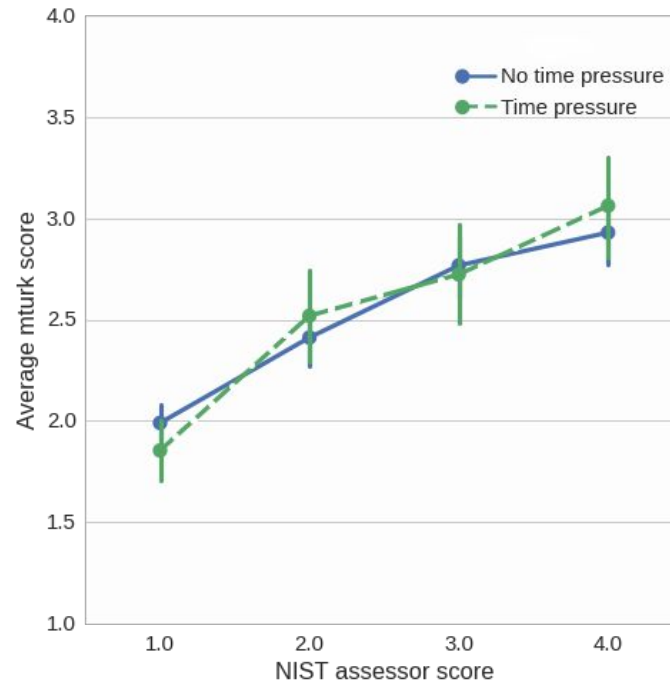**I'm graduating beginning 2017 and open for academic and industry research positions**

[ OK ]

Thank you!

# It's better to present candidates ordered by their predicted quality

Average answer score if presented in different order:

| | |
|---|---|
| Sorted by rank | 2.539 |
| Shuffled | 2.508 |

# [Backup] Crowdsourced labels correlate well with NIST assessor scores ($\rho$=0.52)



✓ Workers prefer to give intermediate scores (2, 3), while NIST assessors gave more extreme scores (1 and 4)

✓ There is no significant difference in quality between groups with and without time pressure

# Features [backup]

| **Answer statistics** |
|---|
| — Length in chars, words and sentences |
| — Average number of words per sentence |
| — Fraction of non-alphanumeric characters |
| — Number of question marks |
| — Number of verbs |
| **Answer source** |
| — Binary feature for each of the search verticals: Web, Yahoo! Answers, Answers.com, WikiHow.com |
| **N-gram matches** |
| — Cosine similarities using uni-, bi- and tri-gram representations of the question title and/or body, and answer text, topic or context |
| — The lengths of longest spans of matched terms between question title and/or body, and answer text, topic or context |
| **Information Retrieval score** |
| — BM25 scores between question title and/or body, and answer text, topic or context |