

Using Crowdsourcing for Labelling Emotional Speech Assets

Alexey Tarasov, Charlie Cullen, Sarah Jane Delany

Digital Media Centre
Dublin Institute of Technology



Project Introduction

- ◆ Science Foundation Ireland funded project
- ◆ Objective:
 - ◆ prediction of levels of emotion in natural speech
- ◆ 2 strands:
 - ◆ acoustic analysis (Dr Charlie Cullen - DIT MmIG)
 - ◆ machine learning prediction (Dr Sarah Jane Delany - DIT AIG)
- ◆ 4 year project, started in October 2009
 - ◆ 2 PhD students

Requirements for Supervised Learning

- ◆ Performance of supervised learning techniques depends on the quality of the training data
- ◆ Requirements:
 - ◆ High quality speech assets
 - ◆ Good labels

Starting Point...

- ◆ **Emotional speech corpus** [Cullen et al. LREC 08]
 - ◆ natural assets
 - ◆ use of *Mood Induction Procedures*
 - ◆ high quality recording
 - ◆ participants recorded in separate sound isolation booths
 - ◆ contextual or meta data is recorded where available
 - ◆ based on *IMDI* annotation schema

Next Steps...

- ◆ Need to rate these assets...
- ◆ Challenges:
 - ◆ manual annotation can be expensive and time consuming
 - ◆ experts often disagree
 - ◆ expertise does not necessarily correlate with experience

Consider Crowdsourcing?

Crowdsourcing

“The act of taking a task traditionally performed by a designated agent and outsourcing it to an undefined, generally large group of people in the form of an open call” [Jeff Howe]

Crowdsourcing

- ◆ June 2006 Wired magazine article by Jeff Howe

...the power of many...

www.wired.com/wired/archive/14.06/crowds.html

ADVERTISEMENT

Making the smartphone brilliant.  Samsung E... GALAXY...

SEE FOR YOURSELF

WIRED SUBSCRIBE >> SECTIONS >> BLOGS >> REVIEWS >> VIDEO >> HOW-TO >>

Sign In | RSS Feeds

Issue: June 2006 - June 2006
Subscribe to WIRED magazine and receive a FREE...

The Rise of Crowdsourcing

Remember outsourcing? Sending jobs to India and China is so 2003. The new pool of cheap labor: everyday people using their spare cycles to create content, solve problems, even do corporate R & D.

By Jeff Howe Page 1 of 4 [next >>](#)

1. The Professional

Story Tools

[PRINT](#) [MAIL](#)

Story Images



Feature:
The Rise of Crowdsourcing

Plus:
[5 Rules of the New Labor Pool](#)
[Look Who's Crowdsourcing](#)

Claudia Menashe needed pictures of sick people. A project director at the National Health Museum in Washington, DC, Menashe was putting together a series of interactive kiosks devoted to potential pandemics like the avian flu. An exhibition designer had created a plan for the kiosk itself, but now Menashe was looking for images to accompany the text. Rather than hire a photographer to take shots of people suffering from the flu, Menashe decided to use preexisting images – stock photography, as it's known in the publishing industry.

Mechanical Turk is a marketplace for work.
We give businesses and developers access to an on-demand, scalable workforce. Workers select from thousands of tasks and work whenever it's convenient.
86,131 HITS available. [View them now.](#)

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

Find an interesting task → **Work** → **Earn money**

[Find HITS Now](#)

or [learn more about being a Worker](#)

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

Fund your account → **Load your tasks** → **Get results**

[Get Started](#)

<https://www.mturk.com/mturk/>



reCAPTCHA™

- WHAT IS reCAPTCHA
- GET reCAPTCHA
- PROTECT YOUR EMAIL
- MY ACCOUNT
- RESOURCES: DOCS & PLUGINS

reCAPTCHA IS A FREE ANTI-BOT SERVICE THAT HELPS DIGITIZE BOOKS.

steamboat train, from New
this **morning** ran off the track.
New-London. Four cars plunged



→ LEARN HOW reCAPTCHA WORKS

USE reCAPTCHA ON YOUR SITE

-  **STRONG SECURITY**
-  **ACCESSIBLE TO BLIND USERS**
-  **30+ MILLION SERVED DAILY**

NEW See how accurate reCAPTCHA is at digitizing content!

[Blog](#) | [About Us](#) | [Contact](#) | [FAQs](#) | [Terms](#) | [Privacy](#)
© 2010 Google, all rights reserved.

The screenshot shows the Gwap website interface. At the top, there is a navigation bar with the Gwap logo and several game categories: ESP Game, Tag a Tune, Verbose, Squigl, Matchin, FlipIt, and PopVideo. Below the navigation bar is a login section with input fields for email and password, a 'Sign In' button, a 'remember me' checkbox, and a 'forgot password?' link.

The main content area features a large banner with the text 'Play the Games, Change the Web.' and 'When you play a game at Gwap, you aren't just having fun.' Below this are 'Learn More' and 'Register' buttons. To the right, there is a spotlight effect on the 'Verbose' game, with the text 'Verbose it's common sense.' and 'It has wheels... It's bigger than a car. Quick! Guess the word!' Below this is a screenshot of the game interface showing a clue 'it contains instruments.' and a 'band?' question, with a 'PLAY NOW' button.

At the bottom of the main content area, there is a 'Top 10' leaderboard with the following entries:

Rank	Player	Score
1	PlasticBiddy	251 K
2	mathman	236 K
3	brasso	186 K
4	Lazer	163 K
5	Jeff	160 K

At the bottom of the page, there are links for 'Blog', 'About', and 'Contact'.

The screenshot shows the ESP Game website interface. At the top, there is a navigation bar with several colored buttons: 'gwap', 'ESP Game', 'Tag a Tune', 'Verbosity', 'Squigl', 'Matchin', 'FlipIt', and 'PopVideo'. Below the navigation bar, the main header features the 'ESP Game' logo with a starburst icon and the tagline 'Concentrate...'. The main content area is green and titled 'How to Play'. It contains two numbered steps: Step 1: 'You and a partner see the same image.' with a small image of a tree. Step 2: 'Each of you must guess what words your partner is typing.' with a text input field containing 'Tree' and a 'make a new' button. To the right of the steps are two large green buttons: 'Got it, Let's Play!' and 'View Top Scores'. The background of the main content area is decorated with a starburst pattern of white stars and green lines.

gwap ESP Game Tag a Tune Verbosity Squigl Matchin FlipIt PopVideo

ESP Game
Concentrate...

How to Play

- 1 You and a partner see the same image.
- 2 Each of you must guess what words your partner is typing.

make a new
Tree

Got it, Let's Play!

View Top Scores

Crowdsourcing

- ◆ Triggered a shift in the way labels or ratings are obtained in variety of domains:
 - ◆ natural language tasks [Snow et al. 2008]
 - ◆ computer vision [Sorokin & Forsyth 2008, vonAhn & Dabbish 2004]
 - ◆ sentiment analysis [Hsueh et al. 2008, Brew et al. 2010]
 - ◆ machine translation [Ambati et al. 2010]

Practical Experiences

◆ Speed

- ◆ 300 annotations from each of 10 annotators in < 11 mins [Snow et al. 2008]
- ◆ evidence that obtaining 'quality' annotations effects time (avg completion time 4 mins vs 1.5 mins) [Kittur et al. 2008]

Practical Experiences

◆ Quality

- ◆ 875 expert-equivalent affect labels per \$1

[Snow et al. 2008]

- ◆ by identifying 'good' annotators accurate labels can be achieved with significant reduction in effort

[Donmez et al. 2008, Brew et al. 2010]

Challenges

- ◆ How to
 - ◆ select which assets are presented for rating?
 - ◆ estimate the reliability of the annotators?
 - ◆ ensure the reliability of the ratings?
 - ◆ select training data for the prediction systems?
 - ◆ maintain the balance between consensus and data coverage?

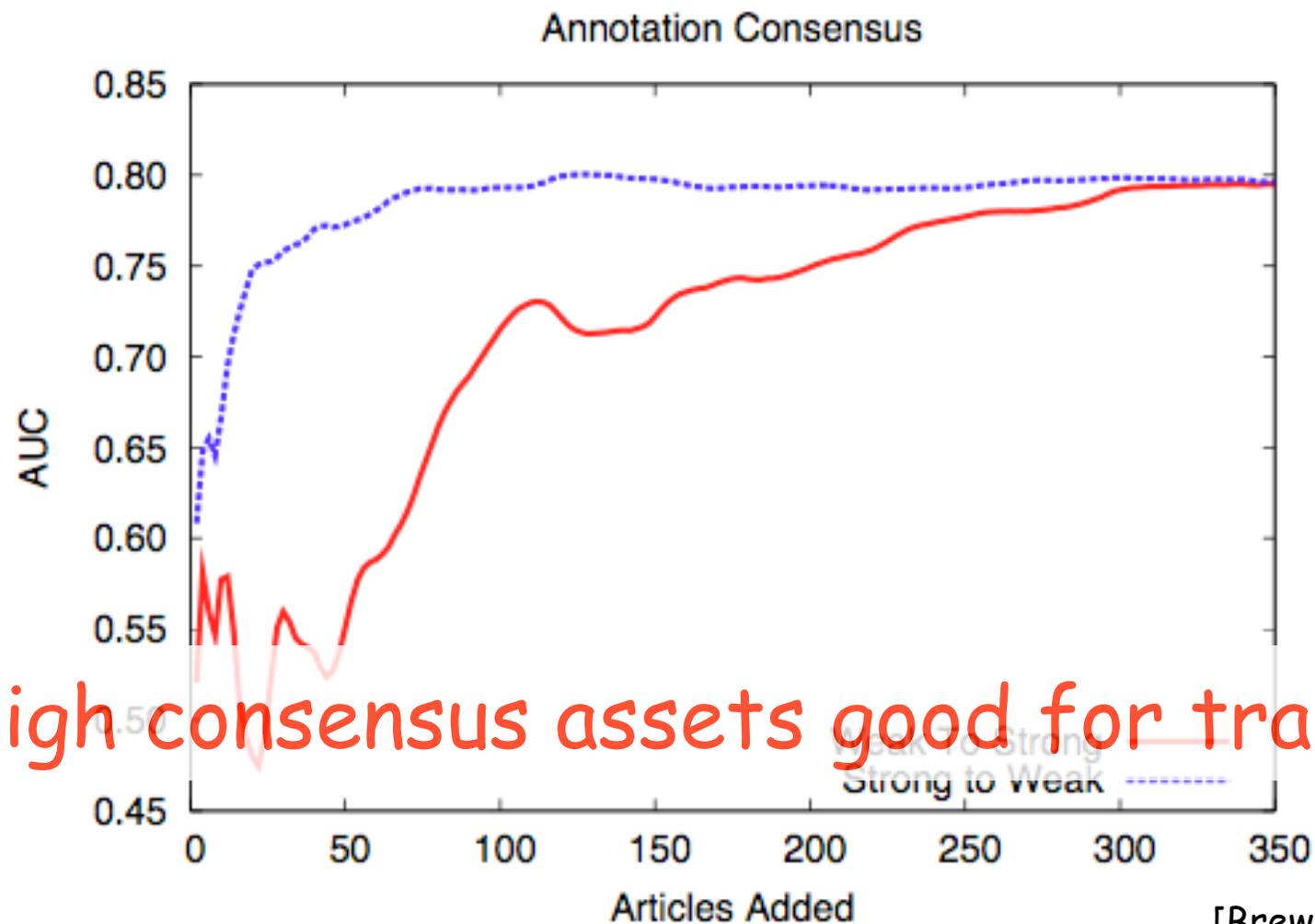
Asset Selection

- ◆ **Active Learning** used by [Ambati et al. 2010, Domnez et al. 2009]
 - ◆ a supervised learning technique which selects the most informative examples for annotation
- ◆ **Clustering** used by [Brew et al. 2010]
 - ◆ grouping examples and selecting representative examples from cluster to annotate

Annotator Reliability

- ◆ Depends on whether annotators are identifiable or not...
- ◆ Strategies for recognising strong annotators
 - ◆ 'Good' Annotators those that 'agree' with the consensus rating [Brew et al. 2010]
 - ◆ Iterative approach to filter out weaker annotators [Domnez et al. 2009]

Good Annotators are Useful...



high consensus assets good for training...

[Brew et al. 2010]

Deriving ratings

- ◆ Use consensus rating [Brew et al. 2010]
 - ◆ select the rating with highest consensus
 - ◆ thresholds can apply
- ◆ Only use good annotators to derive rating [Domnez et al. 2009]
- ◆ Using learning techniques to estimate 'ground truth' from multiple noisy labels [Smyth et al. 1995, Raykar et al. 2009/10]

Consensus vs. Coverage

- ◆ Is it better to label more assets or get more labels per asset?
 - ◆ Research suggests fewer annotations are needed in domains with high consensus [Brew et al. 2010]

Reliability of the Ratings

- ◆ Evidence of 'gaming' with crowdsourcing services
 - ◆ numbers of untrustworthy users is not large
- ◆ Techniques
 - ◆ require users to complete a test first [Ambiati et al. 2010]
 - ◆ use percentage of previously accepted submissions [Hsueh et al. 2008]
 - ◆ include explicitly verifiable questions [Kittur et al. 2008]

Use Case

"Seán has a set of speech assets extracted from recordings of experiments using mood induction procedures.

He wants to get these assets rated on a number of different scales, including activation and evaluation, by a large number of non-expert annotators.

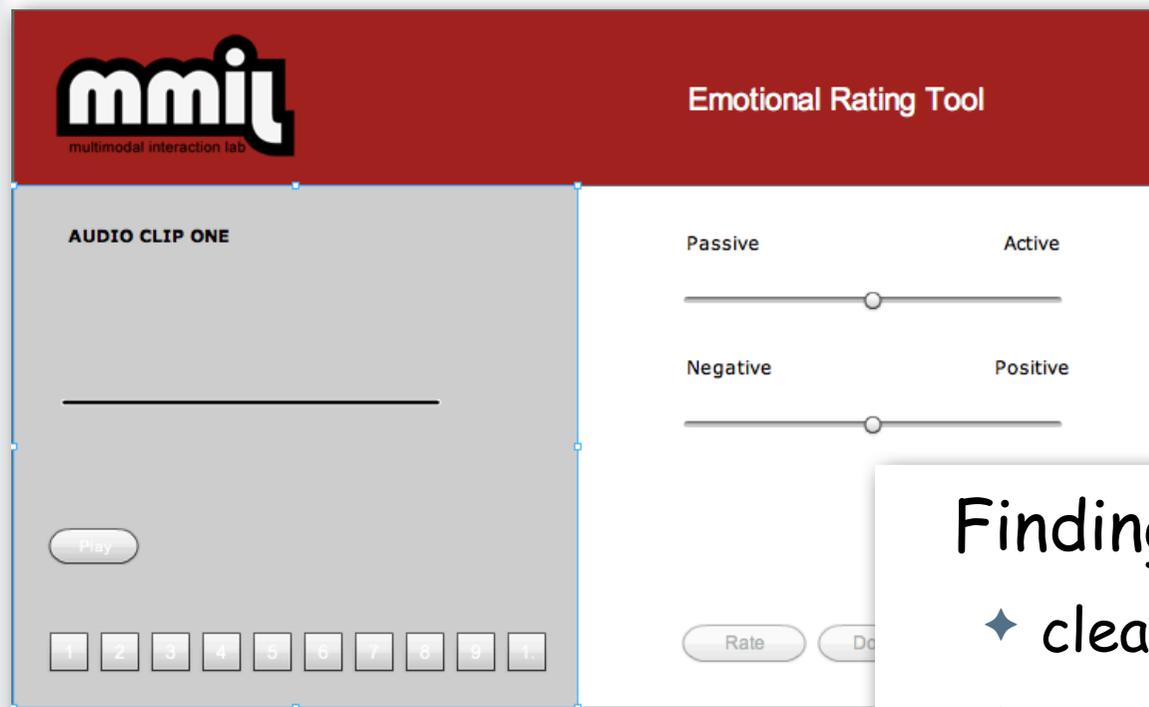
He wants to use a micro-task system such as Mechanical Turk to get these ratings.

Active learning will be used to select the most appropriate assets to present for labels from the annotators.

He will then analyse and evaluate different techniques for identifying good annotators and determining consensus ratings for the assets which will be used as training data for developing prediction systems for emotion recognition."

Experience in our group

- ◆ Preliminary rating using crowdsourcing [Brian Vaughan]



Findings

- ◆ clear instructions
- ◆ asset selection strategy
- ◆ payment amounts

References

- ◆ V. Ambati, S. Vogel, and J. Carbonell. Active Learning and Crowd-Sourcing for Machine Translation. In Proc. of LREC'10, pages 2169–2174, 2010.
- ◆ A. Brew, D. Greene, and P. Cunningham. Using Crowdsourcing and Active Learning to Track Sentiment in Online Media. In Proc. of PAIS 2010, pages 1–11. IOS Press, 2010.
- ◆ J. Howe. Crowdsourcing: Why the Power of the Crowd Is Driving the Future of Business. Crown Business, 2008.
- ◆ A. Kittur, E. Chi, and B. Suh. Crowdsourcing for Usability: Using Micro-Task Markets for Rapid, Remote, and Low-cost User Measurements. Proceedings of CHI 2008.
- ◆ V.C. Raykar, S. Yu, L.H. Zhao, A. Jerebko, C. Florin, G.H. Valadez, L. Bogoni, and L. Moy. Supervised Learning from Multiple Experts: Whom to trust when everyone lies a bit. In Proc. of ICML-2009, pages 889–896, 2009.
- ◆ V.C. Raykar, S. Yu, L.H. Zhao, G.H. Valadez, C. Florin, L. Bogoni, and L. Moy. Learning from Crowds. Journal of Machine Learning Research, 11:1297–1322, 2010.
- ◆ P. Smyth, U. Fayyad, M. Burl, P. Perona, and P. Baldi. Inferring Ground Truth from Subjective Labelling of Venus Images. Advances in neural information processing systems, 7:1085–1092, 1995.
- ◆ R. Snow, B. O'Connor, D. Jurafsky, and A.Y. Ng. Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks. In Proceedings of the Conference on Empirical Methods in Natural Language Processing, pages 254–263. ACL, 2008.
- ◆ A. Sorokin and D. Forsyth. Utility data annotation with Amazon Mechanical Turk. In Proc. of CVPR 2008, pages 1–8, 2008.
- ◆ L. von Ahn and L. Dabbish. Labelling Images with a Computer Game, In Procs of CHI 2004