# Two Tools are Better Than One: Tool Diversity as a Means of Improving Aggregate Crowd Performance

JEAN Y. SONG, RAYMOND FOK, ALAN LUNDGARD, FAN YANG, JUHO KIM, WALTER S. LASECKI

# Crowdsourcing Platforms

# Crowdsourcing for Human Computation



https://playment.io/

https://www.crowdguru.de/en/

# Crowdsourcing Strategy: Microtasking

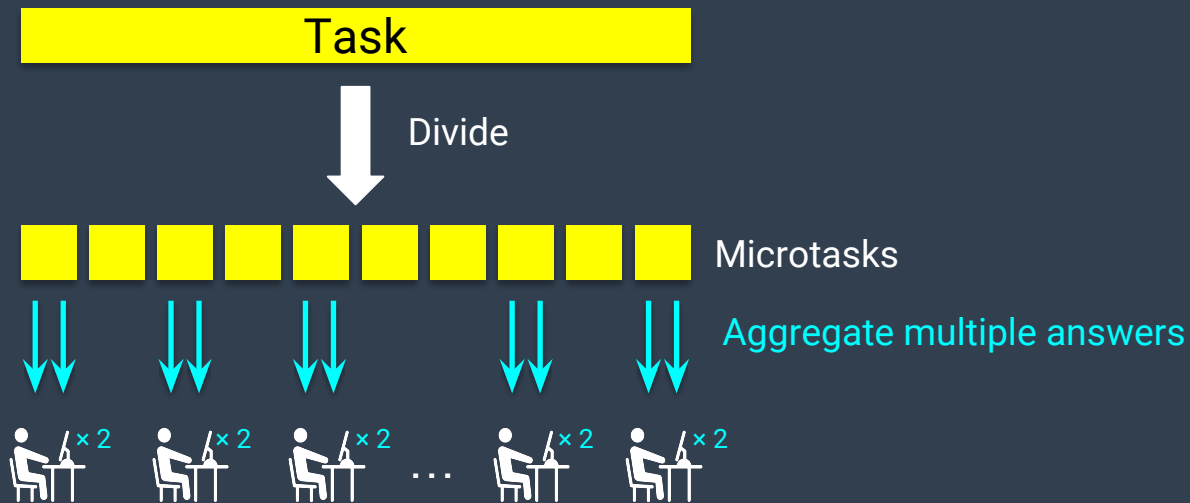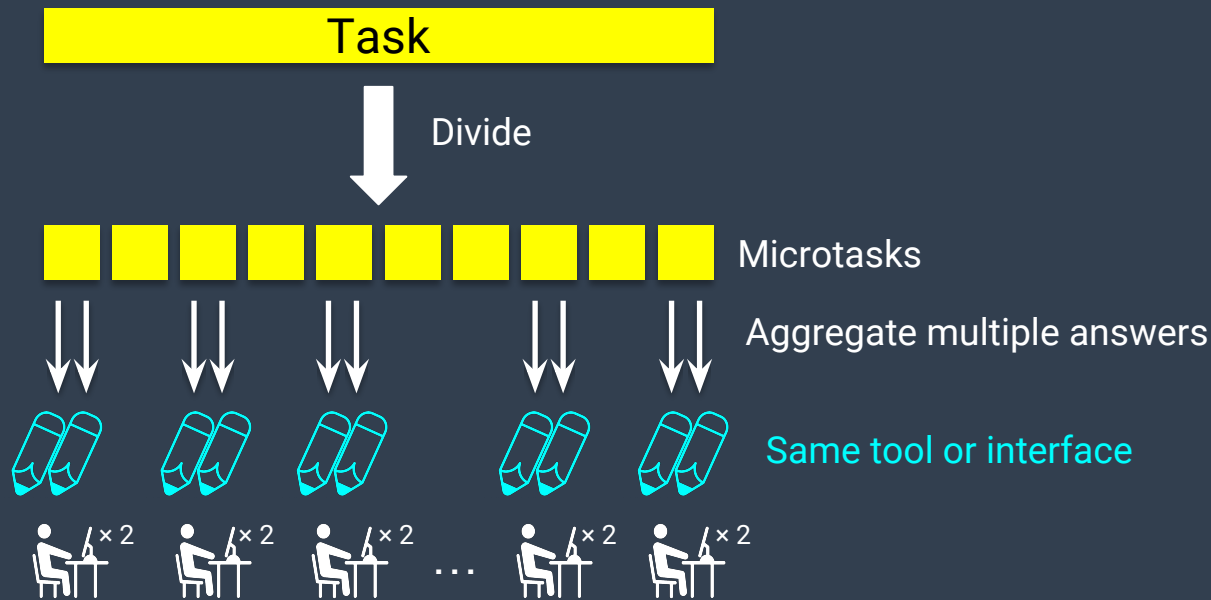# Crowdsourcing Strategy: Aggregation

# Crowdsourcing Strategy: Using Single Tool

Task

Divide

Microtasks

Aggregate multiple answers

Same tool or interface

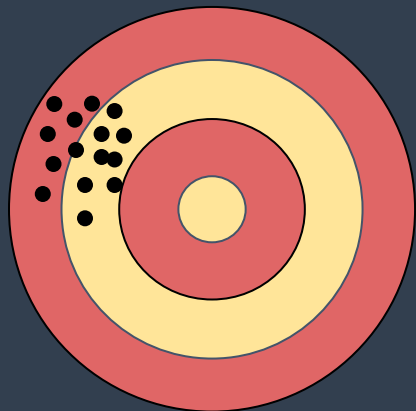× 2    × 2    × 2    ...    × 2    × 2

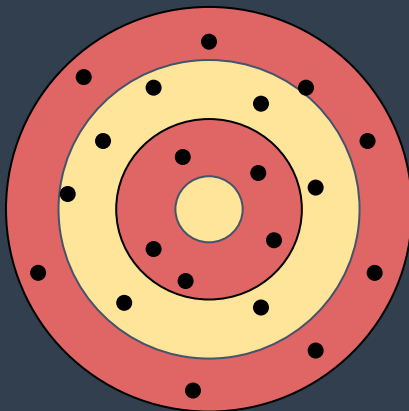**Problem with using a single tool:**

**Systematic bias** **can be accumulated, resulting in inaccurate aggregated result.**

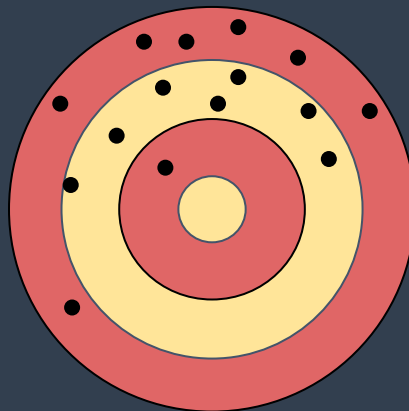# Q. What is Systematic Bias?

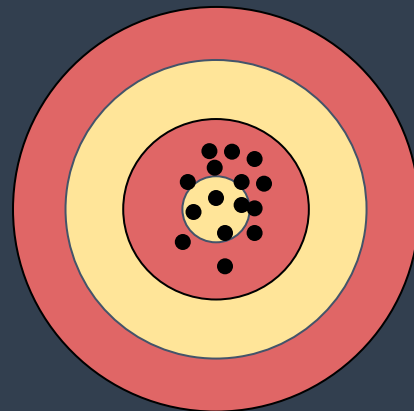A. Reliable, but not valid performance



Reliable,
not Valid
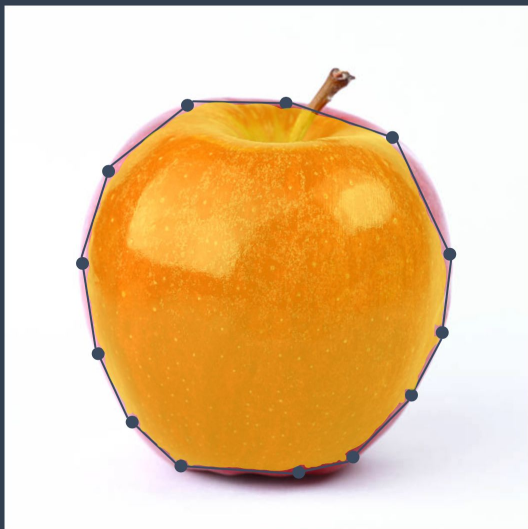
Not Reliable,
But Valid

Not Reliable,
not Valid

Reliable,
Valid

# Example of Systematic (Error) Bias

## Tool 1: Opensurfaces (TOG 2013)

Bell, Sean, et al. "**Opensurfaces**: A richly annotated catalog of surface appearance." *ACM Transactions on Graphics (TOG)*32.4 (2013): 111.



## Tool 2: Click'n'Cut (CrowdMM 2014)

Carlier, Axel, et al. "**Click'n'Cut**: Crowdsourced interactive segmentation with object candidates." *International ACM Workshop on Crowdsourcing for Multimedia*. 2014.
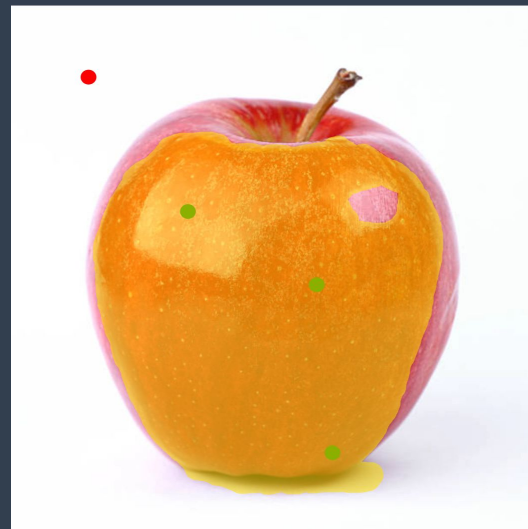
# Example of Systematic (Error) Bias

**Tool 1: Opensurfaces (TOG 2013)**

Bell, Sean, et al. "**Opensurfaces**: A richly annotated catalog of surface appearance." *ACM Transactions on Graphics (TOG)*32.4 (2013): 111.



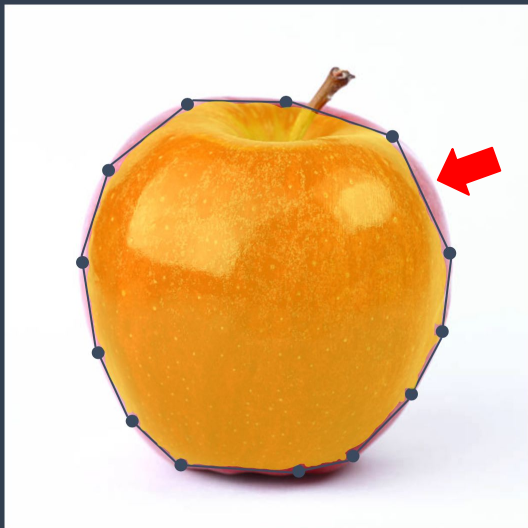**Tool 2: Click'n'Cut (CrowdMM 2014)**

Carlier, Axel, et al. "**Click'n'Cut**: Crowdsourced interactive segmentation with object candidates." *International ACM Workshop on Crowdsourcing for Multimedia*. 2014.

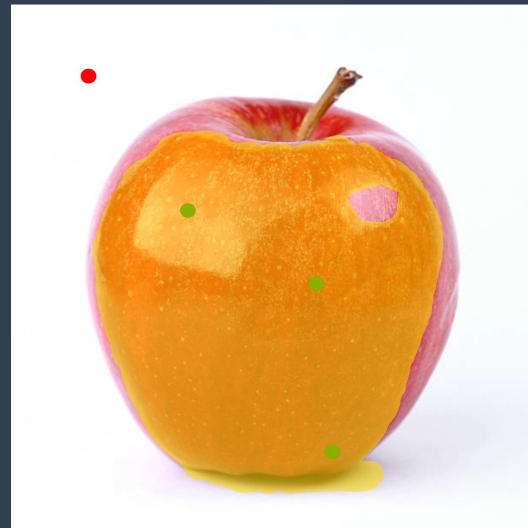# Example of Systematic (Error) Bias

## Tool 1: Opensurfaces (TOG 2013)

Bell, Sean, et al. "**Opensurfaces**: A richly annotated catalog of surface appearance." *ACM Transactions on Graphics (TOG)*32.4 (2013): 111.
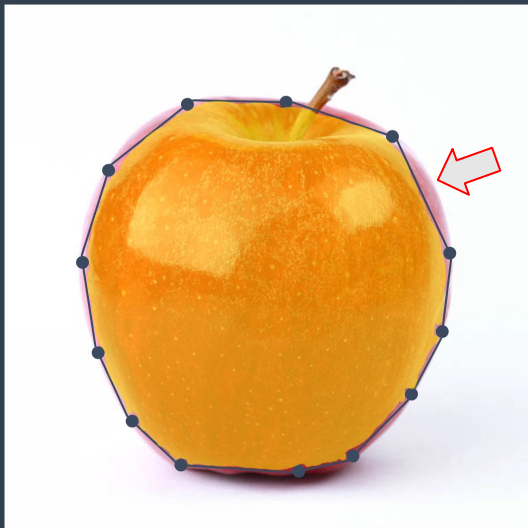


## Tool 2: Click'n'Cut (CrowdMM 2014)

Carlier, Axel, et al. "**Click'n'Cut**: Crowdsourced interactive segmentation with object candidates." *International ACM Workshop on Crowdsourcing for Multimedia*. 2014.
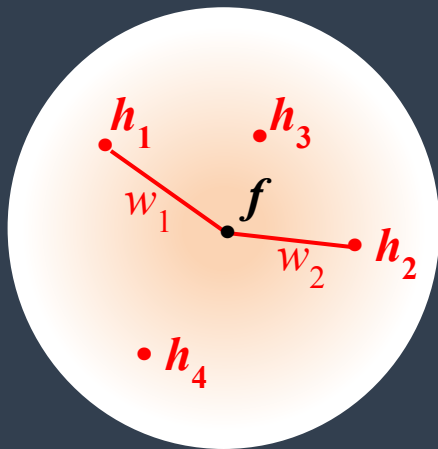
**Proposed Approach:**

**Use tool diversity as a means of improving aggregate crowd performance**

**What is Tool Diversity?**

**A property that measures how different tools can be built in terms of their induced biases.**

# Analogy to Ensemble Learning
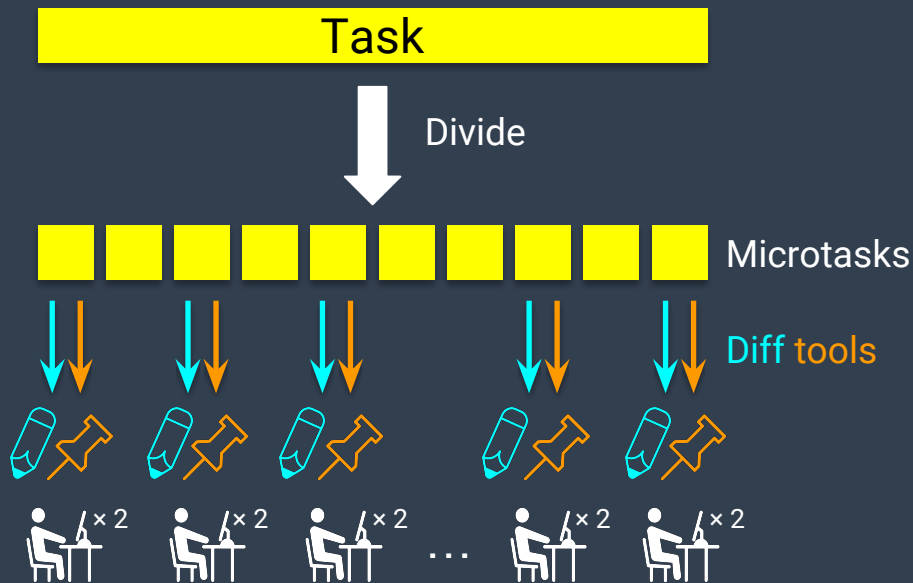
Space of
hypotheses



$f$ : best performing hypothesis
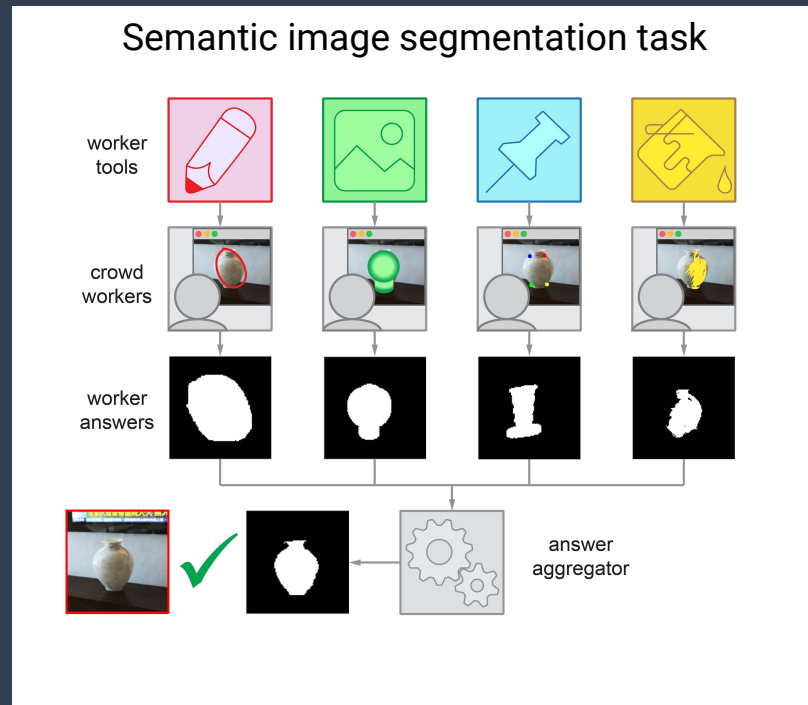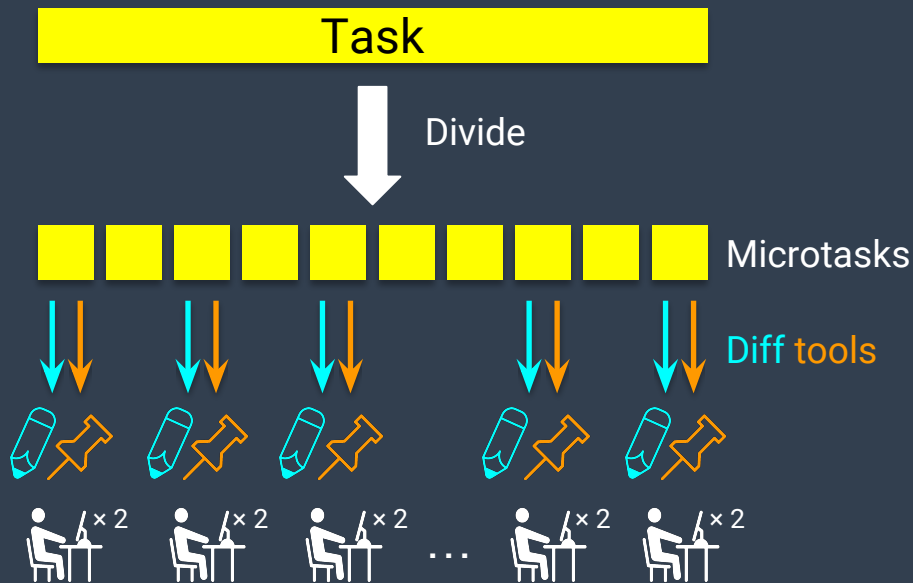$h_i$ : other hypotheses
$w_i$ : weights

Ensemble learning constructs a combination of
two alternative hypotheses $h_1$ and $h_2$ with
proper weights ($w_1$ and $w_2$), and approximates
the best hypothesis $f$ by averaging the two.

# Proposed Method: Leverage Tool Diversity



Task

Divide

Microtasks

Diff tools

× 2  × 2  × 2  ...  × 2  × 2

# Proposed Method: Leverage Tool Diversity



Task

Divide

Microtasks

Diff tools

Semantic image segmentation task

worker tools

crowd workers

worker answers

answer aggregator

# Choosing the Tools

Q. How to diversify errors produced by different tool types?

# Choosing the Tools

Q. How to diversify errors produced by different tool types?

Q. What are different types of objects?

A. General objects, Fuzzy materials, plants, furry objects,

transparent objects, reflective surfaces (intuitive, deformability)

$T_1$ $T_2$ $T_3$ $T_4$

# Instructions and Worker Interface

Worker Interface :

# Instructions and Worker Interface

## Instructions :

# Experiment Settings

- 12 different visual scenes
- Total 51 objects
- Six unique workers for each tool-scene pair (total 288+ workers)
- Total 1224 object segmentations
- Platform: Amazon Mechanical Turk

Each worker was paid between $0.35 and $0.60 per task, depending on the number of objects they had to segment or on the level of difficulty of given tool (a pay rate of ~$10/hr).
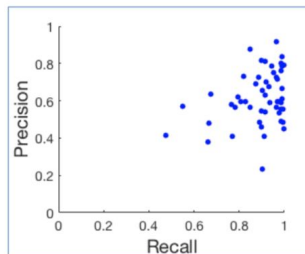
# Results & Discussion

Accuracy Metrics

$$Precision\ (P) = \frac{TP}{TP + FP}$$

$$Recall\ (R) = \frac{TP}{TP + FN}$$

$$F_1\ Score = \frac{2 \times P \times R}{(P + R)}$$

TP = True Positive
FP = False Positive
FN = False Negative

Basic Trace ($T_1$)

Drag-and-Drop ($T_2$)

Pin-Placing ($T_3$)

Floodfill ($T_4$)

Different tools have different error patterns

**Accuracy Metrics**

$$Precision (P) = \frac{TP}{TP + FP}$$
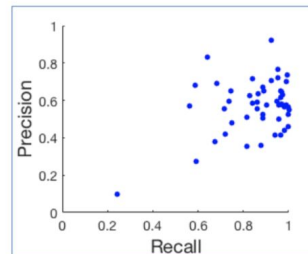
$$Recall (R) = \frac{TP}{TP + FN}$$

$$F_1\ Score = \frac{2 \times P \times R}{(P + R)}$$

TP = True Positive
FP = False Positive
FN = False Negative



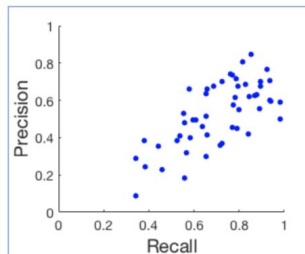Basic Trace ($T_1$)    Drag-and-Drop ($T_2$)

Pin-Placing ($T_3$)    Floodfill ($T_4$)

Different tools have different error patterns

# What we observed

## Accuracy Metrics

$$\text{Precision (P)} = \frac{TP}{TP + FP}$$

$$\text{Recall (R)} = \frac{TP}{TP + FN}$$

$$F_1 \text{ Score} = \frac{2 \times P \times R}{(P + R)}$$

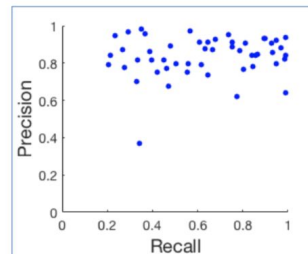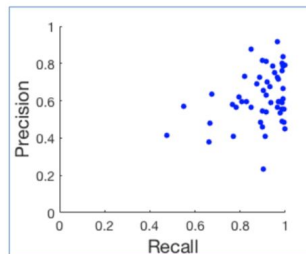TP = True Positive
FP = False Positive
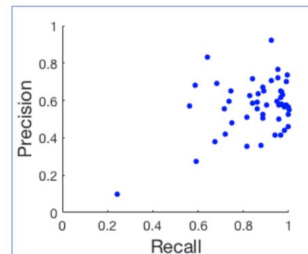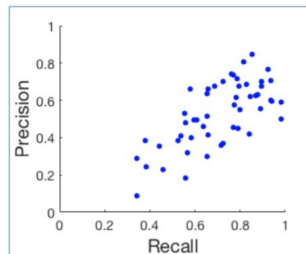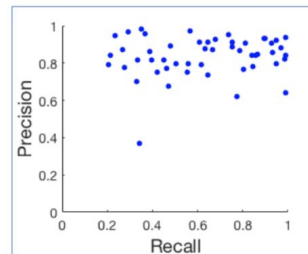FN = False Negative



Different tools have different error patterns

How can we see
the effect of leveraging tool diversity?

# Comparison of Aggregation Methods

Method 1. Single tool aggregation (Uniform majority voting): Baseline

$T_1$                                            → Aggregate

$T_2$                                            → Aggregate

# Comparison of Aggregation Methods

Method 2. Multiple tool aggregation (Uniform majority voting)

$T_1 \times T_2$    $w$    $w$    $w$    $w$    $\rightarrow$ Aggregate

Method 3. Multiple tool aggregation (Expectation maximization)

$T_1 \times T_2$    $w_1$    $w_2$    $w_3$    $w_4$    $\rightarrow$ Aggregate

# Comparison of Aggregation Methods

# Comparison of Aggregation Methods



Basic Trace ($T_1$)   Drag-and-Drop ($T_2$)

**High recall**

Pin-Placing ($T_3$)   Floodfill ($T_4$)

**High precision**

High recall + high precision pairs gave the highest performance improvement.

\* significant at $p < .05$, \*\* significant at $p < .01$ compared to **Multiple (EM)**

Single (majority voting, superior)   Single (majority voting, inferior)   Multiple (majority voting)   Multiple (EM)

# Generalization

# Generalizability: Expected Human Error is Diverse



Tool 1

Tool 2

Aggregate

Reliable,
Valid

# Generalizability: Aggregation Improves Quality

**Quality Improves**

# Generalizability: Objective Correct Answer Exists

**Tasks with objective answers:**

Task with subjective answers:

Creative writing

Image segmentation

Live captioning

Text annotation

Handwriting recognition

This paper presents Soylent, a word processing interface that uses crowd workers to help with proofreading, document shortening, editing and commenting tasks. Soylent is ~~an example of a~~ new ~~kind of~~ interactive user interface in which the end user has direct access to a crowd of workers for assistance with tasks that require human attention and common sense. Implementing these ~~kinds of~~ interfaces requires new software programming patterns ~~for interface software~~, since crowds behave differently than computer systems. We have introduced one important pattern, Find-Fix-Verify, which splits complex editing tasks into a series of identification, generation, and verification stages ~~that use independent agreement and voting~~ to produce reliable results. We evaluated Soylent with a range of editing tasks, finding and correcting 82% of grammar errors ~~when combined with automatic checking~~, shortening text to approximately 85% of original length per iteration, and executing a variety of human macros successfully.

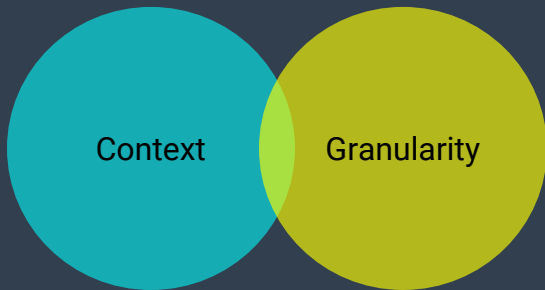# Generalizability: Tolerates Imperfections

Example: Scribe (UIST 2012)

W.S. Lasecki, C.D. Miller, A. Sadilek, A. Abumoussa, D. Borrello, R. Kushalnagar, J.P. Bigham.
Real-time Captioning by Groups of Non-Experts. UIST 2012.

```
1:    learn  g is such        a  suitcase word though right   so         has a lot of there     s a lot
2:  o learning is such                                                                there a are a lot
3:    learning ss such         a  suitcase word though          learning has              is a lot
4:    lea ning is su h         a                      right  so learning                    a lot
5: so learning is such         a  suitcase      though          learning has                  lot
6:    learning is such         a  suitcfse word though right                    this       in a lot
F: so learning is such         a  suitcase word though right  so learning has a lot of there   is a lot
```
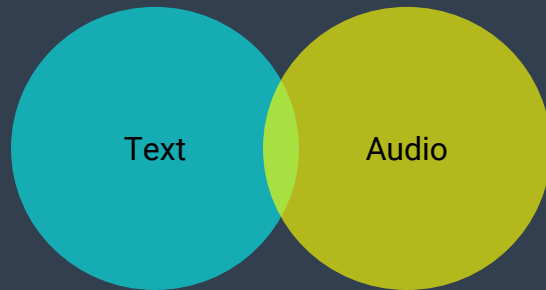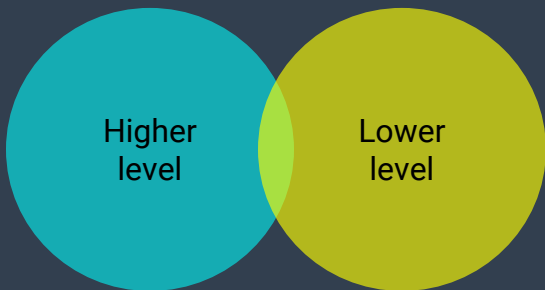
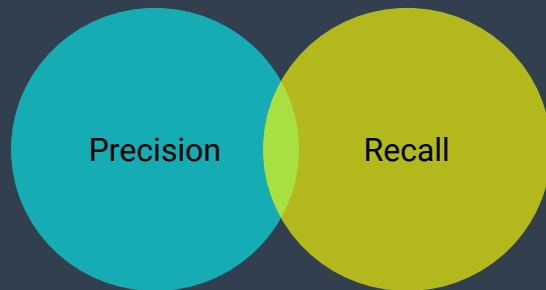# Possible Future Applications



Application1: Tagging Long Videos

Context — Granularity

Application2: Multichannel NLP

Text — Audio

Application3: Complex/Diverse Annotation

Higher level — Lower level

Application4: Computer-Human Integration

Precision — Recall

# Thank you!

Authors:
    Jean Y. Song (jyskwon@umich.edu / jyskwon.github.io),
    Raymond Fok, Alan Lundgard, Fan Yang, Juho Kim, Walter S. Lasecki

# Backup Slides

# Tool 1



Basic Trace (T$_1$)

# Tool 2

# Tool 3

# Tool 4

# Pixel-Level Majority Voting (50% agreement)

# Expectation Maximization (Dawid-Skene Algorithm)

In an image, label a pixel as $1$ if it belongs to a target object, and $0$ if background.

Assume:

- image $A$ having $N$ total pixels
- $M$ crowd workers
- The label a worker $m$ assigns to each pixel is denoted as $z_{mn}$
- all labels from worker $m$ as a vector $Z_m$
- the true labels of $A$ to be estimated are denoted as a vector $Y$
- θ is the confusion matrices set to be estimated.

We can estimate the true labels $Y$ by maximizing the marginal likelihood of the observed worker labels:

$$l(\boldsymbol{\theta}) := \log\Big(\sum_{Y \in \{0,1\}^n} L(\boldsymbol{\theta}; Y, Z)\Big)$$

The EM algorithm works iteratively by applying the 1) expectation step and the 2) maximization step.