# Exploring the Use of Crowdsourcing to Support Empirical Studies in Software Engineering

Kathryn T. Stolee & Sebastian Elbaum

September 16, 2010

**ESQuaReD**

Nebraska UNIVERSITY OF
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
Objective

# Introduction

### Known Issue

Recruiting the right type and number of users for empirical studies in software engineering is hard.

ESQuaReD

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
Objective

# Introduction

## Known Issue

Recruiting the right type and number of users for empirical studies in software engineering is hard.

**Possible Solutions:**

- Use fewer participants of the right type
  - Limits generalizability to larger groups

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
Objective

# Introduction

### Known Issue

Recruiting the right type and number of users for empirical studies in software engineering is hard.

**Possible Solutions:**

- Use fewer participants of the right type
  - Limits generalizability to larger groups
- Relax requirements for participation
  - Limits generalizability to target population

**ESQuaReD**

**Nebraska** UNIVERSITY OF
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
Objective

# Introduction

## Known Issue

Recruiting the right type and number of users for empirical studies in software engineering is hard.

**Possible Solutions:**

- Use fewer participants of the right type
    - Limits generalizability to larger groups
- Relax requirements for participation
    - Limits generalizability to target population

- Crowdsource the study

**ESQuaReD**

**Nebraska** UNIVERSITY of
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
**Background**
Background
Objective

# Background

## Crowdsourcing

Leveraging a global community of users with different talents and backgrounds to help perform a task that would not be feasible without a mass of people behind it.



**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Crowdsourcing Services (examples)



Companies with hard problems connect with people interested in solving. 1,000+ problems, 200,000+ solvers

**ESQuaReD**

UNIVERSITY OF
Nebraska
Lincoln

**Introduction**
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Crowdsourcing Services (examples)

**INNO**CENTIVE®

Companies with hard problems connect with people interested in solving. 1,000+ problems, 200,000+ solvers

iStockphoto®

Photographers collect with people who need stock photography. 3,000,000+ members

**ESQuaReD**

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Crowdsourcing Services (examples)



Companies with hard problems connect with people interested in solving. 1,000+ problems, 200,000+ solvers



Photographers collect with people who need stock photography. 3,000,000+ members



Companies with scientific problems connect with retired scientists. 1,000+ companies, 5,000+ scientists

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Crowdsourcing Services (examples)

**INNO**CENTIVE®

Companies with hard problems connect with people interested in solving. 1,000+ problems, 200,000+ solvers

iStockphoto®

Photographers collect with people who need stock photography. 3,000,000+ members

*your*encore®

Companies with scientific problems connect with retired scientists. 1,000+ companies, 5,000+ scientists

amazon mechanical turk
Artificial Artificial Intelligence

People with many small tasks connect with scalable workforce. 100,000+ tasks, 100,000+ workers
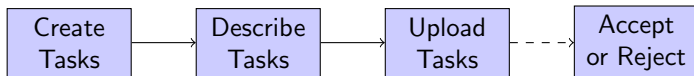
e² ESQuaReD

Nebraska
UNIVERSITY OF
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Workflow in Mechanical Turk

**Requestors:**

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Workflow in Mechanical Turk
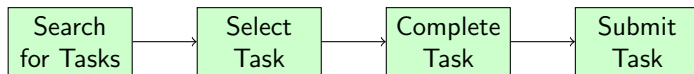
**Requestors:**



**Types of tasks:**

- Short duration (60s. or less)
- Require human intelligence (handwirting analysis, image tagging)
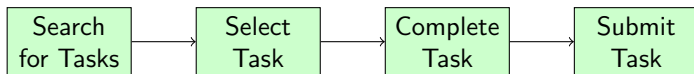- Specialized (requires certain knowledge) or generic

**ESQuaReD**

**Nebraska** Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Workflow in Mechanical Turk

**Workers:**



**ESQuaReD**

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Workflow in Mechanical Turk

**Workers:**

| Search for Tasks | → | Select Task | → | Complete Task | → | Submit Task |

---

**Answer Two Short Questions about Yahoo! Pipes - Easy!**                              View a HIT in this group

| **Requester:** | Katie Stolee | **HIT Expiration Date:** | May 13, 2010 (3 days 8 hours) | **Reward:** | $0.20 |
| | | **Time Allotted:** | 60 minutes | **HITs Available:** | 8 |

**Description:** The task is to answer two short questions, comparing two versions of Yahoo! Pipes programs that have the same output.

**Keywords:** programming, Yahoo, Pipes, survey, mashup, questionnaire, coding, easy

**Qualifications Required:**                                                  **Your Value**
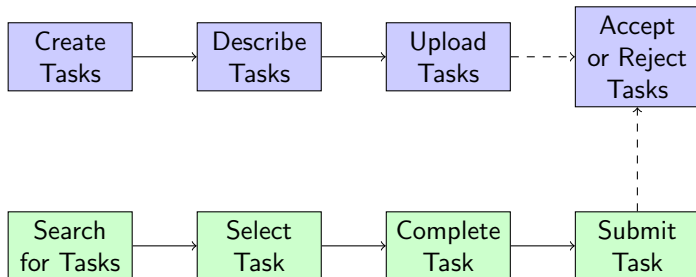
Qualification Quiz for UNL Study on Yahoo! Pipes is greater than 90          100 You meet this qualification requirement

HIT approval rate (%) is greater than 90                                     100 You meet this qualification requirement   Contact the Requester of this HIT

**ESQuaReD**

Nebraska
Lincoln
UNIVERSITY OF

Introduction
Mechanical Turk Study
Summary

Motivation
Background
**Background**
Objective

# Workflow in Mechanical Turk

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
**Objective**

# Goal of This Work

### Conjecture

Crowdsourcing can be a good solution for recruiting the right type and quantity of participants for an empirical study in software engineering.
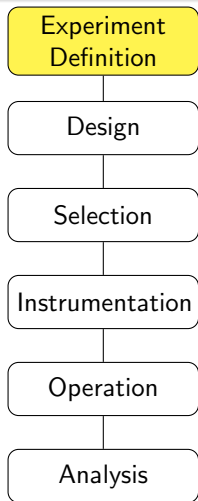
**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Motivation
Background
Background
**Objective**

# Goal of This Work

### Conjecture

Crowdsourcing can be a good solution for recruiting the right type and quantity of participants for an empirical study in software engineering.

In this work, we crowdsource a software engineering experiment using Amazon's Mechanical Turk service, and reflect on our experiences.

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Study Definition

```
┌─────────────────┐
│   Experiment    │
│   Definition    │
└─────────────────┘
         │
┌─────────────────┐
│     Design      │
└─────────────────┘
         │
┌─────────────────┐
│    Selection    │
└─────────────────┘
         │
┌─────────────────┐
│ Instrumentation │
└─────────────────┘
         │
┌─────────────────┐
│    Operation    │
└─────────────────┘
         │
┌─────────────────┐
│    Analysis     │
└─────────────────┘
```

**Purpose:** Evaluate the impact of coding practices (e.g., code smells) on end user's preferences and understanding of web mashups built in Yahoo! Pipes.

**ESQuaReD**

Introduction
**Mechanical Turk Study**
Summary

Definition
**Planning**
Operation
Analysis

# Experimental Task Example
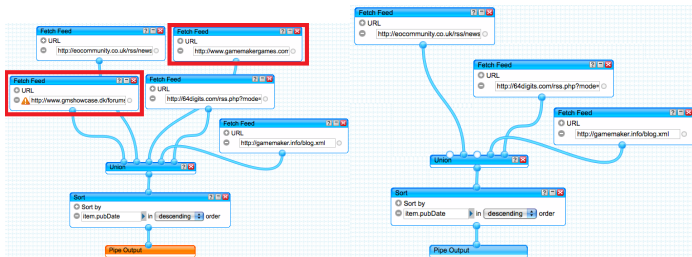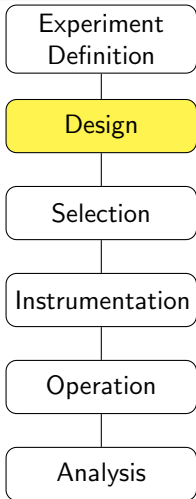
Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**Task Description:** Given two pipes with the same behavior, one with a smell and one without, select the preferable one.



ESQuaReD

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Experimental Task Example



**Task Description:** Given two pipes with the same behavior, one with a smell and one without, select the preferable one.

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Experimental Design

Experiment Definition

Design

Selection

Instrumentation

Operation

Analysis

| Task | Subjects | Pretest | Object | Treatment | Posttest |
|------|----------|---------|--------|-----------|----------|
| 1 | R | $O_1, O_2$ | $Pipe_1$ | $Smell_5$ | $O_3, O_4$ |
| 2 | R | $O_1, O_2$ | $Pipe_2$ | $Smell_4$ | $O_3, O_4$ |
| 3 | R | $O_1, O_2$ | $Pipe_3$ | $Smell_5$ | $O_3, O_4$ |
| 4 | R | $O_1, O_2$ | $Pipe_4$ | $Smell_8$ | $O_3, O_4$ |
| 5 | R | $O_1, O_2$ | $Pipe_5$ | $Smell_7$ | $O_3, O_4$ |
| 6 | R | $O_1, O_2$ | $Pipe_6$ | $Smell_1$ | $O_3, O_4$ |
| 7 | R | $O_1, O_2$ | $Pipe_7$ | $Smell_{5,10}$ | $O_3, O_4$ |
| 8 | R | $O_1, O_2$ | $Pipe_8$ | $Smell_{2,9}$ | $O_3, O_4$ |

$O_1 =$ Education
$O_2 =$ Pipes test score
$O_3 =$ Preference
$O_4 =$ Time to completion

**ESQuaReD**

Nebraska
UNIVERSITY OF
Lincoln

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Experimental Design

Experiment
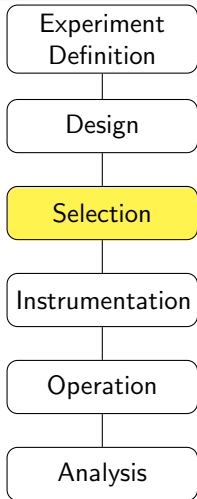Definition

Design

Selection

Instrumentation

Operation

Analysis

**Lessons Learned:**

- Experimental tasks must be modular and independent, but can be longer (ours took 3-4 minutes, on average)
- Qualification tests can be used to capture pretest measures
- Cannot control which tasks are completed by which participants
- Self-selection of tasks may introduce bias that needs to be accounted for in the analysis

**ESQuaReD**

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Selection and Recruitment

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**Desired Participant Characteristics:**

- Limited computer science education (end users)
- Familiar with Yahoo! Pipes

**Mechanical Turk:**

- Facilitates recruitment by hosting tasks
- Allows for qualification tests to be administered prior to participation (pretest measures)

**ESQuaReD**

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Selection and Recruitment

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**Lessons Learned**

- 50 qualification tests submitted in two weeks, 38 passed
- 22 participants in total, 14 were considered "end users"
- More variation and unknowns in participants (e.g., age, gender, education, experimental context)

e² **ESQuaReD**

Nebraska
Lincoln

Introduction
**Mechanical Turk Study**
Summary

Definition
**Planning**
Operation
Analysis

16 / 18

# Experimental Task in Mechanical Turk

Experiment
Definition

Design

Selection

**Instrumentation**

Operation

Analysis



**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Instrumentation

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**ESQuaReD**

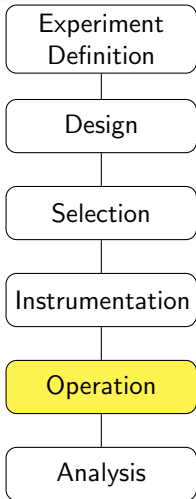**Lessons Learned**

- Need to learn how to use a new tool and/or API
- Need to adjust presentation of tasks to fit the Mechanical Turk interface
- All tasks are in competition with other tasks for participants, so the task description must be enticing.
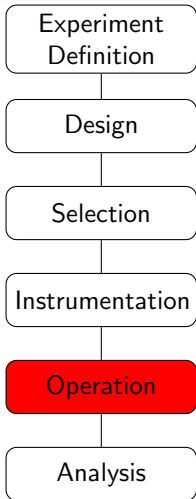
Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Definition
Planning
**Operation**
Analysis

# Experiment Operation

Experiment Definition

Design

Selection

Instrumentation

**Operation**

Analysis

**Mechanical Turk:**

- Hosts tasks for a custom time period (2 weeks)
- Administers qualification tests (50 requests)
- Maintains user anonymity
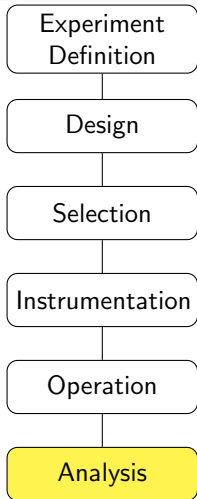- Collects results and metrics (188 tasks submitted)

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Experiment Operation

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**ESQuaReD**

**Lessons Learned:**

- Hand-grading qualification tests introduce delay, and may discourage further participation
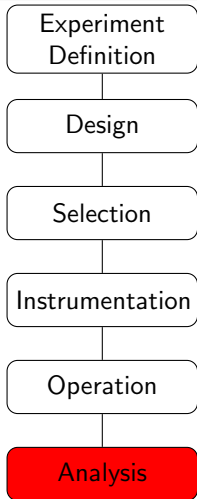- Time to completion is reported, but is suspicious

Nebraska
Lincoln

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Analysis

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**Response Quality:**

- Qualitative responses were detailed and demonstrated understanding (Average length was 31 words, only 10 were required)

- Did not need to reject any responses

**ESQuaReD**

Introduction
Mechanical Turk Study
Summary

Definition
Planning
Operation
Analysis

# Analysis

Experiment
Definition

Design

Selection

Instrumentation

Operation

Analysis

**Lessons Learned:**

- We were able to validate our hypotheses (for only $42)
- May need to throw away some data due to learning (we threw away 28 responses)
- Too many responses from a small group of participants could skew results

**ESQuaReD**

# Summary

**Crowdsourcing allowed us to:**

- Obtain a sufficient number of participants with the desired characteristics
- Evaluate our research questions using an empirical study for low cost

**However...**

- Requires careful experimental design to work within the Mechanical Turk infrastructure
- Due to the "unknowns" about the subjects and environment, crowdsourcing may not be appropriate for all studies

**ESQuaReD**