# Crowdsourced Data Management Survey

Group 13
The University of Melbourne

# Overview

1. Introduction
   a. Our Team
   b. Our Project Goal
   c. Our Project Timeline

2. Approach
   a. Taxonomy
   b. Literature Review
      i. Application & Platforms
      ii. Techniques

3. Conclusion
   a. Discussion
   b. Future Directions

# Team & Project Goal

Project Goal:

1. Compose a survey to help the general public understand Crowdsourced Data Management。

2. Illustrate the survey results through presentation.

# Communication & Research Tools
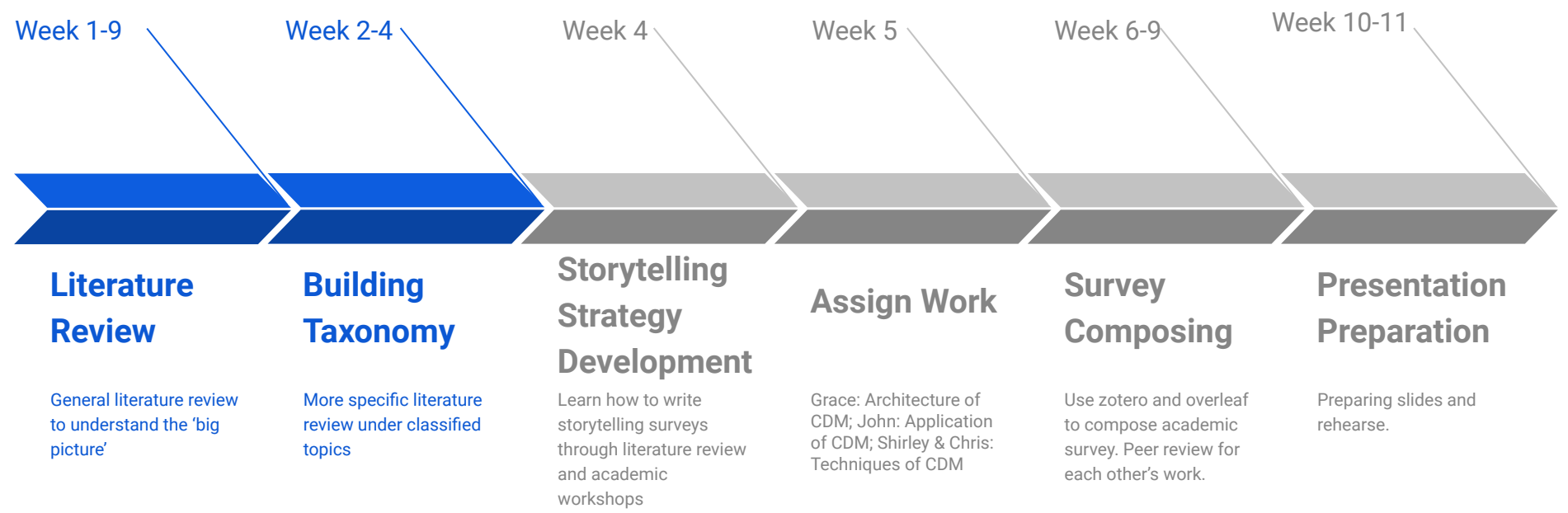
Manage weekly meetings & tasks
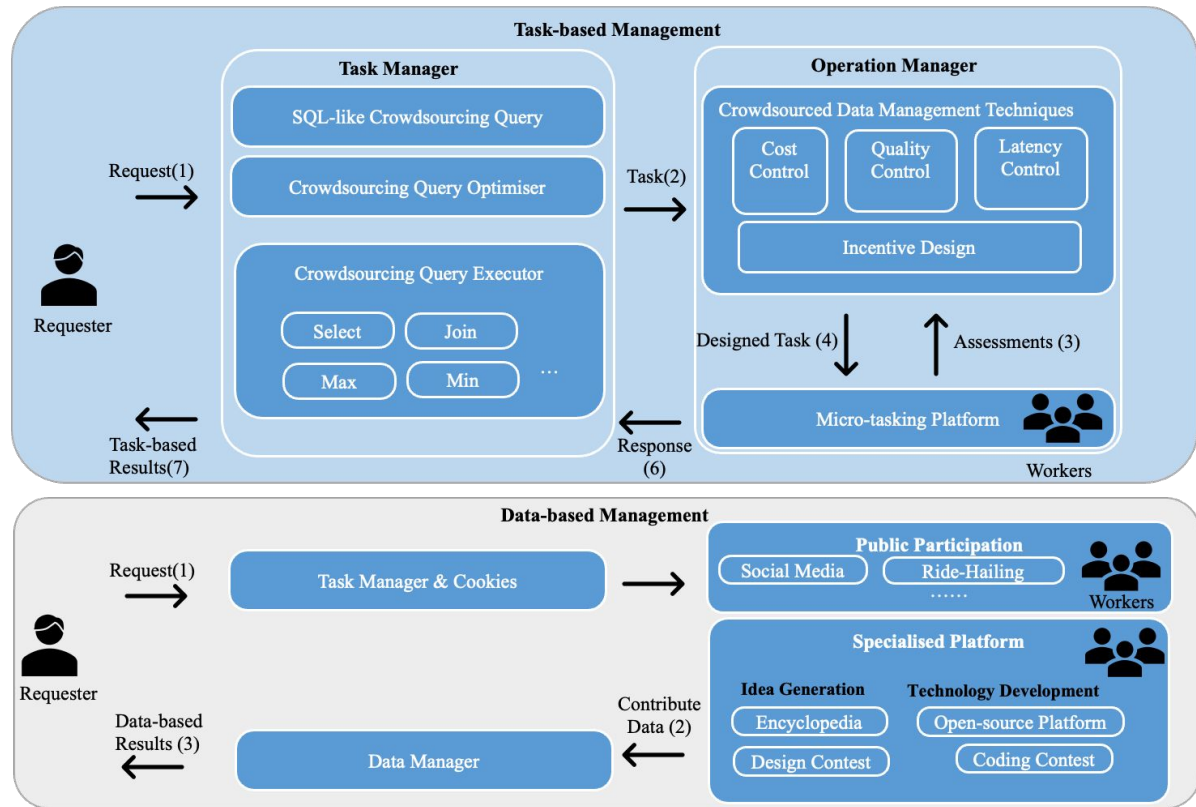
Manage report and reference

# Project Timeline

Approx.10 hours per week

| Week 1-9 | Week 2-4 | Week 4 | Week 5 | Week 6-9 | Week 10-11 |

**Literature Review**

General literature review to understand the 'big picture'

**Building Taxonomy**

More specific literature review under classified topics

**Storytelling Strategy Development**

Learn how to write storytelling surveys through literature review and academic workshops

**Assign Work**

Grace: Architecture of CDM; John: Application of CDM; Shirley & Chris: Techniques of CDM

**Survey Composing**

Use zotero and overleaf to compose academic survey. Peer review for each other's work.

**Presentation Preparation**

Preparing slides and rehearse.

# What is Crowdsourced Data Management? (Ellis, 2004; Howe, 2006)

| | |
|---|---|
| 1714 | Longitude Prize |
| 1936 | Toyota Logo Design |
| 2005 | Launch of Amazon Mechanical Turk |
| 2006 | Conceptualisation of "Crowdsourcing" |
| 2000s | Crowdsourcing Applications & Technology Explosion |
| **2010s** | **Crowdsourced Data Management Applications & Techniques** |
| 2020s | Nowadays Application |

**Task-based Management**

**Task Manager**

- SQL-like Crowdsourcing Query
- Crowdsourcing Query Optimiser
- Crowdsourcing Query Executor
  - Select
  - Join
  - Max
  - Min
  - …

Requester — Request(1) →

Task(2) →

**Operation Manager**

Crowdsourced Data Management Techniques
- Cost Control
- Quality Control
- Latency Control
- Incentive Design

Designed Task (4) ↓  Assessments (3) ↑

Micro-tasking Platform

Workers

← Task-based Results(7)

← Response (6)

**Data-based Management**

Requester — Request(1) →

Task Manager & Cookies →

**Public Participation**
- Social Media
- Ride-Hailing
- ……

Workers

**Specialised Platform**

Idea Generation
- Encyclopedia
- Design Contest

Technology Development
- Open-source Platform
- Coding Contest

Workers

← Data-based Results (3)

Data Manager

← Contribute Data (2)

# Taxonomy



**Crowdsourced Data Management (CDM)**

- CDM Overview
  - Historical Development
    - Reason of using CDM
    - Difference between CDM and traditional DM
  - Architecture
    - Requesters
      - Financial Constraint
      - Temporal Constraint
    - Workers
      - Uncertain Reliability
      - Motivated by Incentive
    - Platform
      - Genral Platform
      - Specialised Platform
- CDM Application
  - Micro-tasking
    - Data Cleaning, Data Labelling, Data Integration
      - Amazon Mechanical Turk — Genral Platform
  - Other Application
    - Idea Generation — Wikipedia
    - Public Participation — Google Map (Point of Interest)
    - Technology Development — Linux
      - Specialised Platform
- CDM Techiniques
  - Cost Control
    - Task pruning and Task selection
    - Sampling and Answer Deduction
  - Quality Control
    - Enhance Data Quality
    - Modelling Worker's Quality
    - Evaluate Worker's Quality
    - Aggregating Outputs
    - Task Assignment
  - Latency Control
    - Latency Modeling
    - Reduce Recruitment Time
  - Incentive Design
    - Monetary Incentive
    - Non-monetary Incentive

# Platform (Chittilappilly et al., 2016)

## General-purpose platform

- Amazon Mechanical Turk



## Specialised platform

- Kaggle
- Lego Ideas

# General-purpose Platform

## Amazon Mechanical Turk

- Requester
- Worker
- Developer

Cost Control
- Multiple assignments
- Iterative Learning
- …

Latency Control
- Monetary Rewards

mTurk System

Quality Control

Aggregating Results
Major voting
…
Qualification Test

Work & Contribution

Rewards

Large Pool of Workers

HIT > 10,000    HIT > 10,000    HIT > 10,000

# General-purpose Platform

## Microtasks in AMT

- Data Labelling

- Data Cleaning

- Data integration

Bounding Box

Image Labelling

Survey

# Specialised Platform

## Kaggle

- Sponsors
- Workers

**Kaggle System**

Contribution | Reward

**Competition**

Quality Control
- Leaderboard

Cost Control
- Multiple assignments

Latency Control
- Monetary Rewards

Work

Reward

**Workers with expertise**

# Quality Control

- Data quality may be influenced by several factors

- Various approaches have been proposed to guarantee high quality results

  - ❖ Enhancing data quality

  - ❖ Selection of workers

  - ❖ Aggregating outputs of different workers

  - ❖ Assigning suitable tasks to high quality workers

COMPARISON

| Quality Control | Pros | Cons |
|---|---|---|
| Aggregating Result | ● Produce more accurate and reliable overall results<br>● Assist identifying malicious or low-quality workers | ● Requires more workers than other methods<br>● Higher cost |
| Majority Voting | ● Easy to implement<br>● Robust to noise | The professional level of different workers and task difficulty is ignored |
| Confusion Matrix | Capture more information | ● Not applicable in all situations<br>● Biases in ground truth. |
| Golden Task | Simplify the process of assessing worker quality. | ● Higher cost from hiring experts<br>● Difficulty in deciding sufficient number of golden tasks<br>● If the answer is leaked, or many requesters use the same golden task, the mechanism will fail |
| Qualification Task | Improves confidence of result's quality from reliable workers | ● Many workers are unwilling to answer "extra" tasks for free.<br>● Poses potential risk that spammers could carefully label these golden tasks to increase their reputation. |

Table 3: Pros and Cons for Task Assignment Techniques

| Techniques | Pros | Cons |
|---|---|---|
| Multiple assignment | • No need to know or infer worker's reliability.<br>• Easy to implement. | • Not efficient: workers didn't get potentially acquainted tasks.<br>• Increasing cost. |
| Iterative Learning | • Worker's reliability estimated based on comparing with other's answers. | • High computational cost. |
| Dual Task Assigner | • Leveraged resource usage.<br>• Worker's reliability can be estimated from their previous performance. | • Difficult to infer worker's reliability when there's no enough historical data. |
| Real-time Task Assigner | • Leveraged resource usage.<br>• Reduced cost on multiple assignments.<br>• Flexible with dynamic task assignments. | • Only near-optimal results can be achieved. |
| Collaborative Task Assignment | • Worker's reliability is known.<br>• More flexible task assignment with worker's domain knowledge.<br>• Can be applied in Knowledge-Intensive and collaborative crowdsourcing settings.<br>• Application is narrow in real context, as most times worker's skill set is unknown. | |
| Task Assignment with AI Planning | • Leveraged task standardisation process<br>• Enabled testing of task allocation strategies with different scenario variables | • Did not perform optimisation under consider budget. |

# Cost Control

- Keeping costs under control without compromising the quality of the results

- Techniques:

  ❖ Fix the number of workers per task based on the requester budget

  ❖ Reduce the number of tasks performed by workers

    - Task Pruning

    - Task Selection

    - Sampling

    - Answer Deduction

# Comparison

| Cost Control | Pros | Cons |
| --- | --- | --- |
| Fixing number of workers | Easy to implement | Wasted expenses |
| Task Pruning | Significantly saves labour costs while maintaining high quality | • Cost cannot be reduced on a per-task level<br>• Limited to certain types of tasks |
| Task Selection | Sufficiently reduce the number of tasks | • Sacrifice some quality<br>• Introduce some delay |
| Answer Deduction and Sampling | Avoid crowds doing a lot of unnecessary work | • Introduce human error<br>• Sampling fail under certain situation |

# Latency Control

- If the requester has a time limit, controlling latency is important

- Several strategies to address this issue:

  ❖ Adjusting price

  ❖ Latency modelling

  ❖ Reducing Recruitment Time

Comparison:

  ❖ Increasing task price will greatly increase the cost

  ❖ Dynamic budget allocation may cause confused or dissatisfied among workers

  ❖ Latency modelling provides more objective decisions but

  ❖ The retainer model can efficiently reduce latency, but may introduce low quality results and more costs

# Incentive

- **Monetary incentives**
  - ❖ Financial Benefits
  - ❖ Straightforward & Easy Implementation

- **Non-monetary incentives**
  - ❖ Individual Development
  - ❖ Public Good
  - ❖ Reputation
  - ❖ Gamification

# Comparison

| Incentive types | Pros | Cons |
|---|---|---|
| Monetary | <ul><li>Straightforward way to motivate individuals to participate.</li><li>Compensation for contributions</li><li>Can attract a consistent and steady involvement of participants.</li><li>Better accuracy and speed compared to voluntary work</li></ul> | <ul><li>Pay rate should be carefully considered to match time and effort required.</li><li>Overpaying or underpaying can lead to issues.</li><li>Participants may attempt the project multiple times, leading to poor data quality.</li><li>Not feasible for project starters with no or low budget</li></ul> |
| Non-Monetary | <ul><li>Cost-effective</li><li>Increase motivation.</li><li>Provide better data quality.</li><li>Can offer opportunities for personal skill development.</li><li>Can contribute to public good.</li><li>Enhance an individual's reputation</li></ul> | <ul><li>Fewer people may be willing to participate.</li><li>Not applicable for all workers</li><li>More time and effort required for task design.</li><li>Different people value different types of incentive</li></ul> |

# Discussion & Future Direction

- Uncertainty

- Bias

- Ethical issue

# Reference

Chittilappilly, A. I., Chen, L., & Amer-Yahia, S. (2016). A Survey of General-Purpose Crowdsourcing Techniques. *IEEE Transactions on Knowledge and Data Engineering*, *28*(9), 2246–2266. https://doi.org/10.1109/TKDE.2016.2555805

Ellis, S. (2014). A History of Collaboration, a Future in Crowdsourcing: Positive Impacts of Cooperation on British Librarianship. *Libri*, *64*(1), 1–10. https://doi.org/10.1515/libri-2014-0001

Jeff, H. (2006, June). *The Rise of Crowdsourcing*. *14.06*. https://www.wired.com/2006/06/crowds/