

Confidentiality Protection in Crowdsourcing

Simran Saxena

Dr. Ponnurangam Kumaraguru (Chair)
Dr. Alpana Dubey (Co-chair)



[linkedin/in/simransaxena](https://www.linkedin.com/in/simransaxena)



[@simran_s21](https://twitter.com/simran_s21)



[fb.com/simran.s21](https://www.facebook.com/simran.s21)

Thesis Committee

- ◆ Dr. Arun Balaji Buduru, IIIT Delhi
- ◆ Dr. Niharika Sachdeva, InfoEdge
- ◆ Dr. Alpana Dubey, Accenture Technology Labs
- ◆ Dr. Ponnurangam Kumaraguru, IIIT Delhi

Demo

◆ Proof of Concept

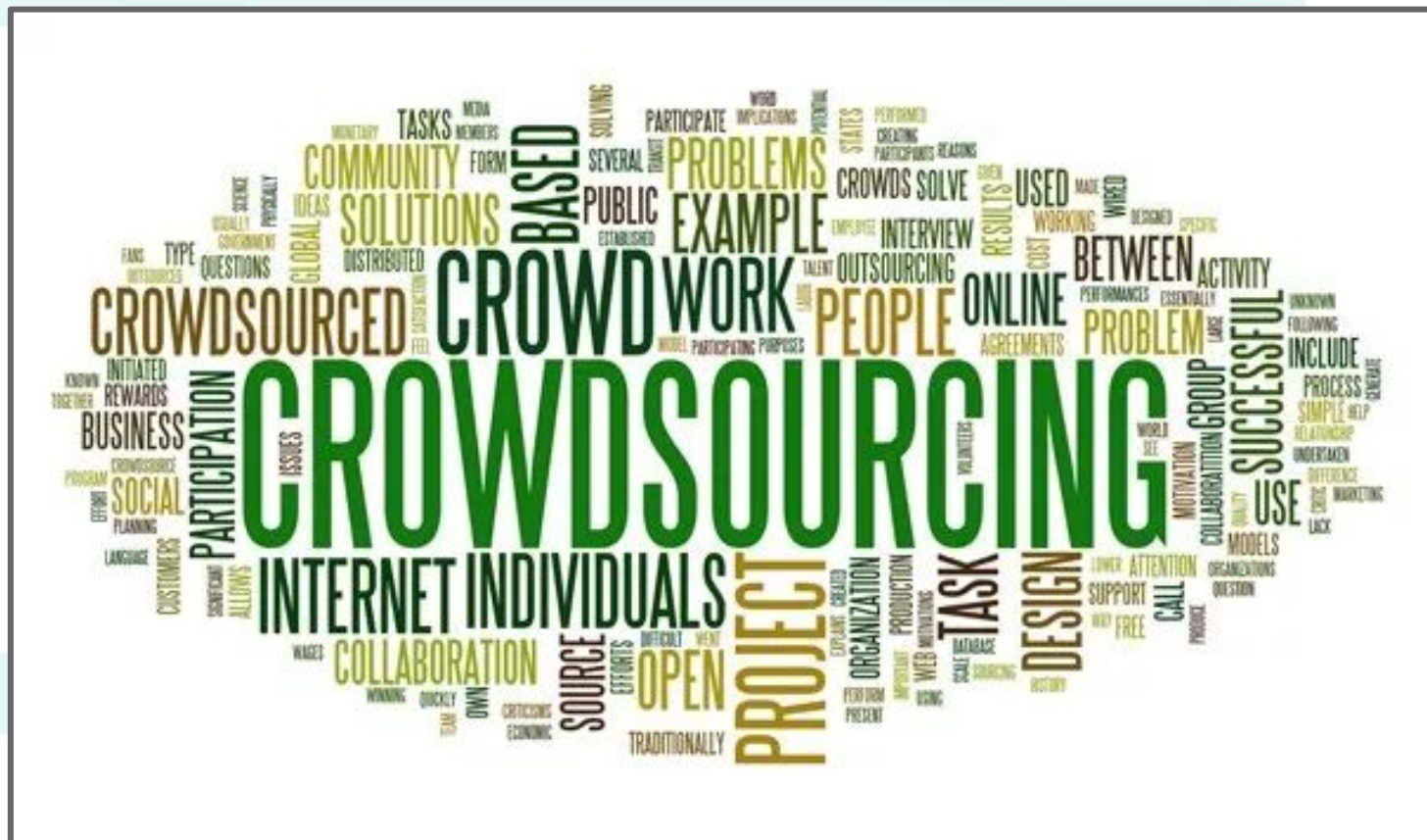
Task Analyzer

Task Title:

Task Description:

Task Resources:

Analyze



Outline

- ◆ Research Motivation
- ◆ Research Aim
- ◆ Crowdsourcing at the level of organizations
- ◆ Survey to understand confidentiality
- ◆ Unboxing a typical crowdsourcing task
- ◆ Understanding conversations between workers and task posters
- ◆ Protecting confidentiality loss in crowdsourcing
- ◆ Conclusion

Shift Towards a Gig-Economy: Benefits

Cost
Effective



Access to
a diverse
talent pool

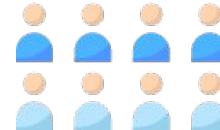


Faster
delivery



The fact that there is this quality of engineers all over the world, it really made me re-think how I resource the projects.

Marc Kocher
Sr. Manager, Engineering

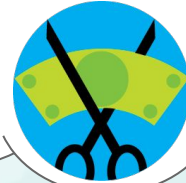


Benefits of Crowdsourcing

Cheaper and faster
60% of the times



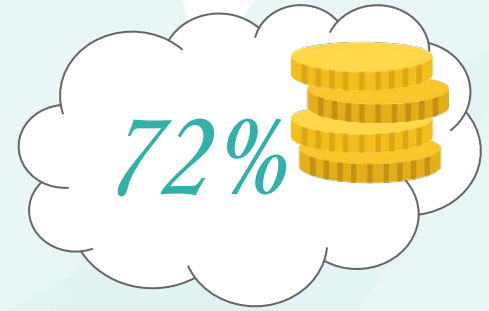
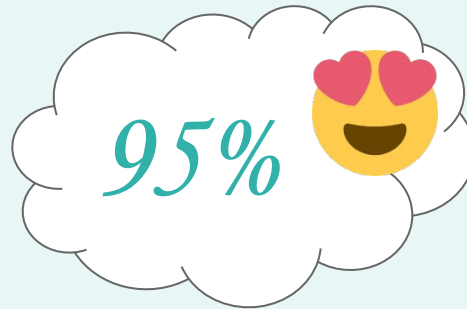
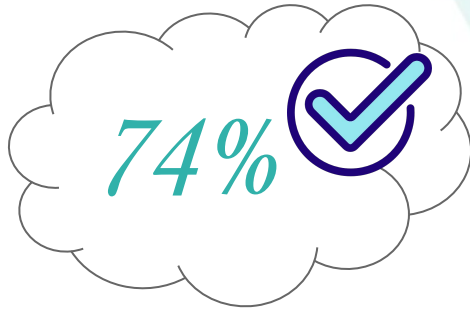
Reduces cost by
30 - 80%



Reduced **time-to-market**
because of Parallelism



Happiness and Satisfaction Levels



Crowdsourcing as the Next Big Revolution



By **2020**, freelancers
will form **43%** of the
workforce in USA



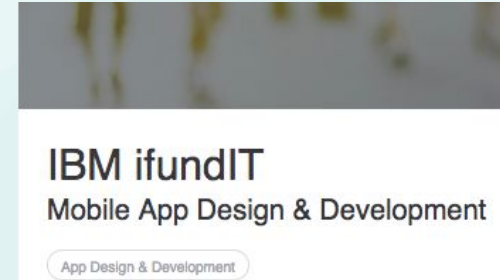
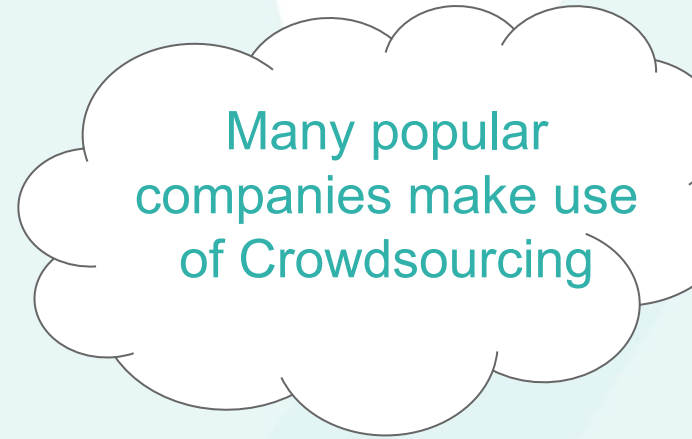
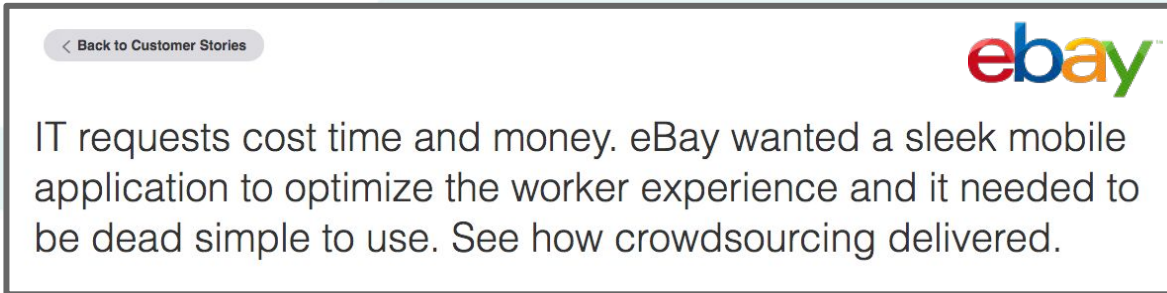
[HOME](#) [SERVICES](#) [NEWS](#) [EDUCATION](#) [ABOUT US](#)

Search

Intuit Forecast: 7.6 Million People in On-Demand Economy by 2020

QuickBooks Survey Reveals New Era of Entrepreneurship

Crowdsourcing in Organizations



Crowdsourcing and Software Development

Coding



Data
Analysis



Designing,
Prototyping



Web
Development



Testing



Mobile
Development



Crowdsourcing Platforms

upwork™

Find Freelancers

Web Dev Mobile Dev Design Writing Admin

Get it done with
freelancer

Grow your business with the top freelancers

What type of work do you need?



MyCrowd QA
A QASource Company

Sign Up Now

Home Features

MyCrowd Studio is an extraordinary platform
for crowdtesting websites and mobile apps



Startups and DIY



Large Ecommerce



Managed Tests and
Enterprise Solutions

topcoder™

Deliver faster for your business
through crowdsourcing.

With a community of over 1,000,000 design and
technology experts, Topcoder provides on-demand
capability, bandwidth, and velocity so you can do more.

I want to get work done

I want to join Topcoder

And many more..

Stakeholders in Crowdsourcing

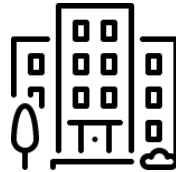
Task Poster -> ***tp***



Worker -> ***w***



Company/Organization -> ***c***

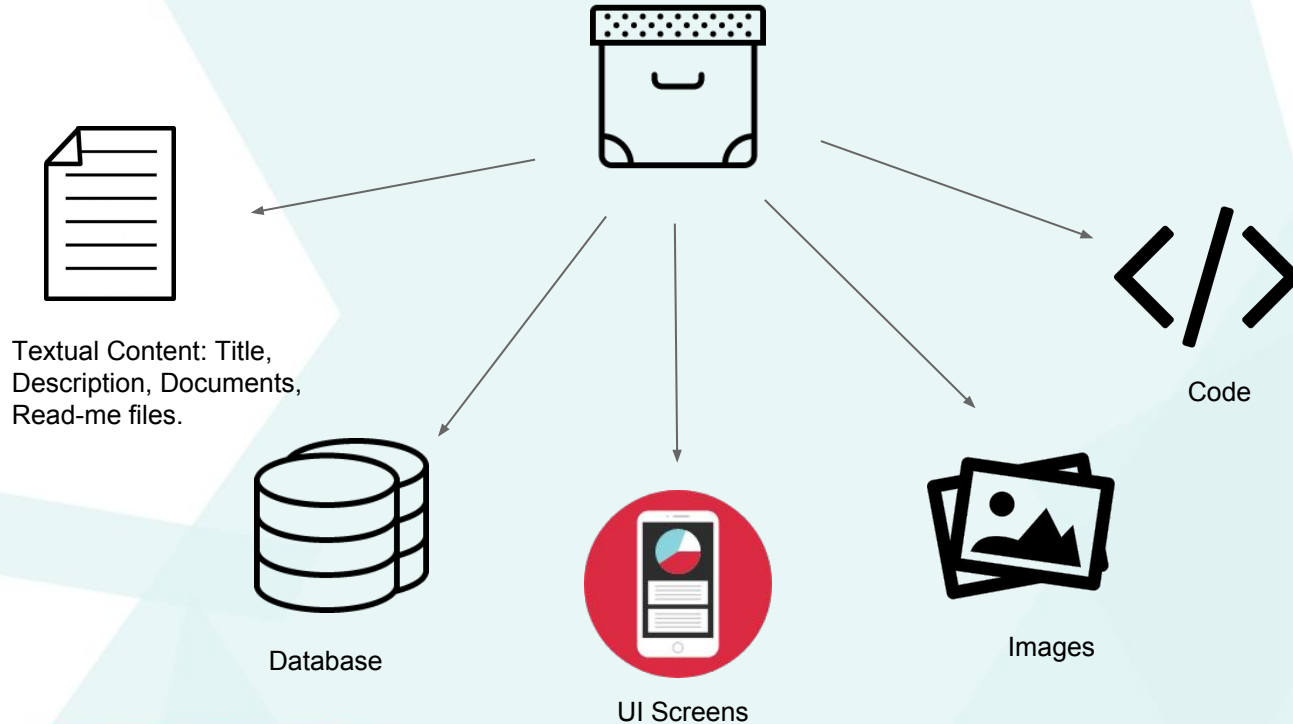


Components of a Crowdsourcing Task

◆ ~ 20 attributes per task

- Task ID
- Task Title
- Task Description
- Task Category
- Task Type
- Task Workload
-

Resources shared commonly



Sample Software Development Task ^[1]

Password generator with barcode

Desktop Software Development

Hi,

I'm looking for a JAVA software that can generate strong password with barcodes.

The software should allow me to set:

- password length (max 40)
- options to include symbols i.e @#\$%
- options to include numbers
- options to include lowercase characters and upper case characters.

And this software cannot produce the same password.

Project Stage: N/A

Operating systems: Windows

Hours to be determined
Hourly

Less than 1 week
Project Length

Intermediate Level
I am looking for a mix of experience and value

June 21, 2017
Start Date

Post a job like this

Submit a proposal

About the Client

★★★★★ (5.00) 5 reviews

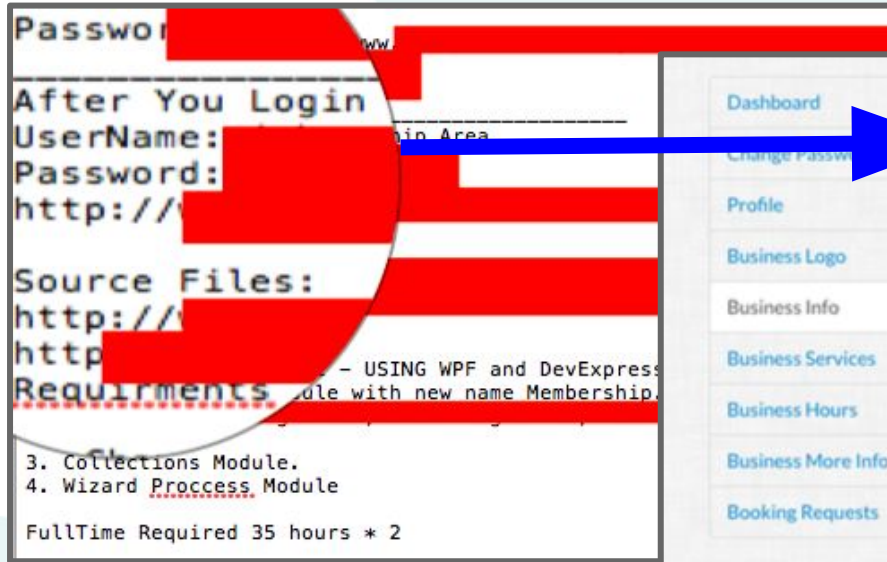
Canada
Montreal 05:46 AM

22 Jobs Posted
60% Hire Rate, 1 Open Job

\$5k+ Total Spent
13 Hires, 0 Active

\$12.93/hr Avg Hourly Rate Paid
74 Hours

Why is it a problem?



(A)

Dashboard

Change Passw

Profile

Business Logo

Business Info

Business Services

Business Hours

Business More Info

Booking Requests

All fields are optional, except mentioned.

Business Name required

Slogan

Email required

You will receive booking requests at this email address.

Phone

You will receive booking notifications at this number, if purchased.

Mobile required

You will receive booking notifications via SMS at this number, if purchased.

Address

Suburb

State / Territory

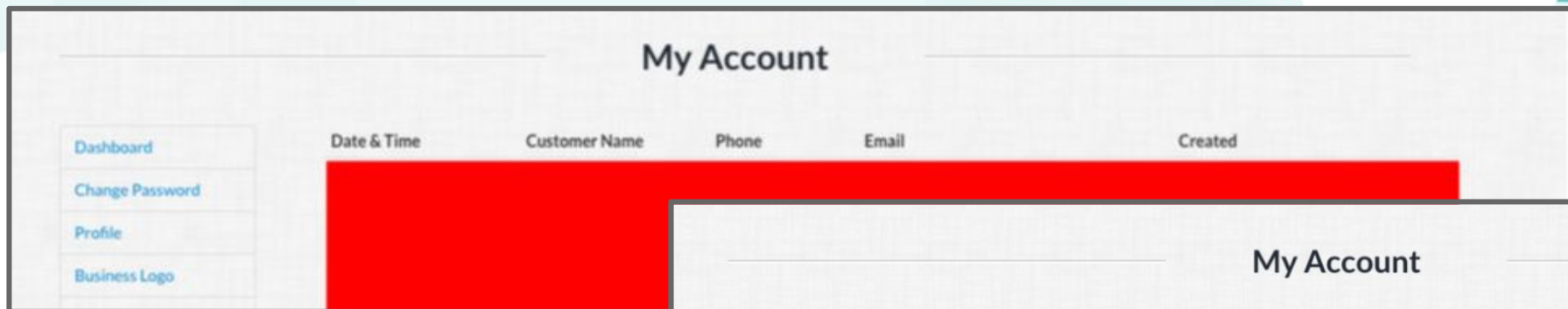
Post Code

Country

Website

Welcome Message

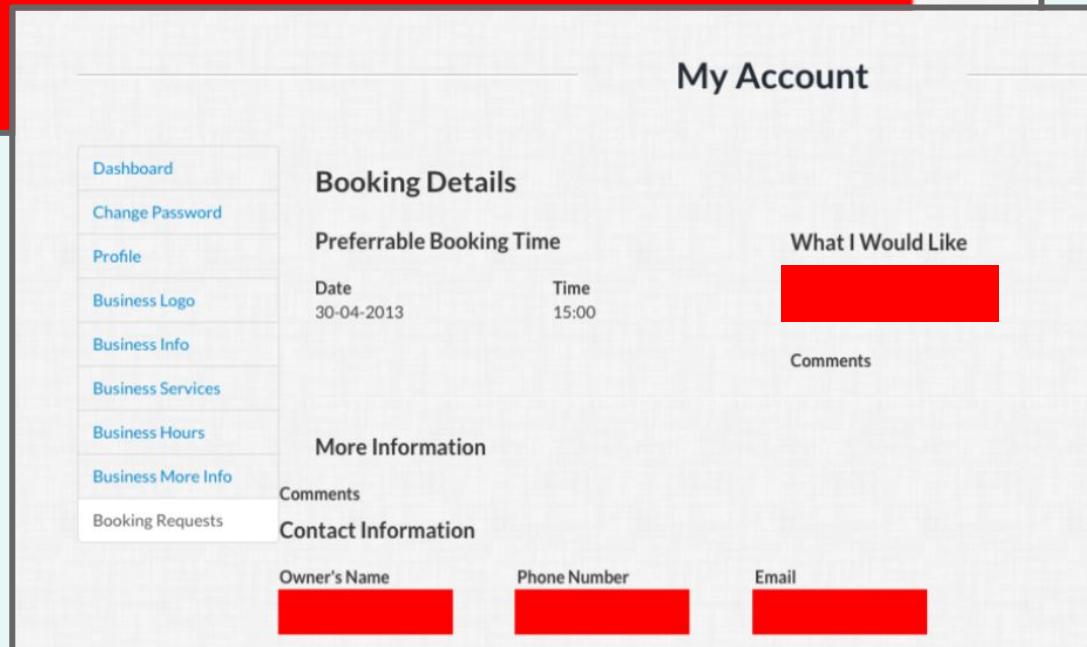
(B)



(A)

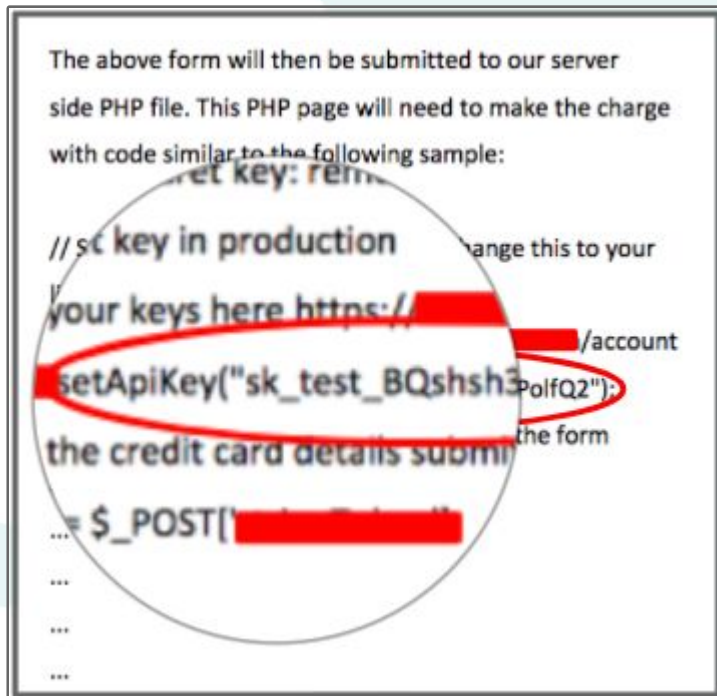
Exposed data of the Customers:

- Name
- Address
- Email ID
- Phone number
- Details of availed services
 - Date
 - Time
 - Cost



(B)

Why is it a problem?



(A)

My \$2375 Amazon EC2 Mistake 📌

Jan 7th 2015 in Commentary ©

A word of warning: Know what your modules/extensions/pods/plugins are doing, **especially** if they use any of your credentials.

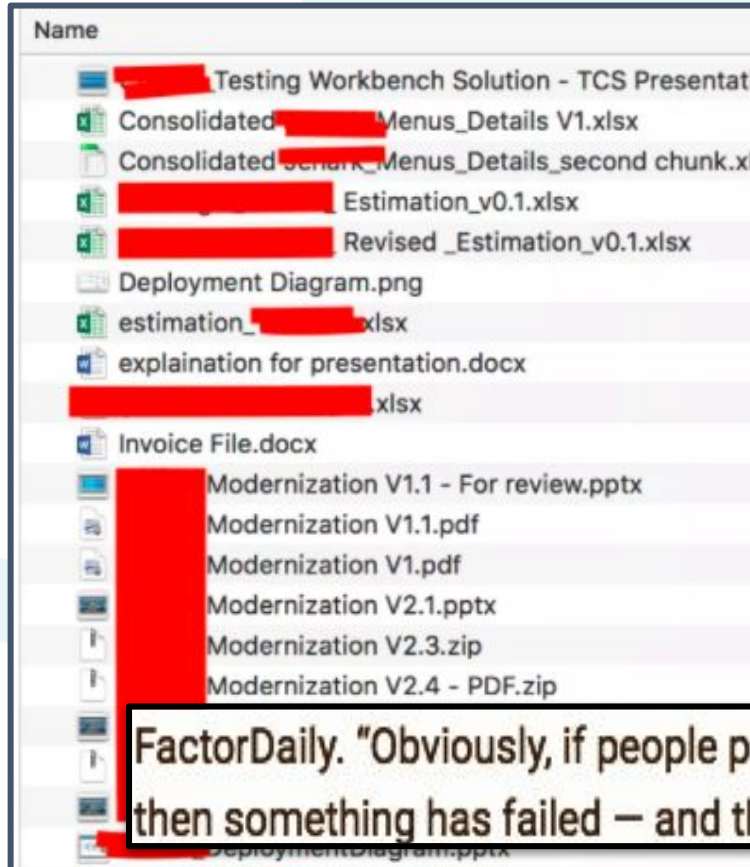
I'm actually surprised that this actually was up that long. I accidentally did this once and Amazon was on the phone with me 10 mins later.

Turns out through the S3 API you can actually spin up EC2 instances, and my key had been spotted by a bot that continually searches GitHub for API keys. Amazon AWS customer support informed me this happens a lot recently, hackers have created an algorithm that searches GitHub 24 hours per day for API keys... Once it finds one it spins up max instances of EC2 servers to farm itself bitcoins...

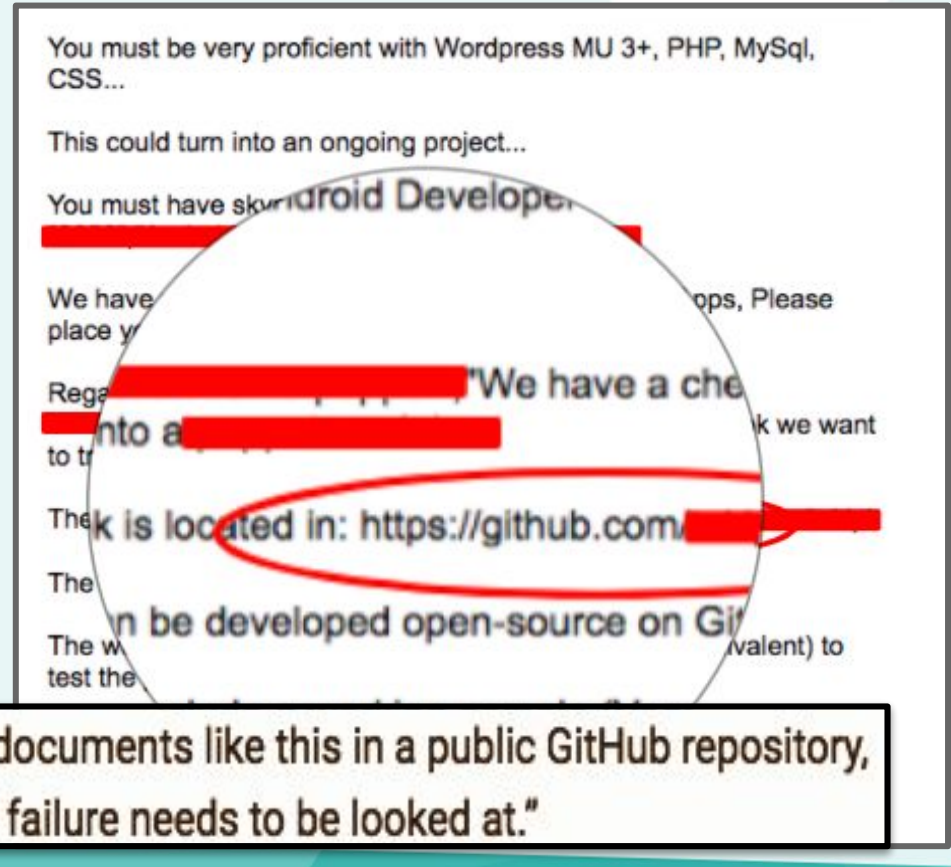
Boom! A \$2375 bill in the morning. Just for trying to learn rails.

(B)

Why is it a problem?



FactorDaily. "Obviously, if people put documents like this in a public GitHub repository, then something has failed – and that failure needs to be looked at."



(A)

(B)

20

Impacts of Exploited Vulnerability

- ◆ Economic Exploitation
- ◆ Information Theft
- ◆ Intrusion on personal privacy
- ◆ Social engineering
- ◆ System penetration/attack
- ◆ Information Bribery
- ◆ Sale of personal information
- ◆ System sabotage
- ◆ Unauthorized system access
- ◆ CIA Loss

Research Aim:

Given Crowdsourced Software Development:

- ◆ Identify
 - stages involved in confidentiality loss
 - critical attributes of a task
 - type of critical data
- ◆ Develop techniques to analyze tasks



Outline

- ◆ Research Motivation
- ◆ Research Aim
- ◆ Crowdsourcing at the level of organizations
- ◆ Survey to understand confidentiality
- ◆ Unboxing a typical crowdsourcing task
- ◆ Understanding conversations between workers and task posters
- ◆ Protecting confidentiality loss in crowdsourcing
- ◆ Conclusion

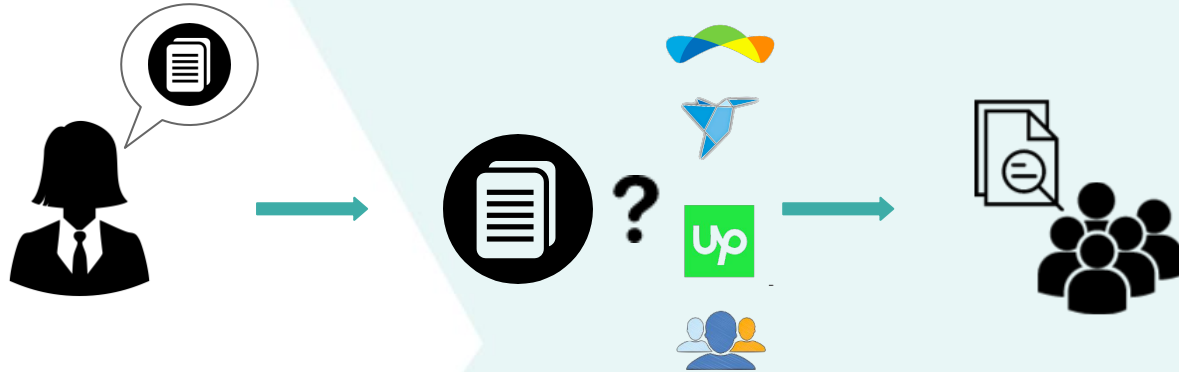
The Crowdsourcing Cycle



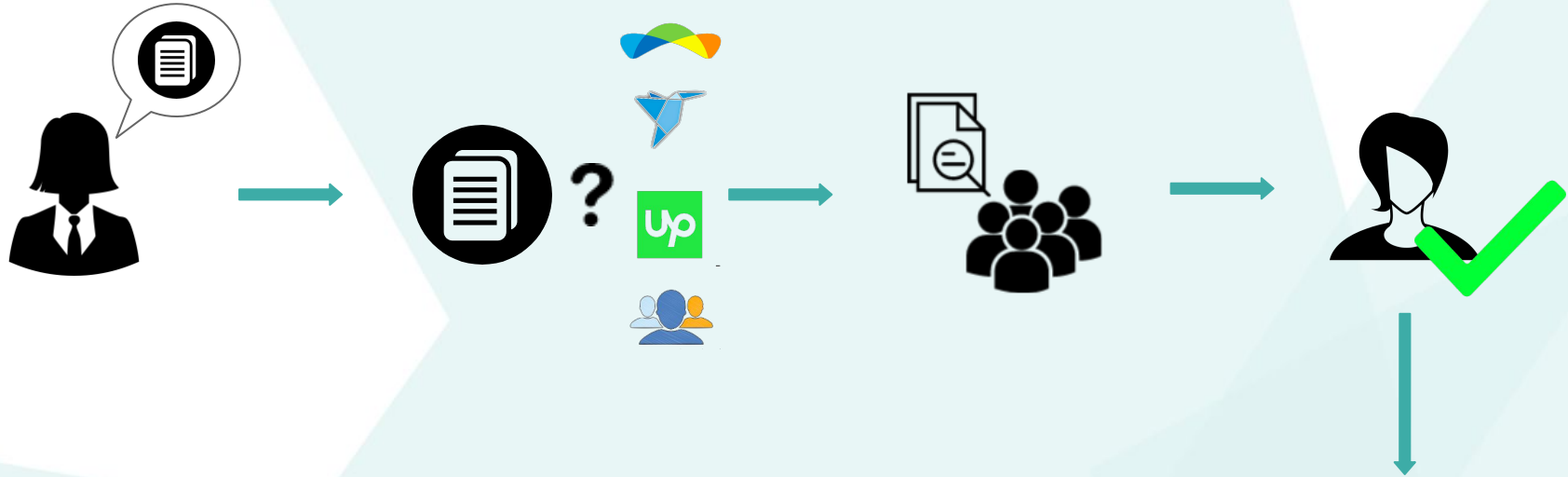
The Crowdsourcing Cycle



The Crowdsourcing Cycle



The Crowdsourcing Cycle



The Crowdsourcing Cycle



The Crowdsourcing Cycle



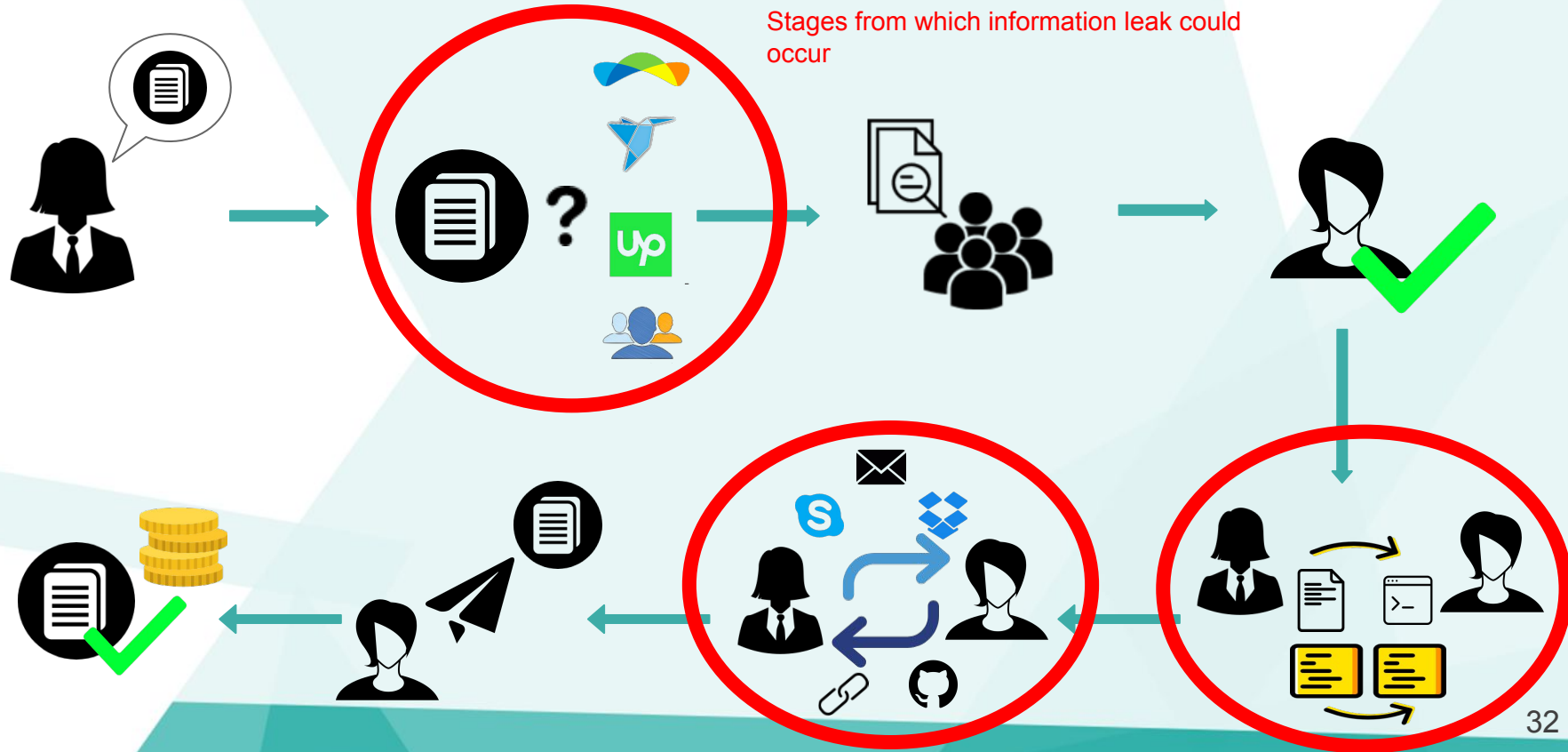
The Crowdsourcing Cycle



The Crowdsourcing Cycle



The Crowdsourcing Cycle



Survey to validate some assumptions

- ◆ PII details are sensitive
- ◆ Sharing a database publicly may compromise with confidentiality
- ◆ Sensitive information in code
- ◆ Sensitive components in the comments



Survey to validate some assumptions

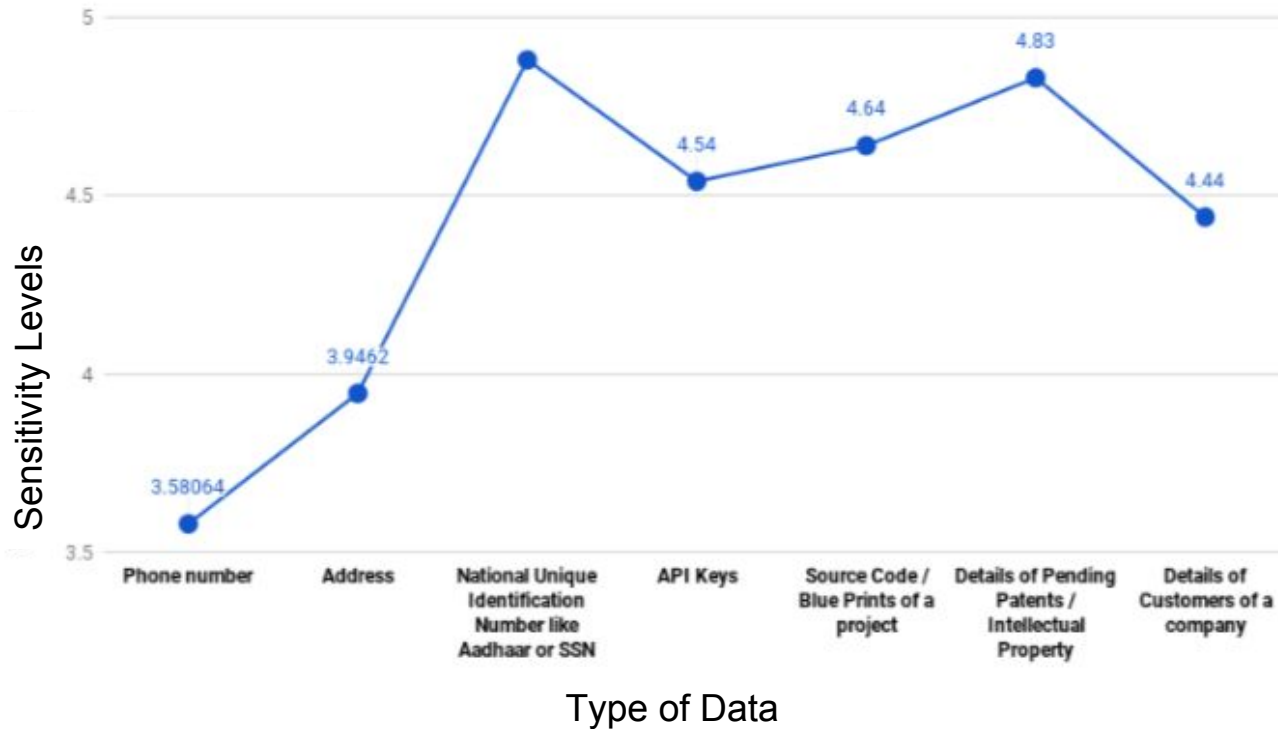
- ◆ Wireframes and UX examples may give out information about the company posting the job: logo
- ◆ Market of the end result of the task can be inferred given a task
- ◆ Domain of the task can be identified from the task, given the type details it is asking for

Survey to understand Confidentiality

Suppose you are a project manager in a multinational organization and you need to get some job done by an external freelancer. These jobs require sharing some details like sample code, documents, UX wireframes, etc. and some information might be sensitive in these resources. You need to answer the following questions considering the above scenario.

Likert scale based Questions

Average Sensitivity Levels of Sharing Various Details



Some observations

- ◆ 6.4% of the respondents believed that sharing a database publicly is not a sensitive activity
- ◆ 2.2% people think that sharing API Keys is okay

Situational Questions

Sample Task Question

This is a sample task posted by a Task Poster, according to you, what all sensitive data does it contain?

I am using a online platform called tomakeanapp.com, there is a module which you can upload a picture but I HAVE TO DO IT one by one.

I need a imacro or anyscript that could upload a batch of photo from a folder. No other requirement. That's it. It is dead simple and first come first serve. Please apply

If you want to know what i am saying.

Go to testmyapp.com
Username: trialappdemo
Password: password134

>Click my project> Click "Edit", then you will see "TraiApp" in YOUR ACTIVITY

Then click "Edit" and you will see the upload screen I talked above

Sample Task Question

This is a sample task posted by a Task Poster, according to you, what all sensitive data does it contain?

I am using a online platform called tomakeanapp.com, there is a module which you can upload a picture but I HAVE TO DO IT one by one.

I need a imacro or anyscript that could upload a batch of photo from a folder. No other requirement. That's it. It is dead simple and first come first serve. Please apply

If you want to know what i am saying.

Go to ~~testmyapp.com~~

Username: trialappdemo

Password: password134

>Click my project> Click "Edit", then you will see "TraiApp" in YOUR ACTIVITY

Then click "Edit" and you will see the upload screen I talked above

Sample Code Question

Following is a sample code snippet, point out the parts in the code that you think are sensitive:

```
import pynder

# Replace the given Facebook username and password with yours to get your access token
# My username -> jakemichael@gmail.com
# My password -> ThisIsMyPass123
fb_username = ''
fb_password = ''
token = get_access_token(fb_username, fb_password)

#put in your FB ID here
fb_id = '7003251371123'

# Replace with api_key
api_key = 'qpefsey31cs3c24dtahd7sbfs5fshs5w3jESFVB6527ARVJ31LJsglsjsajKbsgdbh426bFvsh'

# To process the paypal payment for further authentication to get access to unlimited API calls use my
# Account Number and IPin for payment at www.paypal.com
# Account Number: GH1728K
# IPin: angela52

session=pynder.Session(facebook_token=fb_auth_token, facebook_id=fb_id)

a = 0
for user in session.nearby_users():
    print a
    a += 1
    try:
```

Sample Code Question

Following is a sample code snippet, point out the parts in the code that you think are sensitive:

```
import pynder

# Replace the given facebook username and password with yours to get your access token
# My username -> jakemichael@gmail.com
# My password -> ThIsIsMyPass123
fb_username = ''
fb_password = ''
token = get_access_token(fb_username, fb_password)

#put in your FB ID here
fb_id = '7003251371123'

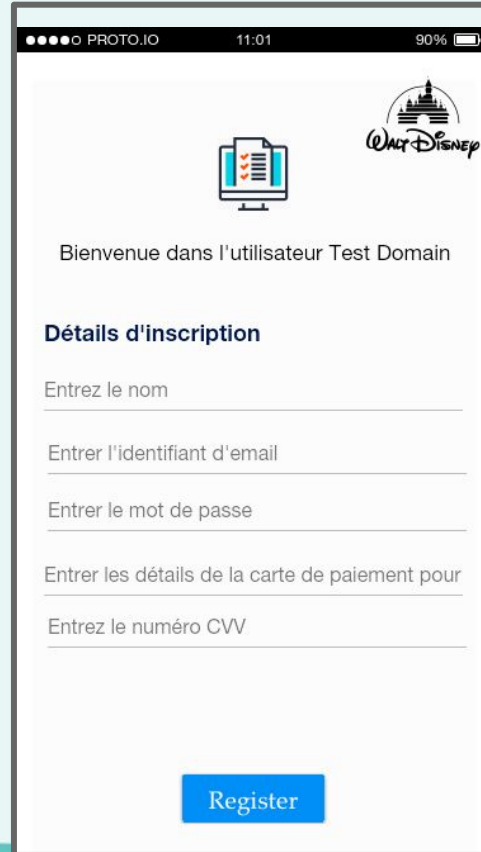
# Replace with apli_key
apli_key = 'qpefsey31cs3c24dtahd7sbfs5fshs5w3jESFVB6527ARVJ31LJsglsjsajKbsgdbh426bFvsh'

# To process the paypal ppayment for further authentication to get access to unlimited API calls use my
# Account Number and IPin for payment at www.paypal.com
# Account Number: GH1728K
# IPin: angela52

session=pynder.Session(facebook_token=fb_auth_token,facebook_id=fb_id)



a = 0
for user in session.nearby_users():
    print a
    a += 1
    try:
```

A wireframe of an App has been attached, what all can you infer from it?



The wireframe shows a mobile app interface for registration. At the top, the status bar displays 'PROTO.IO', '11:01', and '90%'. The app header features a document icon with checkmarks and the 'Walt Disney' logo. The main content area includes a welcome message, a section title, and five input fields for user details, followed by a 'Register' button.

PROTO.IO 11:01 90%

Bienvenue dans l'utilisateur Test Domain

Détails d'inscription

Entrez le nom

Entrez l'identifiant d'email

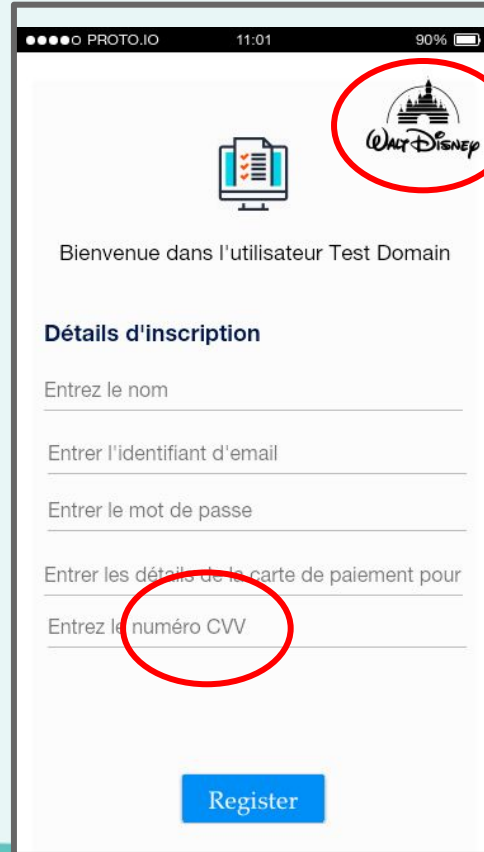
Entrez le mot de passe

Entrez les détails de la carte de paiement pour

Entrez le numéro CVV

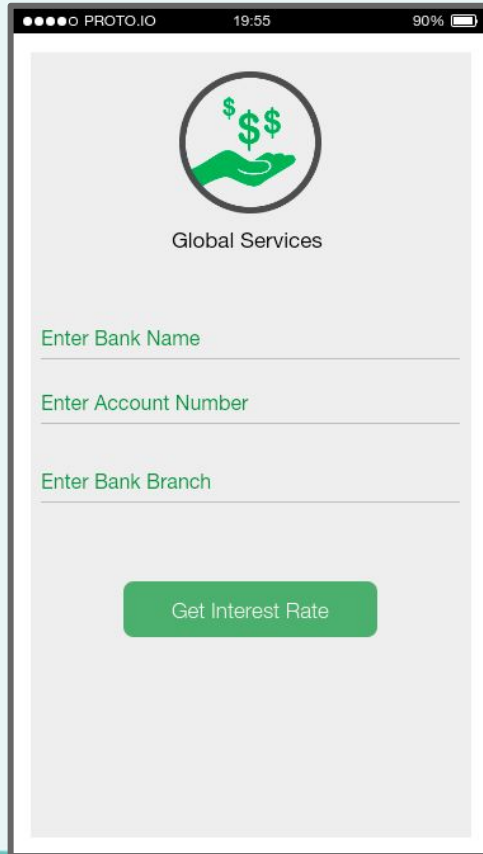
Register

A wireframe of an App has been attached, what all can you infer from it?



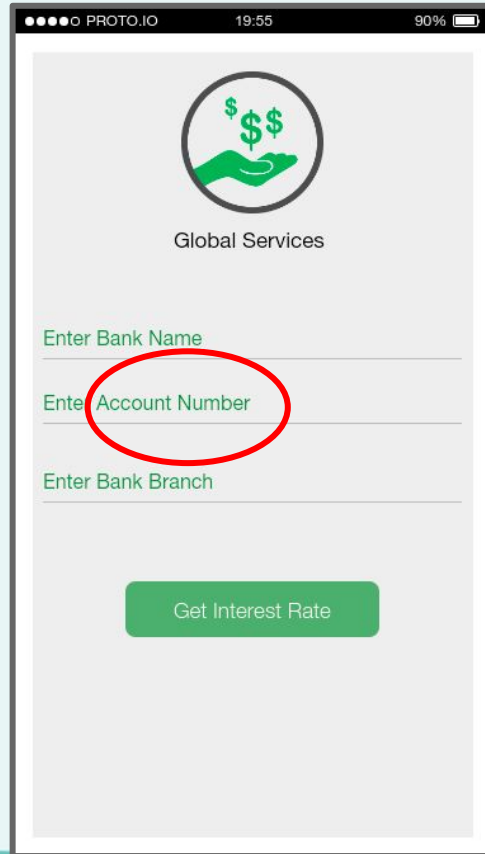
The wireframe shows a mobile app interface for registration. At the top, the status bar displays 'PROTO.IO', '11:01', and '90%'. Below the status bar, there is a header area with a document icon on the left and a 'Walt Disney' logo on the right, which is circled in red. The main content area starts with the text 'Bienvenue dans l'utilisateur Test Domain'. Below this is a section titled 'Détails d'inscription'. This section contains five input fields: 'Entrez le nom', 'Entrez l'identifiant d'email', 'Entrez le mot de passe', 'Entrez les détails de la carte de paiement pour', and 'Entrez le numéro CVV'. The 'Entrez le numéro CVV' field is circled in red. At the bottom of the form is a blue button labeled 'Register'.

What can you infer about the type of industry that this app is for?



The image shows a mobile application interface on a smartphone screen. The status bar at the top displays 'PROTO.IO', the time '19:55', and a battery level of '90%'. The app's header features a circular logo with a green hand holding three dollar signs, with the text 'Global Services' centered below it. The main content area contains three input fields, each with a green label and a horizontal line for text entry: 'Enter Bank Name', 'Enter Account Number', and 'Enter Bank Branch'. At the bottom of the form is a green button with the text 'Get Interest Rate'.

What can you infer about the type of industry that this app is for?



The image shows a mobile application interface on a screen. At the top, the status bar displays 'PROTO.IO', '19:55', and '90%' battery. The app's header features a logo of a green hand holding three dollar signs, with the text 'Global Services' below it. The main form area contains three input fields: 'Enter Bank Name', 'Enter Account Number' (which is circled in red), and 'Enter Bank Branch'. At the bottom of the form is a green button labeled 'Get Interest Rate'.

Other Confidential Entities

1. Master Database of a company
2. Personal Information:
 - a. Name
 - b. Email ID
 - c. SSN / Unique Identification Number
 - d. Address
 - e. Date of Birth
 - f. Employee ID
 - g. Health related information
3. Context for the given task
4. Company related details:
 - a. Details of the clients of a company
 - b. Details of a company
 - c. Details of a company's future projects
 - d. Any sort of financial data
 - e. Proprietary Data

Other Confidential Entities

7. Credentials of any sort: Usernames and Passwords
8. Bank Account Number / Credit or Debit Card Details
9. API Keys
10. Server Details
11. Google Play keys / App Secrets / OAuth Tokens
12. Named Entities
13. Variable names / Class Names
14. Test Inputs
15. Sensitive Datasets where the users can be uniquely identified
16. Images and Diagrams with confidential information

Dataset of 65,000 tasks

- ◆ 4 unique Task Categories:
 - Data Science & Analytics
 - IT & Networking
 - Web, Mobile & Software Dev
 - Writing

- ◆ 22 Unique Task Subcategories
 - Game Development
 - Data Visualization
 - Network and System Administration
 - Machine Learning
 - ...

Analysis of attributes of a task

◆ ~ 20 fields per task

- Task ID
- Task Title
- Task Description
- Task Category
- Task Sub-Category
- Date of creation of Task
- Task Type
- Task Workload
- Task Duration
- Task Budget
- Preferred Location of the worker
- Preferred range of feedback score
- Level of English Skills required
- Payment range
- Range of working hours per day
- Requirement of Cover Letter
- Requirement of Portfolio
- Candidates Registered for the task
- Skills required for the task

Analysis of Tasks

- ◆ Crucial attributes:
 - Task Title
 - Task Description
- ◆ Observations:
 - **23%** tasks contained links
 - **2.2%** tasks had credentials
 - **2.46%** tasks had GitHub links

Outline

- ◆ Research Motivation
- ◆ Research Aim
- ◆ Crowdsourcing at the level of organizations
- ◆ Survey to understand confidentiality
- ◆ Unboxing a typical crowdsourcing task
- ◆ Understanding conversations between workers and task posters
- ◆ Protecting confidentiality loss in crowdsourcing
- ◆ Conclusion

Understanding Conversations between Workers and Task Posters

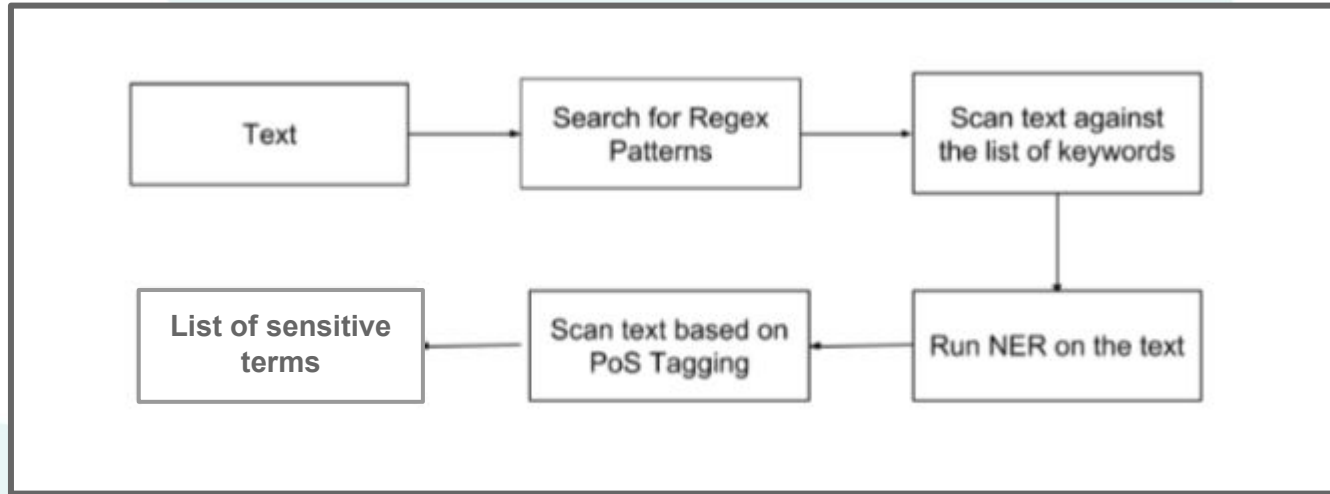
- ◆ ~ 214 lines per task (average)
- ◆ ~30 files exchanged [min: 0, max: 90]
- ◆ 42% tasks involved Skype calls
- ◆ 14% involved sharing data over cloud



Analysis of a Task for Sensitivity



Pipeline



Algorithm

Algorithm 1 Analysis of a task

Result: Array of keywords found sensitive in t

Given t and its components $tt, td, r_1, r_2, r_3, \dots, r_n$;

res = [$tt, td, r_1, r_2, r_3, \dots, r_n$];

for each item in res **do**

 tokenize(res)

 removeStopWords(res)

 detectRegex(res)

 detectPasswordsAndAPIKeys(res)

 detectSensitiveKeywords(res)

 detectNER(res)

 add sensitive terms detected to result

end

- Task Title tt
- Task Description td
- Task Resources r_1, r_2, r_3, \dots

Algorithm 2 Detecting Password and API Key like instances

Result: Return a list of terms identified

list = []

for each *term* **do**

if *term is not an English Word* **then**

 | list.append(term)

end

if *term is a digit and length(term) > 3* **then**

 | list.append(term)

end

if *term is alphanumeric* **then**

 | list.append(term)

end

if *term contains special characters* **then**

 | list.append(term)

end

if *term is in CamelCase and has >1 Capital Letters* **then**

 | list.append(term)

end

 Return list

end

Testing

Induced sensitivity in the 65,000 tasks dataset

- ◆ Curated a list of sensitive terms:
 - Names of companies
 - Email Addresses from the Enron Dataset
 - Dummy API Keys, Passwords, etc
 - List of URLs
 - List of countries and cities
 - Randomly generated usernames
 - ...

Performance Metrics

$$\text{Precision} = \frac{\text{Number of sensitive terms correctly identified}}{\text{Number of sensitive terms identified by the algorithm}}$$

$$\text{Recall} = \frac{\text{Number of sensitive terms correctly identified}}{\text{Number of sensitive terms in the ground truth}}$$

Precision	Recall	F1 Score
0.68	0.82	0.74

Real World Application

◆ Analyze

- Tasks
- Content before putting on cloud
- Content before sharing

Conclusion

- ◆ Narrow down on the crucial stages in the crowdsourcing cycle
- ◆ Enumerate critical attributes in a task
- ◆ Highlight type of confidential data
- ◆ NLP and Rule based algorithm to detect confidentiality loss

Challenges, Limitations and Future Work

- ◆ Dataset availability
- ◆ Lack of labeled data
- ◆ Expand the algorithm
 - Cater to images, databases, etc
- ◆ Incorporate Machine Learning techniques for classification
- ◆ Sanitization
- ◆ More fine grained analysis

Acknowledgement

- ◆ Committee Members
- ◆ Abhinav and Sakshi, Accenture Labs Bangalore
- ◆ Indira, Gurpriya, Arpit, Shubham, Sonu, Divyansh
- ◆ Members of Precog family
- ◆ Family and friends

References

- ◆ <http://www.timbroder.com/2015/01/my-2375-amazon-ec2-mistake.html>
- ◆ <https://www.topcoder.com/about-topcoder/customer-stories/>
- ◆ <https://www.merriam-webster.com/dictionary/crowdsourcing>
- ◆ <https://www.figure-eight.com/>
- ◆ <https://www.flaticon.com/authors/freepik>



Thanks!



simran13104@iiitd.ac.in



[@simran_s21](https://twitter.com/simran_s21)