# Crowdsourcing Historical Tabular Data – 1961 Census of England and Wales

CHRISTIAN CLAUSNER, JUSTIN HAYES AND APOSTOLOS ANTONACOPOULOS

PATTERN RECOGNITION AND IMAGE ANALYSIS RESEARCH LAB, UK

HIP'19, SYDNEY, AUSTRALIA

# The 1961 Census Digitisation Project

- Millions of data items trapped in 100,000+ pages (tables)
- Main part of project in 2018/2019
  - For Office for National Statistics
- Automated processing pipeline
  - About 98% correct results
  - Requires post-correction
- Two other publications, this one is an experience paper concentrating on the crowdsourcing aspects

# Challenges

- ▶ Inconsistent scan quality (illumination, warping, skew, scaling, placement)

- ▶ Faint print, handwritten corrections

- ▶ Microfilm scratches and general degradation

- ▶ Missing parts, printing errors

- ▶ Unorganised data (pages not in any particular order)

- ▶ Dense tables, sometimes with no separation between columns

# Workflow

- Complete digitisation workflow from image to structured data in database
  - Simplified workflow in the right
- Validation of data is crucial
- Identify errors by
  - Visual checks
  - Automated crosschecks
- Manual intervention
  - In part in-house
  - Mostly by crowd

# Zooniverse

- ► We used Zooniverse for crowdsourcing
- ► Public platform (also open source)
- ► Big base of volunteers
- ► Free for projects that benefit the public good
- ► Easy to use
- ► Good support

https://www.zooniverse.org/

# Micro Tasks

- Task for volunteers as simple as possible

- "Enter text for highlighted table cell"
  - We don't even show the OCR result

- Problematic or unclear cases can be tagged (Talk section with hashtags)



Number of volunteers

Task complexity

# Census Zooniverse Project

- One of the most active projects in the time period
- No promotion
- Difficult to provide enough data

**Live Workflows**

**20190517_sc13**
Retirement limit: 3
Images retired: 7,423 / 7,423
ETC* 0 days

**20190517_sh06_wd**
Retirement limit: 3
Images retired: 7,517 / 7,517
ETC* 0 days

**20190517_sh0478_wd**
Retirement limit: 3
Images retired: 7,587 / 7,587
ETC* 0 days

100% Complete    100% Complete    100% Complete

**Number of classifications**

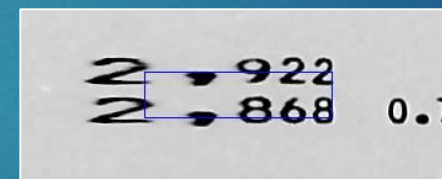| Month | Number of classifications |
|-------|---------------------------|
| Jul-18 | 792,129 |
| Aug-18 | 568,464 |
| Sep-18 | 524,245 |
| Oct-18 | 479,130 |
| Nov-18 | 579,422 |
| Dec-18 | 201,682 |
| Jan-19 | 302,043 |
| Feb-19 | 471,446 |
| Mar-19 | 664,131 |
| Apr-19 | 408,776 |
| May-19 | 513,463 |

# User Activity

# Great Participation

HIP'19
23/09/2019

- Large user base with auto-promotion of new/active/stagnant projects on Zooniverse

- High interest in historical projects (and UK)

- Micro-tasking (mindfulness?)

- User engagement

- Consistency in data provision

- Power users (special attention)

# Discussion

- Crowdsourcing was very successful for the Census 1961 project
- Accuracies
  - OCR about 98%
  - Cell recognition in total about 95%
  - Correctness after crowdsourcing about 99.5%
    - Rest corrected by expert

# Problems

- ▶ Malicious users
  - ▶ Needs vigilance from our side
  - ▶ Can be blocked from Zooniverse side
- ▶ Bugs in the Zooniverse platform
  - ▶ We had a nasty one where text entered by users was incomplete
  - ▶ Fast fix
- ▶ Problems with data upload at busy times
  - ▶ Need to work around it

# Conclusion

- ▶ Worth it

- ▶ Over 5 million corrections in a few months

- ▶ Volunteers liked it (even demanded more data)

- ▶ Possibly more to come in near future

# Questions?

- zooniverse.org/projects/dataliberation/1961-census
- primaresearch.org/publications