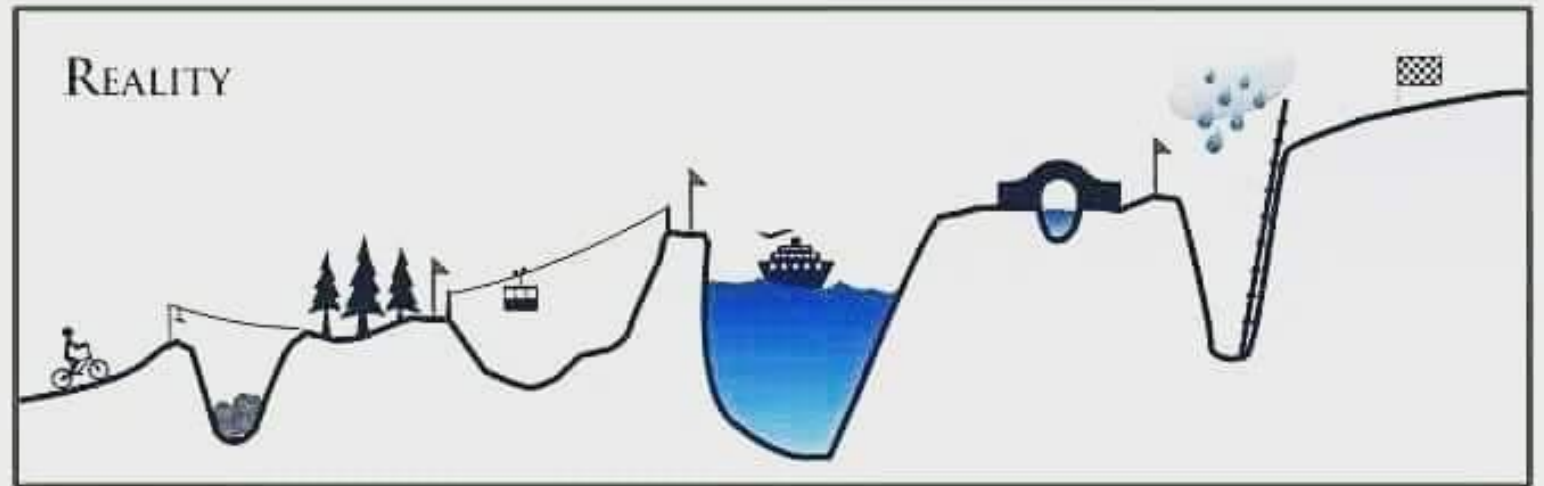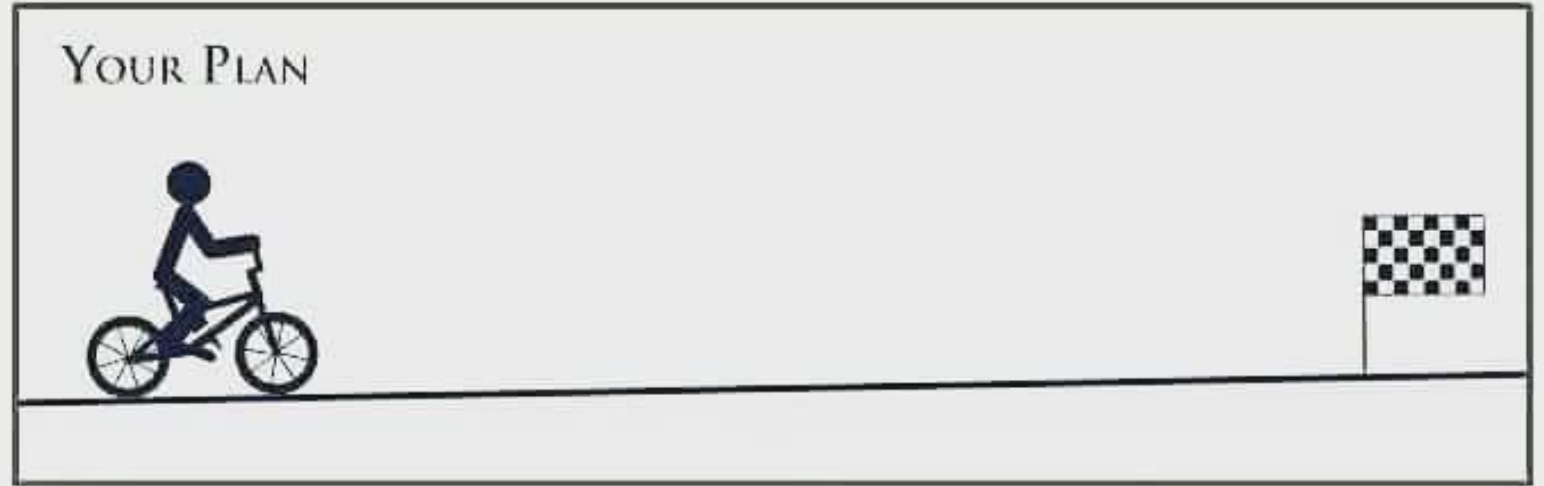# Issues in optical character recognition for statistical data capture

Arij Dekker

arijdekker@gmail.com

# There will be problems



YOUR PLAN

REALITY

# The choices

- Central keyboarding from paper
- Central scanning/reading from paper (OCR/ICR)
- Portable keyboard and store devices – no paper
- On-line enumeration or self-enumeration

# OCR: Mark sensing

# OCR: Number recognition

Form.No within HH

5. Completed Age

*If age greater than or equal to 98, write "98". If less than one write "00".*

In Years

# OCR: Working with text

Recognition

Computer Assisted Coding

# One or several forms?

## Single form:

- Cheaper
- Serves comfort
- Contractor choice

## Dedicated forms:

- Fewer "not applicables"
- Pre-numbering in first household sheet
- Improved control

# OCR versus manual capture: issues

- Duplicate questionnaire Id's
- Range errors
- Inconsistencies, within form
- Inconsistencies between forms

*Reduced "hands-on" operator involvement during data capture places particular requirements on questionnaire design for OCR*

# Line or column diversion

# Repeating person id & Correcting mistakes

# Duplicate household id's



*Duplicate household id's are the scourge of census files*

# Concluding remarks

- Involve IT staff in questionnaire design
- Test, test, test

- Prevent the risk of catastrophic failure