

## What does this program do?

This program extracts barcode and handwritten values from images.

### Template of Image files:

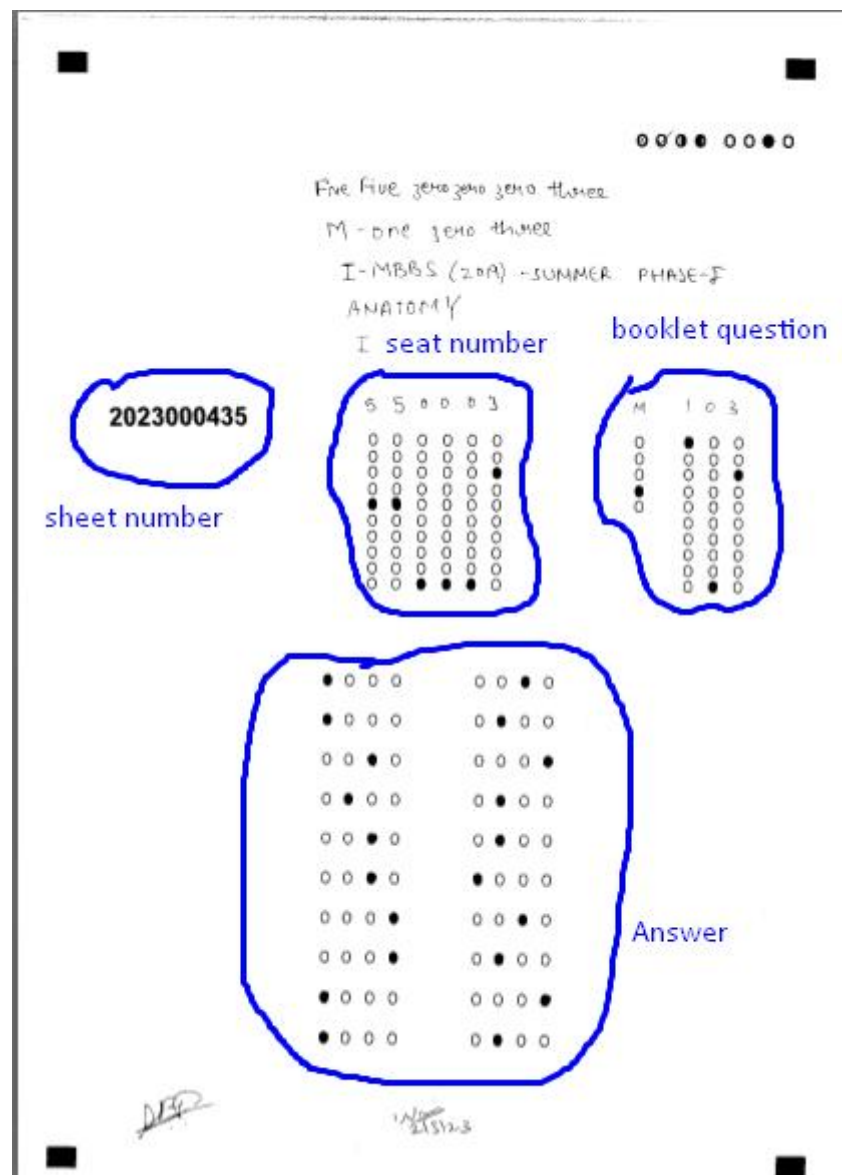


Fig.1 Template file for inputs of program

### Quality of files

- Include skewed some images
- Lots of Noise in getting total value.

## What is important in this program?

- Sheet number extraction
- Seat number extraction
- booklet question extraction
- Answer extraction.

## What engine does this program use ?

Opencv, tesseract, CNN

## How does this program work?

### Algorithm of this program

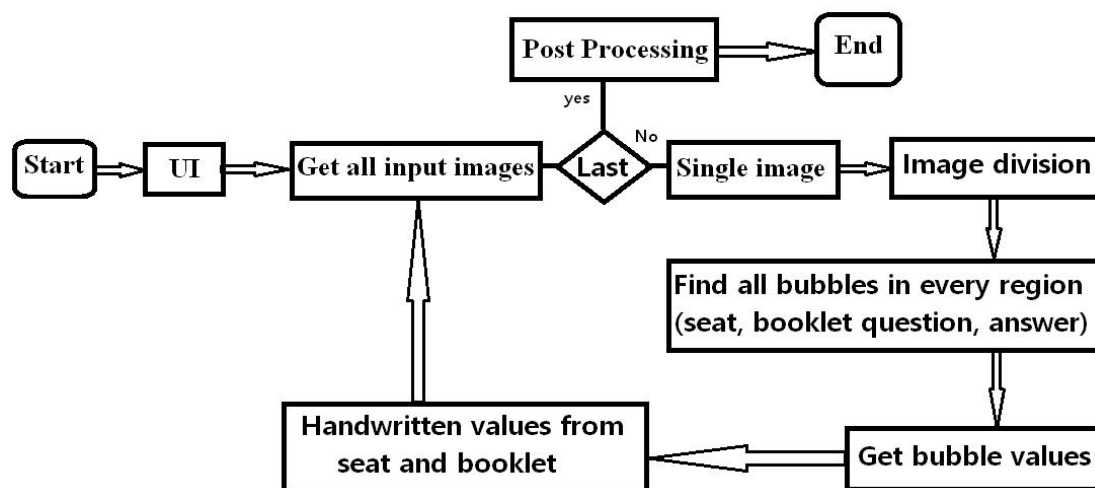


Fig. 2 Overall Algorithm

In code,

MCQ.py: take charge of parts of program starting, UI and several pre-processing such as getting input images.

main\_mcq.py: reflect logics of overall algorithm.



Booklet image:

M 1 0 3

0	●	0	0
0	0	0	0
0	0	0	●
●	0	0	0
0	0	0	0
	0	0	0
	0	0	0
	0	0	0
	0	0	0
	0	0	0
	●	0	

Answer image:

●	0	0	0	0	0	0	0	0	0
●	0	0	0	0	0	●	0	0	0
0	0	●	0	0	0	0	0	0	●
0	●	0	0	0	0	0	●	0	0
0	0	●	0	0	0	0	●	0	0
0	0	●	0	0	0	●	0	0	0
0	0	0	●	0	0	0	0	●	0
0	0	0	●	0	0	0	0	0	0
●	0	0	0	0	0	0	0	0	●
●	0	0	0	0	0	0	●	0	0

## Find all bubbles in every region

```

376     ### Finding all bubble in every region ###
377     id_loc, id_img, idW, idH = Match(template_path/f"0_2.jpg", id_img, threshold, 50)
378     q_loc, que_img, qW, qH = Match(template_path/f"0_2.jpg", que_img, threshold, 50)
379     a_loc, ans_img, aW, aH = Match(template_path/f"0_2.jpg", ans_img, threshold, 50)
380     ### seat number consideration ###
381     seat_num, seat_num_hand, que_num, que_num_hand = 'xxxxxx', 'xxxxxx', 'xxxx', 'xxxx'
382     try:
383         id_coors = getCoors(id_img, id_loc, idW, idH)
384         id_rows, id_cols = getRow_Col(id_coors, 1)
385         id_rows, id_cols = RowsColsCheck(id_rows, id_cols, 10, 6)
386         seat_num = SfindMark(id_rows, id_cols, id_img, idW, idH)
387         seat_num = ''.join(seat_num)
388         idCheckImg = id_img[id_rows[0]-120:id_rows[0]-40, id_cols[0]-30:id_cols[-1]+30]
389         sheetNumImg = id_img[id_rows[0]-120:id_rows[0]+40, 0:id_cols[0]-50]
390         seat_num_hand = next(Recognize_Digit(idCheckImg, "seat"))
391     except Exception as e:
392         print(e)

```

Function **Match()**, **getCoors()**, **getRow\_Col()**, **RowsColsCheck()** are used for this.

**Match()**: return all bubbles similar to "0\_2.jpg".

**getCoors()**, **getRow\_Col()**, **RowsColsCheck()** get exact bubble coordinates.

To find all bubbles in every region (seat, booklet, answer images), **Match()**-> **getCoors()**-> **getRow\_Col()**-> **RowsColsCheck()** procedures is applied.

And, cropped sheetnumber image and cropped seat image is gotten by seat number bubble coordinate.

**2023000435** 5 5 0 0 0 3

Also, cropped booklet image is gotten by booklet bubble coordinates.

M 1 0 3

Here, cropped sheetnumber image is not handwritten. This can be extracted by using pytesseract. Its accuracy will be high. But cropped seat image and booklet image can not be by using pytesseract because they are handwritten.

Therefore, sheet number value will be extracted by using tesseract.

## Get bubble values

Function SfindMark() is used.

Principle: Based on bubble coordinate (row and column coordinates), when average values of binary image by coordinate and calculated width and height is greater than certain threshold, it will be filled bubble. To get bubble value means finding filled bubbles.

## Handwritten values from seat and booklet.

Input images for this have been obtained in the part of “Find all bubbles in every region”.

To get a handwritten value of seat number, model.h5 trained in CNN is used as well as barcode handwritten digit detection.

And, to get the handwritten value of the Booklet image, two steps are needed.

- Handwritten alphabet recognition
- Handwritten digit recognition.

We can use model.h5 for handwritten digits.

And, for alphabet recognition, we use the pretrained model “Alphabet\_Recognition”.

## Post-processing

- This part makes xlsx file and json of output.