# Portfolio
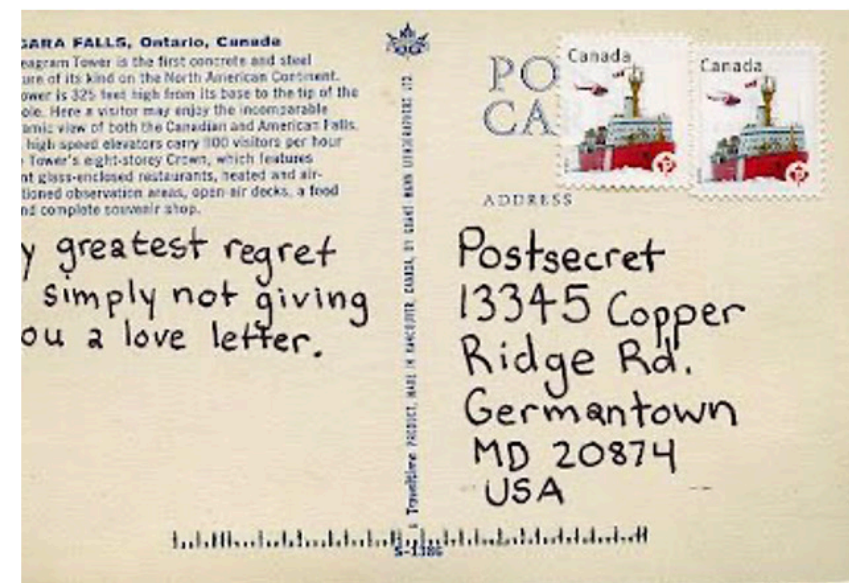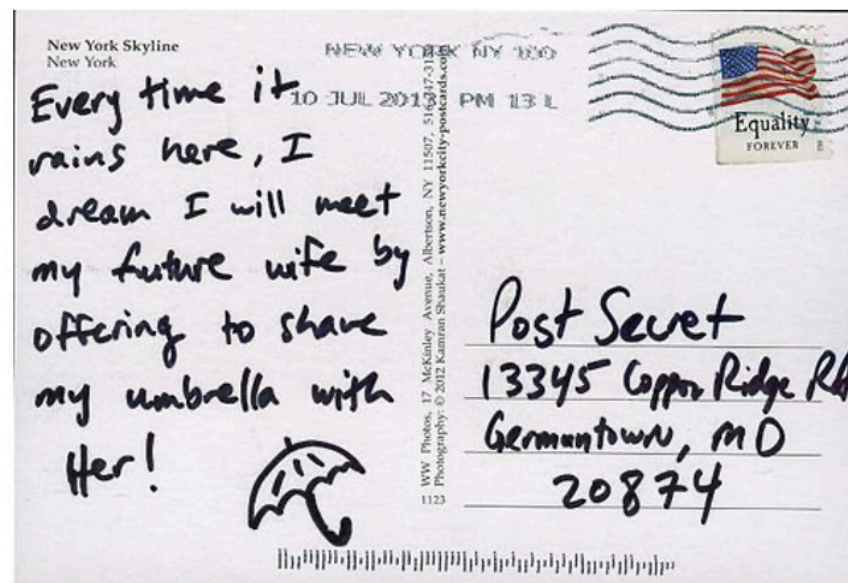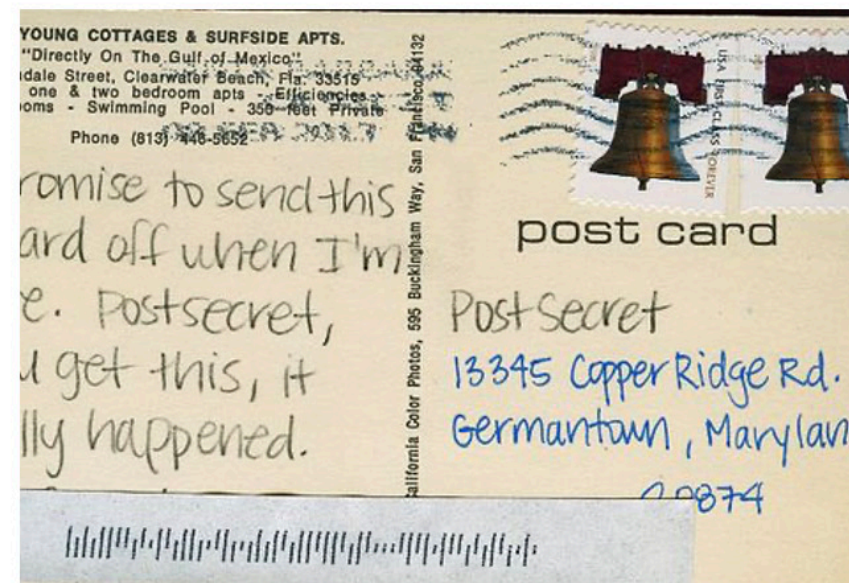
BIG DATA, THE SELF AND SOCIAL
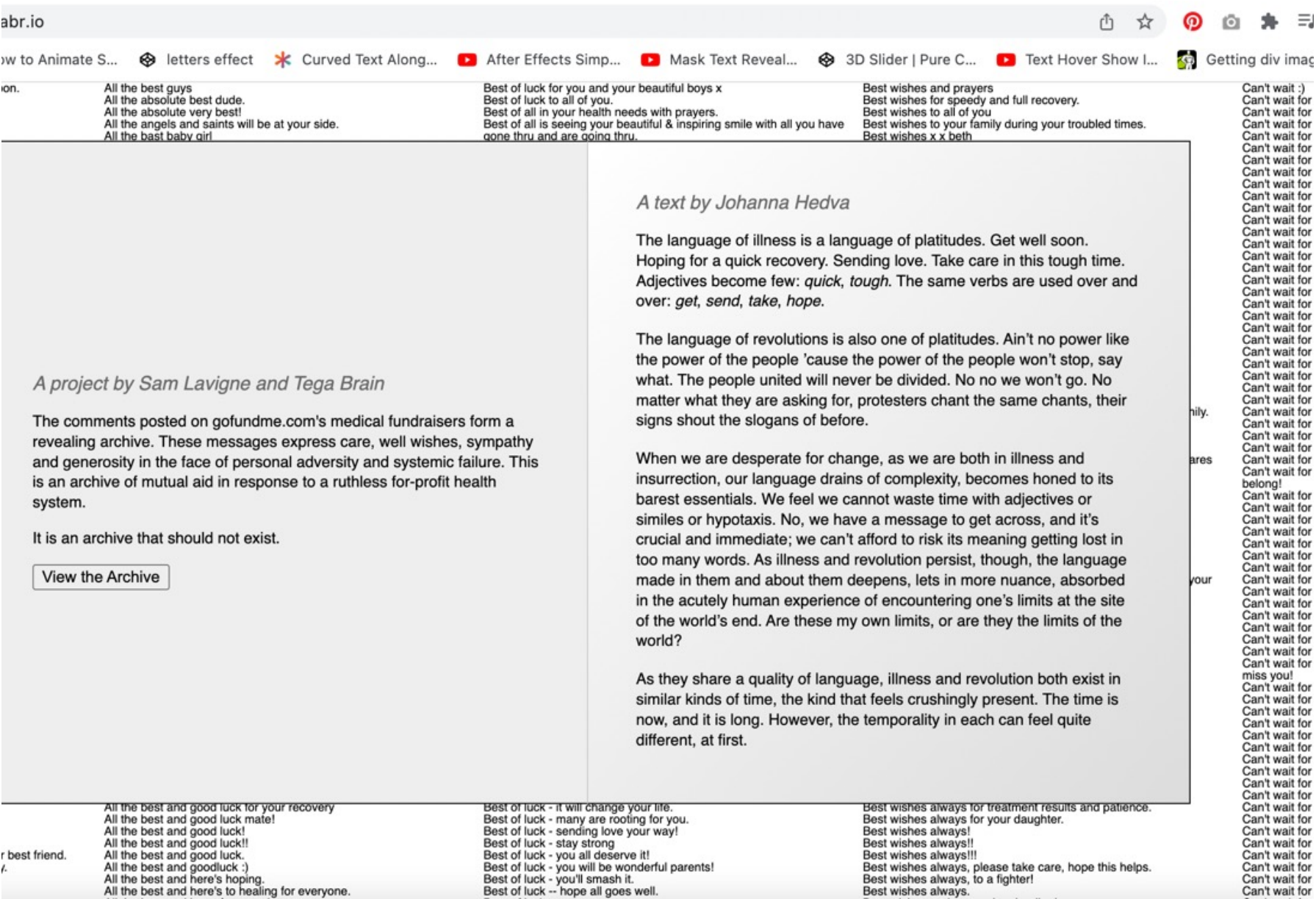
PLATFORMS

2022

# Abstract

The brief explains to explore an idea in depth, both technically and conceptually. It had to involve Big Data in some way, as well as include a web-based experiment.

The web based experiment chosen was a dynamic website, and the project itself is to visualise the handwritten messages on the back of lost postcards. Using OCR to extract the text and then convert it into text one can manipulate. How interesting it is to see a message written to another person out there; the emotions and thoughts revealed.

# Dilemmas

The process started with the idea of a dynamic website where a new poem/ message is presented on every refresh. A possible addition would be to add a voice to read it out loud. Some dilemmas were:

- To keep the handwriting or convert it? The presentation is more personal if the handwriting is kept.
- If the messages were to be read out loud, would it be a real person or automated; again, different impression.
- To include data as in how many letters sent and lost a year?
- Possible automated messages where one enters a phone number or email and a random message is sent.

# Development



There were multiple options as to how the project could be coded. Each with problems:

- Option 1 Write html in python
    - Problem adding css
- Option 2 Write python in html
    - Problem getting code to work
- Option 3 Write js in html
    - Problem linking code

The first attempts were done in Python. Different OCR's were tried out, but pytesseract proved to work the best. However, as can be seen by the bottom pictures, the accuracy wasnt the best. It was at that point that the handwritten postcards were temporarily changed to the dedications from books as to more easily interpret the converted text.

# Experimentation
## Python





```
</div>
      </div>
  </div>

<script type="text/python">
    from browser import document, alert

    def hello(ev):
        alert(" Open Alert")

    document["button-click"].bind("mouseenter", hello)

  </script>
    <button id="button-click"> Open Alert </button>
<!-- document['shownDedication'] <= "Hey" -->


<script>
```

```
website > python > 🐍 ocr.py > ...
  1
  2    import pytesseract
  3    import PIL.Image
  4    import random
  5    from gtts import gTTS
  6    from playsound import playsound
  7
  8    myconfig = r"--psm 6 --oem 3"
  9
 10    pytesseract.pytesseract.tesseract_cmd = r'/opt/homebrew/Cellar/tesseract/5.1.0/bin/tesseract'
 11
 12    first_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_1.jpeg"), config=myconfig)
 13    second_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_2.jpeg"), config=myconfig)
 14    third_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_3.jpeg"), config=myconfig)
 15    fourth_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_4.jpeg"), config=myconfig)
 16    fifth_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_5.jpeg"), config=myconfig)
 17    seventh_dedication = pytesseract.image_to_string(PIL.Image.open("dedications/dedication_7.jpeg"), config=myconfig)
 18
 19
 20    dedications = [
 21        first_dedication,
 22        second_dedication,
 23        third_dedication,
 24        fourth_dedication,
 25        fifth_dedication,
 26        seventh_dedication,
 27    ]
 28
 29    message = (random.choice(dedications))
 30    print(message)
 31
 32    language = 'en'
 33    audio_message = gTTS(text=message, lang=language, slow=False)
 34
 35    audio_message.save("dedicationSound.mp3")
 36    playsound("dedicationSound.mp3")
 37
 38
```



The next step was to add text to speech to play in the background, which was done with gtts and playsound. The python code was then pretty much complete. That's when there was an attempt to link the python code to the HTML file using pyscript, however errors either occured or the page would be stuck loading. It was later attempted to be done using brython as well, however that wasn't very successful either.

# Experimentation
## Visual HTML



**hello, there**

amaliemorch@pc-207-181
website % /usr/local/bin/python3
/Users/amaliemorch/Docum
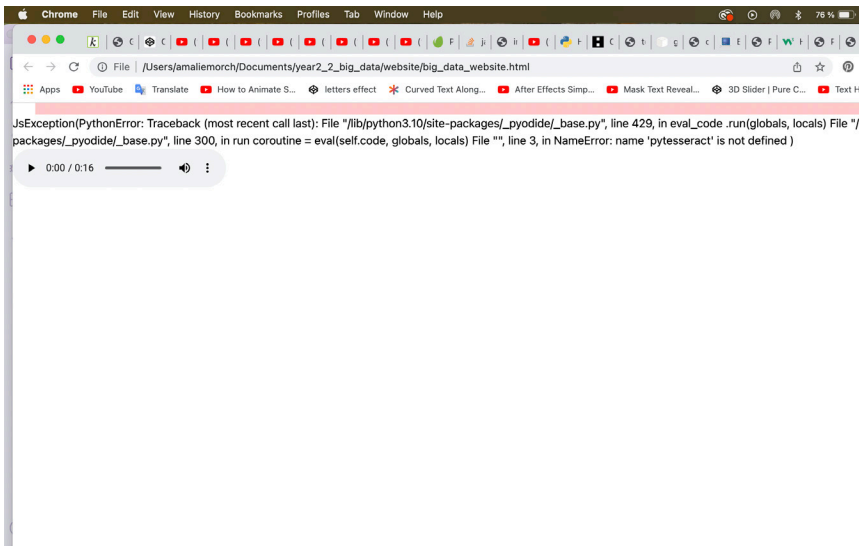ents/year2_2_big_data/website/ocr
This book is dedicated to
everybody you hate. Sorry. Life's
like that sometimes.
amaliemorch@pc-207-181
website % /usr/local/bin/python3 /
Users/amaliemorch/Documen
ts/year2_2_big_data/website/ocr
To all children who have ever felt
different amaliemorch@pc-207-
181 website %



*Dedicated to -*

**hello, there**

amaliemorch@pc-207-181
website % /usr/local/bin/python3
/Users/amaliemorch/Docum
ents/year2_2_big_data/website/ocr
This book is dedicated to
everybody you hate. Sorry. Life's
like that sometimes.
amaliemorch@pc-207-181
website % /usr/local/bin/python3 /
Users/amaliemorch/Documen
ts/year2_2_big_data/website/ocr
To all children who have ever felt
different

The original idea was to have an envelope opening, a letter slide out with the converted text. However that idea changed to a book opening instead. The book started out as closed with the text printed at the front like a title. However, the dedications in books are inside so it was then changed to be presented as open. Furtehr, it was found a little boring, so an animation was created on hover so the book opens and closes.

```
1
2    // export default App;
3
4    const tesseract = require("tesseract.js")
5
6    tesseract.recognize("dedication_4.jpeg", "eng", {logger:m => console.log(m)})
7  ∨  .then(result => {
8        console.log(result.data.text);
9     })
0  ∨  .catch((err) => {
1        console.log(err.message);
2     });
```

```
website > JS tesseract_trial.js > …
27
28   const getImageText = async(fileName, lang, logger) => {
29       let recognizeResult = null
30       try {
31           if(fs.existsSync(fileName)){ //check if file existss
32               if(logger){ //if user has provided logger
33                   const myLogger = {
34                       logger: m => console.log(m)
35                   }
36                   recognizeResult = await tesseract.recognize(fileName, lang, myLogger);
37               } else {
38                   recognizeResult = await tesseract.recognize(fileName, lang);
39               }
40               if(recognizeResult) {
41                   return recognizeResult.data.text
42               }
43           }
44       } catch (error) {
45           return error.message
46       }
47   };
48
49   getImageText("dedications/dedication_4.jpeg", "eng").then(content => {
50       console.log(content)
51   }).catch (errMsg => {
52       cconsole.log(errMsg);
53   });
```

```
// .then(({data: {text}}) => {
//   console.log(text);
// })
        document.addEventListener('DOMContentLoaded', function(){
            var input_image = document.getElementById('input_image');
            input_image.addEventListener('change', handleInputChange);
        });

        function handleInputChange(file){
            var input = event.target;
            var file = input.files[0];
            console.log(file);
            Tesseract.recognize(file)

                .then(function(result){
                    var contentArea = document.getElementById('dedication-text');
                    contentArea.innerHTML = result.text;
                    console.log(result);
                })
                // .catch(function(err){
                //    console.error(err);
                // });
        }
```
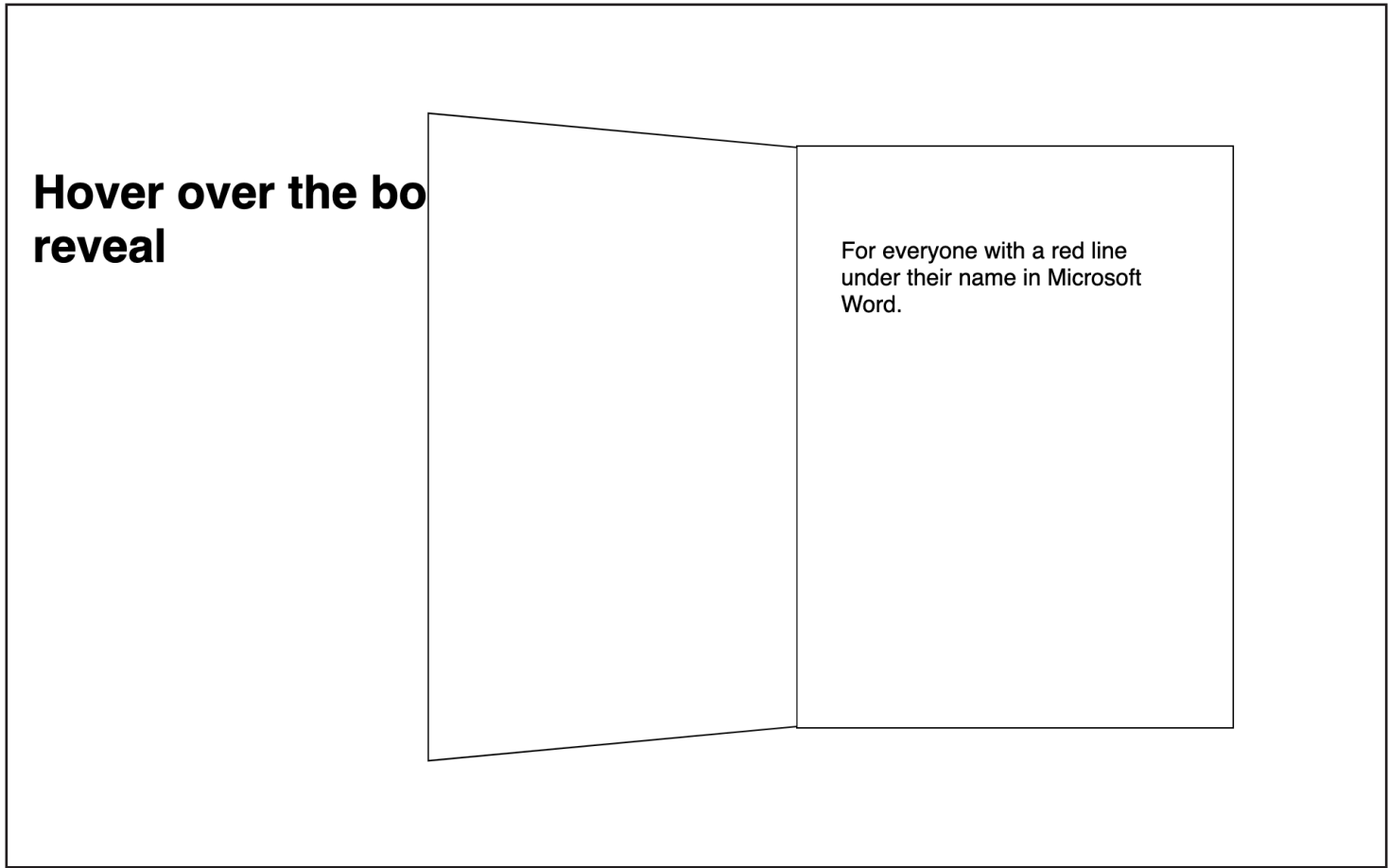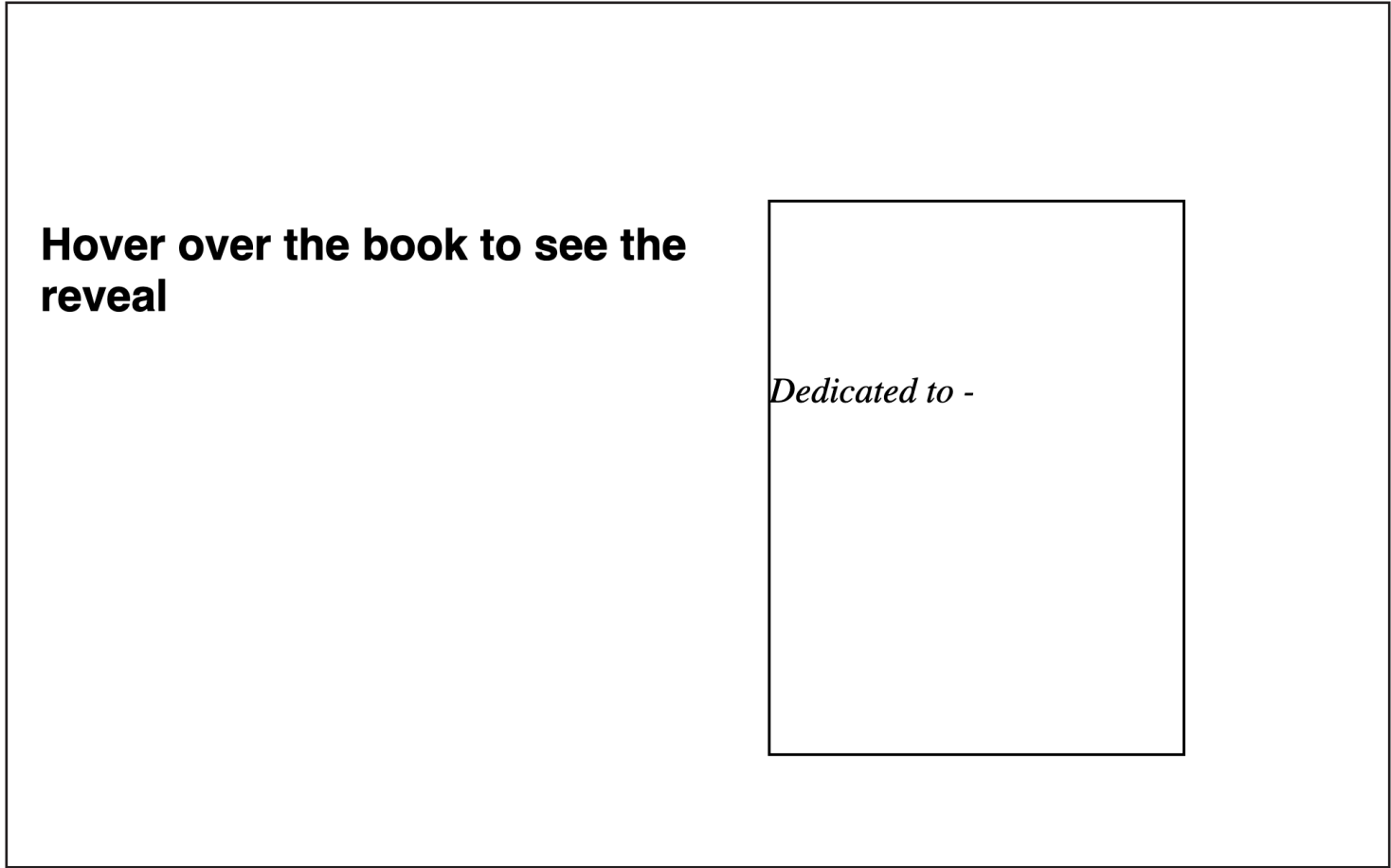
```
31   import { useState} from 'react';
32   import Tesseract from 'tesseract.js';
33   import './App.css';
34
35   function App() {
36       const [imagePath, setImagePath] = useState("");
37       const [text, setText] = useState("");
38
39       const handleChange = (event) => {
40           setImagePath(URL.createObjectURL(event.target.files[0]));
41       }
42
43       const handleClick = () => {
44
45           Tesseract.recognize(
46               imagePath,'eng',
47               {
48                   logger: m => console.log(m)
49               }
50           )
51           .catch (err => {
52               console.error(err);
53           })
54           .then(result => {
55               // Get Confidence score
56               let confidence = result.confidence
57
```

```
const T = require("tesseract.js");
// const [text, setText] = useState("");
T.recognize(chosenDedication, "eng")
.then(({data: {text}}) => {
    var output = text;
    document.getElementById("msgg").innerHTML = output;
    // document.write(output);
});
// .then(function (result) {
//     res.send(result.text);
// })

// .then(function (doneResult) {
//     var text = doneResult.text;
//     setText(text);
// })

//    .then(({data: {text}}) => {
//        console.log(text);
//    });
//    .then(({data: {text}}) => {
//    result.innerHTML = `<p>OCR Result:</p><p>${text}</p>`;
//    });

// .then(convertedText => {
// let text = convertedText.text;
// setText(text);
// })
```
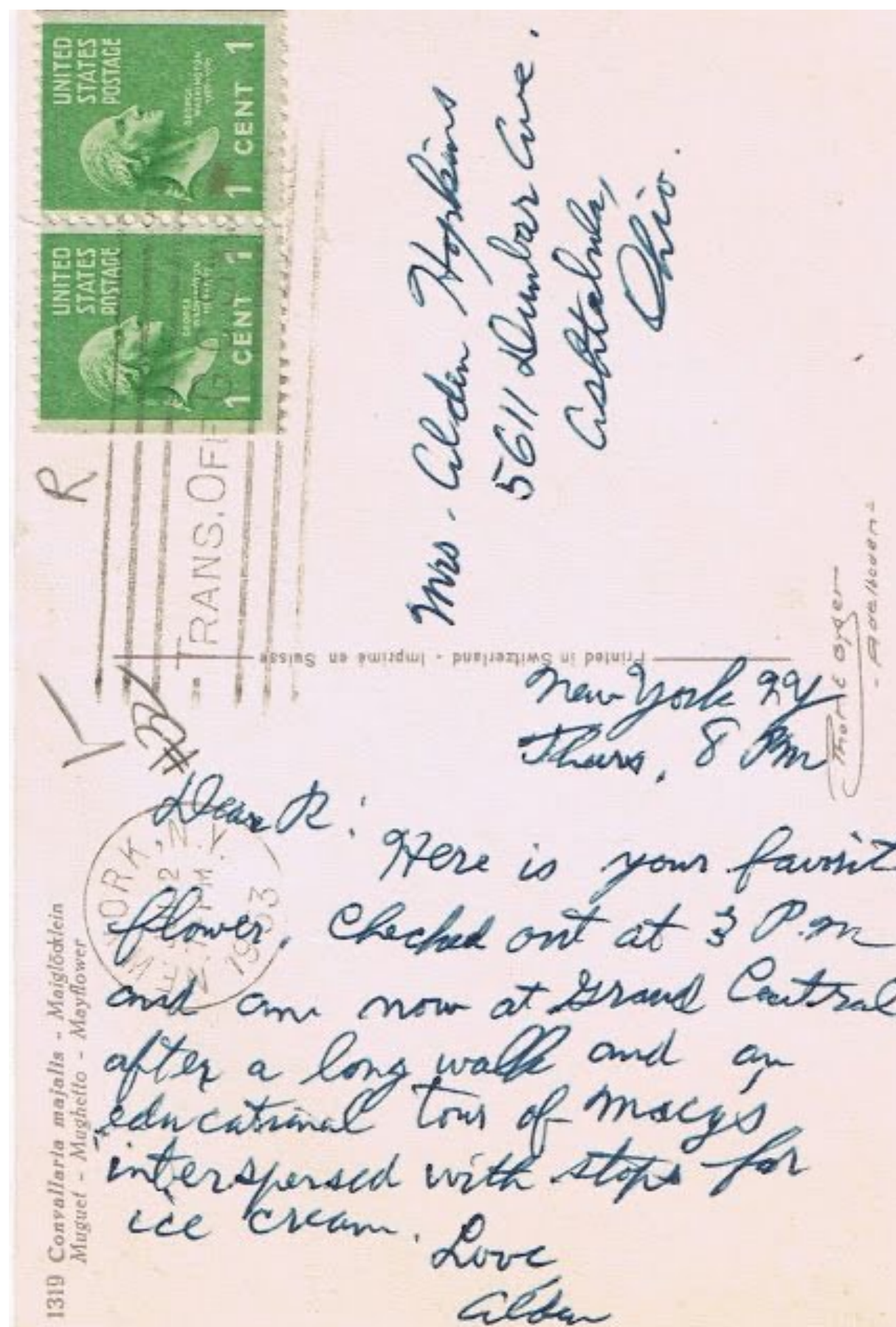
The next step was to try to implement the OCR in Javascript, instead of Python. This was a big challenge as there wasn't much experience with this. There were several attempt at implementing the OCR, as it simply wouldnt visualise in the HTML. It would however print fine in the console, so the code was indeed working.

Many possible solutions were implemented to try to solve the problem:

- To replace innerHTML with innerText.
- To move the id element before and after the script element (before the JS was moved to a seperate file).
- To change the linked HTML element.
- To replace require() with import, as well several other attempts to remove the error occuring in the browser as to items not being defined.
- Turned out the solution was a missing CDN link in the head.

# Finished project

**Hover over the book to see the reveal**

*Dedicated to -*

**Hover over the bo... reveal**

For everyone with a red line under their name in Microsoft Word.

The finished project is a dynamic website which shows a book; on hover, the book opens and presents one random dedication out of an array. There is also a text to speech feature playing in the background.

# Possible improvements
## Further development

- Accuracy: If further developed, there would be improvements made as to the accuracy of the OCR, as the one currently used includes errors.
- Handwriting: There would be an attempt to figure out the original idea of using postcards by making the conversion more accurate, as well as within specific area.
- Visualise a map (as to loctaions) or timeline (as to dates) postcards were posted; possible using D3.
- If continued with dedications, the option to choose a genre of book to see differences in dedications.

# Biggest challenge

Dedicated to -

The biggest challenge was simply visualising the OCR. It was managable to code it in both Python and Javascript, but to actually show the output in HTML was extremely difficult for some reason. Python and HTML wasn't very compatible, and for some reason Javascript simply didn't want to link up to the HTML file. However, with much frustration, and a few tears, it was managable.