# Appendix

To the manuscript 'OCR with Tesseract, Amazon Textract, and Google Document AI: A Benchmarking Experiment'
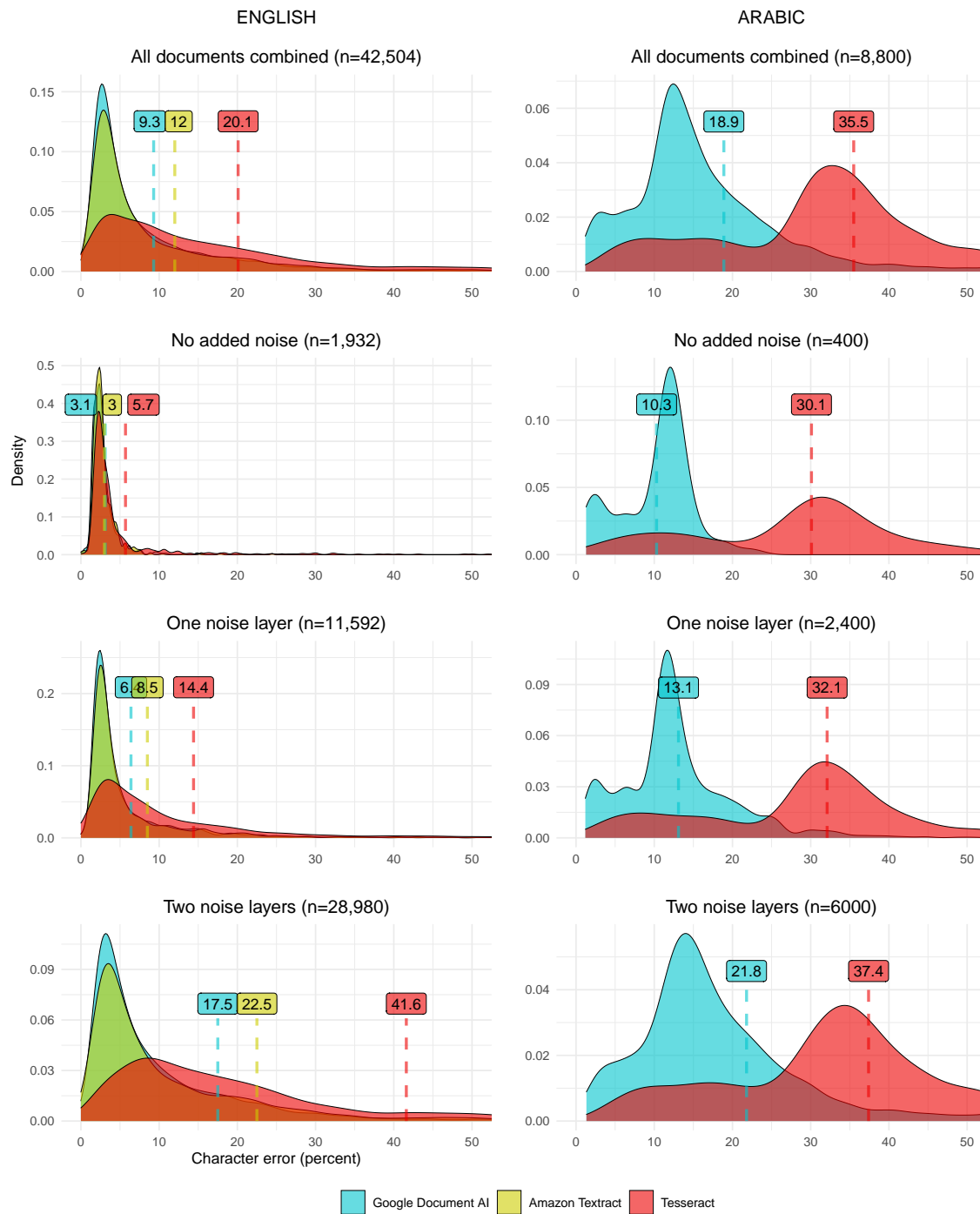
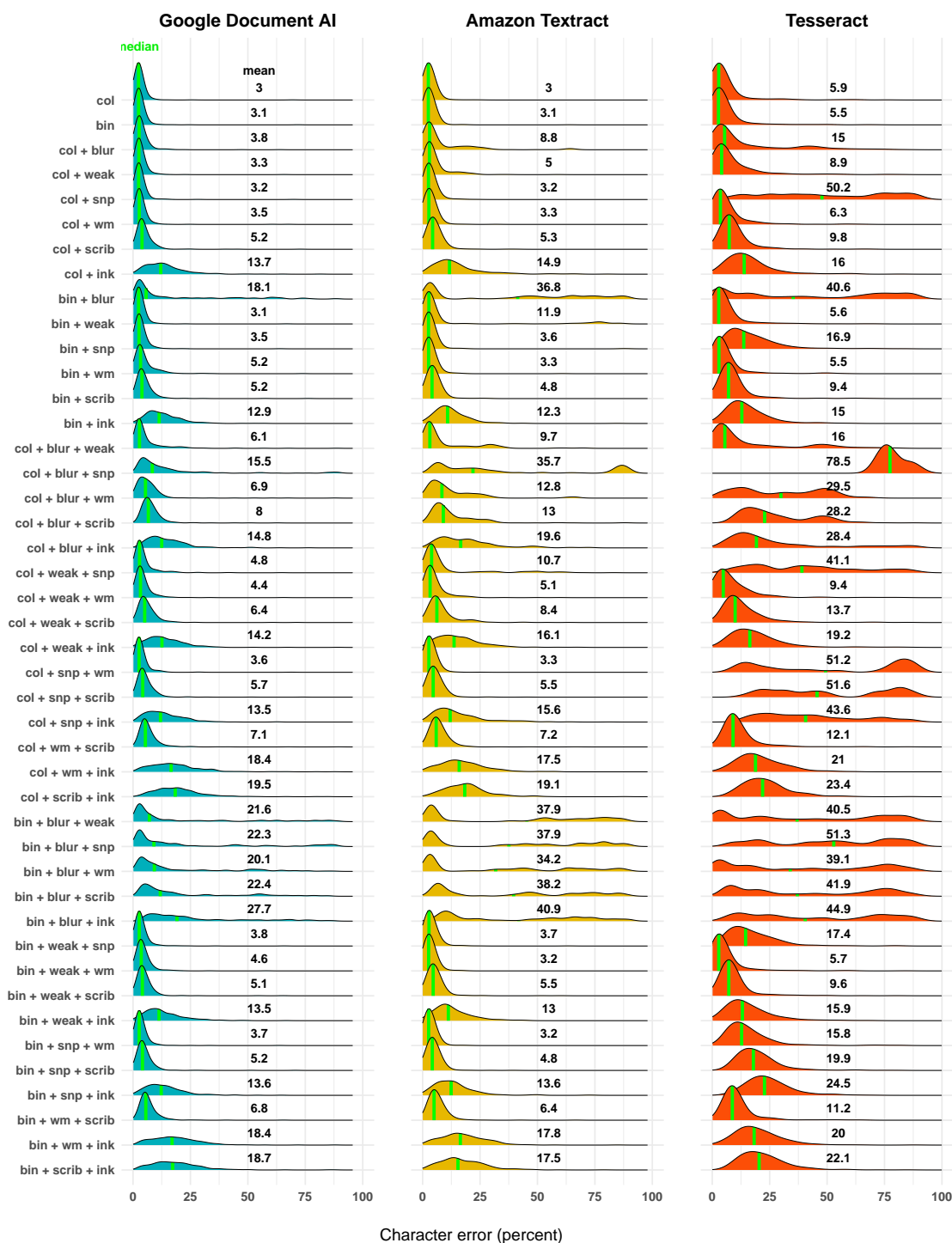## Contents

# 1 Character accuracy results

## 1.1 Character error rates by engine and noise level for English and Arabic documents

Mean error rates in coloured boxes. X axes cropped for visibility, leaving out the tails of the distributions.
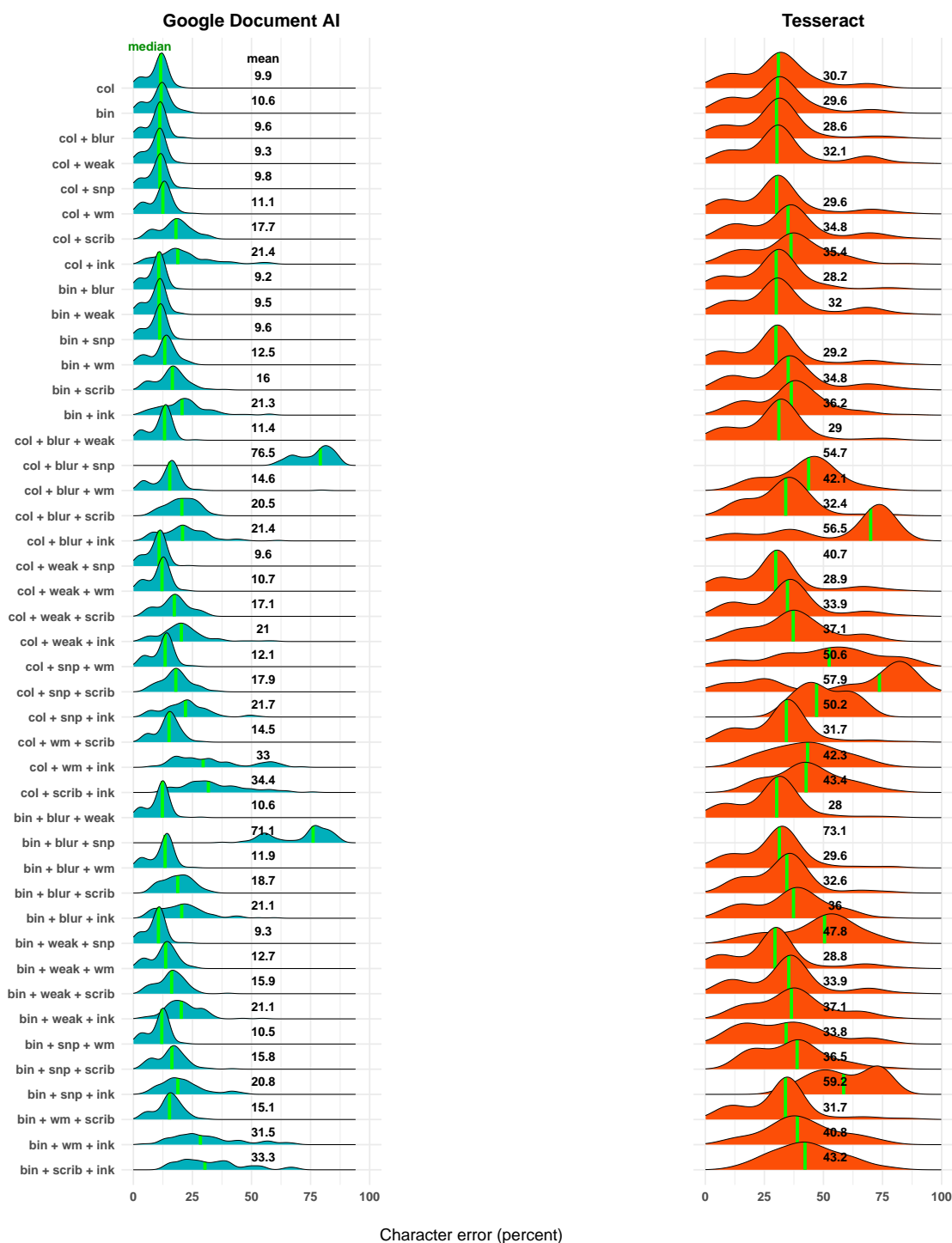
## 1.2 Character error rates by engine and noise type for English-language documents

Data: Single–column text in historical book scans with noise added articifially (n=42,504; 322 per engine and noise type).
Noise codes: 'col'=colour, 'bin'=binary, 'blur'=blur, 'weak'=weak ink, 'snp'=salt&pepper, 'wm'=watermark, 'scrib'=scribbles, 'ink'=ink stains.

| Noise type | Google Document AI (mean) | Amazon Textract (mean) | Tesseract (mean) |
|---|---|---|---|
| col | 3 | 3 | 5.9 |
| bin | 3.1 | 3.1 | 5.5 |
| col + blur | 3.8 | 8.8 | 15 |
| col + weak | 3.3 | 5 | 8.9 |
| col + snp | 3.2 | 3.2 | 50.2 |
| col + wm | 3.5 | 3.3 | 6.3 |
| col + scrib | 5.2 | 5.3 | 9.8 |
| col + ink | 13.7 | 14.9 | 16 |
| bin + blur | 18.1 | 36.8 | 40.6 |
| bin + weak | 3.1 | 11.9 | 5.6 |
| bin + snp | 3.5 | 3.6 | 16.9 |
| bin + wm | 5.2 | 3.3 | 5.5 |
| bin + scrib | 5.2 | 4.8 | 9.4 |
| bin + ink | 12.9 | 12.3 | 15 |
| col + blur + weak | 6.1 | 9.7 | 16 |
| col + blur + snp | 15.5 | 35.7 | 78.5 |
| col + blur + wm | 6.9 | 12.8 | 29.5 |
| col + blur + scrib | 8 | 13 | 28.2 |
| col + blur + ink | 14.8 | 19.6 | 28.4 |
| col + weak + snp | 4.8 | 10.7 | 41.1 |
| col + weak + wm | 4.4 | 5.1 | 9.4 |
| col + weak + scrib | 6.4 | 8.4 | 13.7 |
| col + weak + ink | 14.2 | 16.1 | 19.2 |
| col + snp + wm | 3.6 | 3.3 | 51.2 |
| col + snp + scrib | 5.7 | 5.5 | 51.6 |
| col + snp + ink | 13.5 | 15.6 | 43.6 |
| col + wm + scrib | 7.1 | 7.2 | 12.1 |
| col + wm + ink | 18.4 | 17.5 | 21 |
| col + scrib + ink | 19.5 | 19.1 | 23.4 |
| bin + blur + weak | 21.6 | 37.9 | 40.5 |
| bin + blur + snp | 22.3 | 37.9 | 51.3 |
| bin + blur + wm | 20.1 | 34.2 | 39.1 |
| bin + blur + scrib | 22.4 | 38.2 | 41.9 |
| bin + blur + ink | 27.7 | 40.9 | 44.9 |
| bin + weak + snp | 3.8 | 3.7 | 17.4 |
| bin + weak + wm | 4.6 | 3.2 | 5.7 |
| bin + weak + scrib | 5.1 | 5.5 | 9.6 |
| bin + weak + ink | 13.5 | 13 | 15.9 |
| bin + snp + wm | 3.7 | 3.2 | 15.8 |
| bin + snp + scrib | 5.2 | 4.8 | 19.9 |
| bin + snp + ink | 13.6 | 13.6 | 24.5 |
| bin + wm + scrib | 6.8 | 6.4 | 11.2 |
| bin + wm + ink | 18.4 | 17.8 | 20 |
| bin + scrib + ink | 18.7 | 17.5 | 22.1 |

Character error (percent)

## 1.3 Character error rates by engine and noise type for Arabic-language documents

Data: Single–column text in image scans of Arabic Wikipedia pages with noise added articifially (n = 8800; 100 per engine and noise type).
Noise codes: 'col'=colour, 'bin'=binary, 'blur'=blur, 'weak'=weak ink, 'snp'=salt&pepper, 'wm'=watermark, 'scrib'=scribbles, 'ink'=ink stains.

| Noise type | Google Document AI (mean) | Tesseract (mean) |
|---|---|---|
| col | 9.9 | 30.7 |
| bin | 10.6 | 29.6 |
| col + blur | 9.6 | 28.6 |
| col + weak | 9.3 | 32.1 |
| col + snp | 9.8 | 29.6 |
| col + wm | 11.1 | 34.8 |
| col + scrib | 17.7 | 35.4 |
| col + ink | 21.4 | 28.2 |
| bin + blur | 9.2 | 32 |
| bin + weak | 9.5 | 29.2 |
| bin + snp | 9.6 | 34.8 |
| bin + wm | 12.5 | 36.2 |
| bin + scrib | 16 | 29 |
| bin + ink | 21.3 | 54.7 |
| col + blur + weak | 11.4 | 42.1 |
| col + blur + snp | 76.5 | 32.4 |
| col + blur + wm | 14.6 | 56.5 |
| col + blur + scrib | 20.5 | 40.7 |
| col + blur + ink | 21.4 | 28.9 |
| col + weak + snp | 9.6 | 33.9 |
| col + weak + wm | 10.7 | 37.1 |
| col + weak + scrib | 17.1 | 50.6 |
| col + weak + ink | 21 | 57.9 |
| col + snp + wm | 12.1 | 50.2 |
| col + snp + scrib | 17.9 | 31.7 |
| col + snp + ink | 21.7 | 42.3 |
| col + wm + scrib | 14.5 | 43.4 |
| col + wm + ink | 33 | 28 |
| col + scrib + ink | 34.4 | 73.1 |
| bin + blur + weak | 10.6 | 29.6 |
| bin + blur + snp | 71.1 | 32.6 |
| bin + blur + wm | 11.9 | 36 |
| bin + blur + scrib | 18.7 | 47.8 |
| bin + blur + ink | 21.1 | 28.8 |
| bin + weak + snp | 9.3 | 33.9 |
| bin + weak + wm | 12.7 | 37.1 |
| bin + weak + scrib | 15.9 | 33.8 |
| bin + weak + ink | 21.1 | 36.5 |
| bin + snp + wm | 10.5 | 59.2 |
| bin + snp + scrib | 15.8 | 31.7 |
| bin + snp + ink | 20.8 | 40.8 |
| bin + wm + scrib | 15.1 | 43.2 |
| bin + wm + ink | 31.5 | |
| bin + scrib + ink | 33.3 | |

Character error (percent)

## 2    R code for noise generation

```
### Prerequisites ###########
#
# a) Install the following R packages: dplyr, glue, magick, yarr, showtext, shape
#
# b) For functions 4-6, set the image dimensions (in pixels) as follows:
# img_width <- #<An integer>
# img_height <- #<An integer>
#
# c) For function 5, install the proprietary font "Shopping Script" (https://www.dafont.com/shopping-sc
#
###########################


### Noise-generating functions

# All functions take two inputs:
# 1) source_ims: a vector of filepaths for the source images
# 2) dest_folder: a folder path for destination directory

# 1. BLUR
blur <- function(source_ims, dest_folder) {
  for (i in source_ims){
    magick::image_read(i) %>%
    magick::image_blur(5, 4) %>%
    magick::image_write(glue::glue("{dest_folder}/{basename(i)}"))
  }
}

# 2. WEAK INK
weaken <- function(source_ims, dest_folder) {
  for (i in source_ims){
    magick::image_read(i) %>%
    magick::image_oilpaint() %>%
    magick::image_write(glue::glue("{dest_folder}/{basename(i)}"))
  }
}

# 3. SALT & PEPPER
snp <- function(source_ims, dest_folder) {
  for (i in source_ims){
    magick::image_read(i) %>%
    magick::image_noise(noisetype = "poisson") %>%
    magick::image_write(glue::glue("{dest_folder}/{basename(i)}"))
  }
}

# 4. WATERMARK
watermark <- function(source_ims, dest_folder) {
  transp_grey <- yarrr::transparent(orig.col = "gray80",
                                    trans.val = 0.4,
                                    maxColorValue = 255)
```

```r
  for (i in source_ims){
    img <- magick::image_read(i)
    tiff(glue::glue("{dest_folder}/{basename(i)}"),
         width=img_width,
         height=img_height,
         units="px",
         res=300)
    plot(img)
    text(1200, 2000, # NB adapt coordinates to image dimensions
         "Watermark",
         cex=20,
         srt=50,
         col=transp_grey)
    dev.off()
  }
}


# 5. SCRIBBLES
scribble <- function(source_ims, dest_folder) {
  showtext::font_add(family = "Shopping Script",
                     regular = "~/.fonts/ShoppingScript-Regular.otf")
  showtext::showtext_auto()
  for (i in source_ims){
    img <- magick::image_read(i)
    tiff(glue::glue("{dest_folder}/{basename(i)}"),
         width=img_width,
         height=img_height,
         units="px",
         res=300)
    plot(img)
    # NB adapt coordinates below to image dimensions
    text(1500, 2000, "fascinating", family = "Shopping Script", cex = 10, srt = 15, col = "gray30")
    text(600, 1000, "__", family = "Shopping Script", cex = 12, col = "gray10", srt = 4)
    text(300, 700, "_ _", family = "Shopping Script", cex = 12, col = "gray20", srt = -4)
    text(1600, 800, "NB!", family = "Shopping Script", cex = 10, srt = -10, col = "gray30")
    text(1650, 1400, "V", family = "Shopping Script", cex =8, srt = 10, col = "gray35")
    text(1700, 500, "V", family = "Shopping Script", cex =12, srt = 30, col = "gray25")
    text(400, 1500, "mmm", family = "Shopping Script", cex=10)
    text(1100, 500, "mnnm", family = "Shopping Script", cex=10, col = "gray35")
    text(1000, 1200, "mmmmmmm", family = "Shopping Script", cex=10, col = "gray25")
    text(50, 1300, "|", family = "Shopping Script", cex=10)
    text(100, 1100, "Z", family = "Shopping Script", cex=10, col = "gray25")
    text(400, 1800, "O", family = "Shopping Script", cex=12, col = "gray35")
    text(1200, 700, "O", family = "Shopping Script", cex=10, col = "gray25")
    text(450, 400, "X", family = "Shopping Script", cex=10, col = "gray15")
    dev.off()
  }
}


# 6. INK STAINS
ink <- function(source_ims, dest_folder) {
  for (i in source_ims){
    tiff(file="ink.tiff",
```

```
            width=img_width,
            height=img_height,
            units="px",
            res=300)
    par(bg=NA)
    shape::emptyplot()
    shape::filledshape(matrix(nc = 4,
                              nr = 4,
                              runif(8)),
                       col = shadepalette(50, "black", "grey20"))
    dev.off()
    main <- magick::image_read(i)
    inset <- magick::image_read("ink.tiff")
    final <- magick::image_composite(main, inset, operator = "atop")
    image_write(final, glue::glue("{dest_folder}/{basename(i)}"))
  }
}
```