

# Comparison of OCR and CNN Classification Results

Ryan Heslin

April 27, 2022

This document compares the results of classifying all 707 images currently available, using both a trained neural network and optical character recognition with `tesseract`.

In the output below, blue text designates hyperlinks that should point to the relevant image in my Dropbox storage.

```
source(here::here("R", "utilities.R"))
source(here::here("R", "classifier_utilities.R"))
source(here::here("R", "dropbox_utilities.R"))
ocr_results <- readRDS(here::here("data", "classifier", "outputs", "ocr_results.Rds"))
cnn_results <- read.csv(here::here("data", "classifier", "outputs", "cnn_all_classifications.csv"))
```

I opted to record results separately, since I planned to record CNN results many times. The datasets can be joined using the `image` column as a key. The columns of `cnn_results` are:

- `actual_class`, `pred_class`: True and predicted classes, with 1 corresponding to non-Trump and 2 to Trump
- `image`: Path to the image corresponding to the row
- `p_trump`: Predicted Trump probability
- `log_p_trump`: Log of predicted Trump probability
- `batch_loss`: Value of loss function for batch the observation belonged to

The columns of `ocr_results` are:

- `image` and `ocr_pred_class`: The same meanings as in `cnn_results`
- `ocr_text`: The raw text recognized by `tesseract`. Prediction was done by case-insensitive matching for the string “trump” in this text.

- `ocr_output`: List column of detailed `tesseract` results. Each entry contains the columns `word` (each distinct recognized word, concatenated to form `ocr_text`), `confidence` (percentage confident in word identification), and `bbox` (bounding box of image coordinates corresponding to identified word). Because this is a list column, I chose to store the data as an RDS.

```
combined <- merge(ocr_results, cnn_results, by = "image")
```

All images are in both result sets.

```
setequal(ocr_results$image, cnn_results$image)
```

```
[1] TRUE
```

## Comparing Misclassifications

A contingency table shows the different classifications, with the neural network on the rows and OCR on the columns. Most agree, with the exception of the network predicting “Trump” far too often.

```
present_table(
  with(
    combined,
    table(
      cnn = CLASS_NAMES[pred_class],
      ocr = CLASS_NAMES[ocr_pred_class]
    )
  )
))
```

	no_trump	trump
no_trump	632	1
trump	72	2

Here is a more detailed breakdown of predicted and actual labels. The network appeared to be more likely to predict an image contained Trump than OCR. I added training and loss weights to compensate for the shortage of Trump images, but I think they may need to be adjusted downward.

```
with(combined, table(
  cnn = CLASS_NAMES[pred_class],
  ocr = CLASS_NAMES[ocr_pred_class],
  true = CLASS_NAMES[actual_class]
)) %>%
  present_table(., FALSE)
```

cnn	ocr	true	Freq
no_trump	no_trump	no_trump	631
trump	no_trump	no_trump	36
no_trump	trump	no_trump	0
trump	trump	no_trump	0
no_trump	no_trump	trump	1
trump	no_trump	trump	36
no_trump	trump	trump	1
trump	trump	trump	2

In most cases, both classifiers were correct. Interestingly, the number of cases where only one was correct is about the same for both.

```
combined <- within(combined, prediction_category <- ifelse(actual_class == ocr_pred_class,
  ifelse(pred_class == actual_class, "both", "ocr"),
  ifelse(pred_class == actual_class, "cnn", "neither"))
))

sort(table(combined[["prediction_category"]]), decreasing = TRUE) |>
  present_table()
```

both	ocr	cnn	neither
633	37	36	1

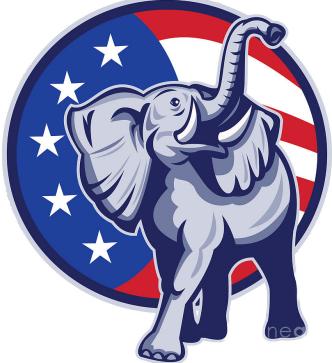
```
disagreements <- combined[["prediction_category"]] != "both"
```

The proportion of cases where one or both classifiers was wrong 0.105. The overall error rate from naively predicting “not Trump” in every case would be 0.057.

We can examine the images where at least one classification was incorrect. This plot shows the first four of those 74 images. I don’t see any obvious pattern. `tesseract` seems to have correctly classified some images with no text, which I have to regard as spurious.

```
to_plot <- seq_len(4)
inspect_images_from_paths(
  paths = combined[disagreements, "image"][to_plot],
  labels = paste("Correctly classified by:", combined[disagreements, "prediction_category"][to_plot]),
  plot_dims = c(2, 2)
)
```

area\_25\_republican\_committee\_bgimg\_20220211.png



Correctly classified by: ocr  
austin\_chenge\_for\_governor\_bgimg\_20220211.jpg

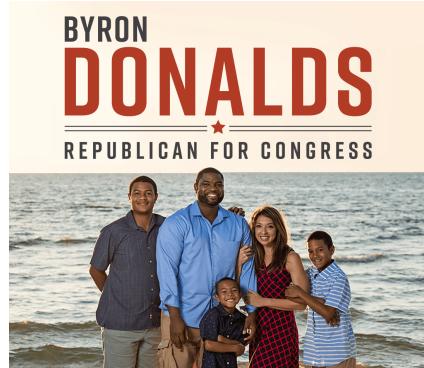


Correctly classified by: ocr

brian\_for\_better\_nj\_inc\_bgimg\_20220211.png



Correctly classified by: ocr  
byron\_donalds\_bgimg\_20220211.png



Correctly classified by: ocr

```
urls <- read.csv(here::here("data", "classifier", "image_urls.csv"))
combined <- merge(combined, urls, by = "image")
```

Here are tables showing the true labels and predictions from both classifiers for the disagreements.

```
# Complicated code to wrangle a table
subset <- combined[disagreements, c(
  "image", "ocr_pred_class",
  "pred_class", "actual_class", "prediction_category"
)]
rownames(subset) <- NULL
subset[["image"]] <-
  tools::file_path_sans_ext(basename(subset[["image"]]))
subset[, endsWith(colnames(subset), "class")] <-
  lapply(
    subset[, endsWith(colnames(subset), "class")],
    function(x) CLASS_NAMES[x]
  )
subset[["prediction_category"]] <- latex_url(
  display = subset[["prediction_category"]],
  url = combined[disagreements, "url"]
)
subset[["image"]] <- gsub("_", "\n", subset$image) %>%
  lapply(strwrap, 20)
```

```

n_lines <- lengths(subset[["image"]])
subset[["image"]] <- lapply(
  subset[["image"]],
  function(x) {
    paste(
      "{\\shortstack{",
      gsub("(?<!\\\\\\\\)\\s(?!\\\\\\\\)", " \\_",
        paste(x, collapse = " \\\\")),
        perl = TRUE
      ),
      "}}"
    )
  }
)

col_spec <- paste0("|", paste(c("l", rep("c", ncol(subset) - 1)),
  collapse = "|"),
), "|")
table_size <- 15
for (spl in split(subset, seq_len(nrow(subset)) %% table_size)) {
  rows <- gsub("_", "\\\\", paste(do.call(paste, append(spl, list(sep = " & "))),
    collapse = " \\\\"\\n\\hline \\n"
  ))
}

display_text <- paste(
  paste0("\\noindent \\begin{tabular}{", col_spec, "}\n"),
  "\\toprule\n",
  "Image & OCR Prediction & CNN Prediction & True Class & Which Correct \\\\"\\n",
  "\\hline\\n",
  rows,
  "\\\\\\\\\\n",
  "\\bottomrule\\n",
  "\\end{tabular}",
  collapse = "\\n"
)
cat(display_text, "\\n\\r\\n")
}

```

Image	OCR Prediction	CNN Prediction	True Class	Which Correct
area_25_republican_committee_bgimg_20220211	no_trump	trump	no_trump	ocr
austin_chenge_for_governor_bgimg_20220211	no_trump	trump	no_trump	ocr
brian_for_better_nj_inc_bgimg_20220211	no_trump	trump	no_trump	ocr
byron_donalds_bgimg_20220211	no_trump	trump	no_trump	ocr
carlos_gimenez_bgimg_20220211	no_trump	trump	no_trump	ocr
citizens_for_amy_summers_bgimg_20220211	no_trump	trump	no_trump	ocr
commette_to_elect_barbara_v_johnson_for_new_mexico_court_of_appeals_bgimg_20220211	no_trump	trump	no_trump	ocr
committee_to_elect_kerry_j_morris_to_new_mexico_supreme_court_bgimg_20220211	no_trump	trump	no_trump	ocr
committee_to_elect_stephen_cicak_bgimg_20220211	no_trump	trump	no_trump	ocr
committee_to_elect_thomas_arlinger_bgimg_20220211	no_trump	trump	no_trump	ocr
douglas_county_republican_party_bgimg_20220211	no_trump	trump	no_trump	ocr
greg_lirette_bgimg_20220211	no_trump	trump	no_trump	ocr
james_risch_bgimg_20220211	no_trump	trump	no_trump	ocr
jasinski_for_senate_bgimg_20220211	no_trump	trump	no_trump	ocr

Image	OCR Prediction	CNN Prediction	True Class	Which Correct
jerry_carl_bgimg_20220211	no_trump	trump	no_trump	ocr
john_barrasso_bgimg_20220211	no_trump	trump	no_trump	ocr
jordan_m_mackey_for_ks_house_33rd_bgimg_20220211	no_trump	trump	no_trump	ocr
joshua_morris_for_state_representative_bgimg_20220211	no_trump	trump	no_trump	ocr
julian_acciard_bgimg_20220211	no_trump	trump	no_trump	ocr
karl_allred_for_house_19_bgimg_20220211	no_trump	trump	no_trump	ocr
kevin_cramer_bgimg_20220211	no_trump	trump	no_trump	ocr
lindsey_graham_bgimg_20220211	no_trump	trump	no_trump	ocr
luke_hellier_for_city_council_bgimg_20220211	no_trump	trump	no_trump	ocr
matthew_burkhart_bgimg_20220211	no_trump	trump	no_trump	ocr
mike_bost_bgimg_20220211	no_trump	trump	no_trump	ocr
milliman_for_colorado_bgimg_20220211	no_trump	trump	no_trump	ocr
parson_for_missouri_bgimg_20220211	no_trump	trump	no_trump	ocr
re_elect_william_fowler_bgimg_20220211	no_trump	trump	no_trump	ocr
senator_tom_james_bgimg_20220211	no_trump	trump	no_trump	ocr

Image	OCR Prediction	CNN Prediction	True Class	Which Correct
team_doyle_bgimg 20220211	no_trump	trump	no_trump	ocr
the_boldyga committee_bgimg 20220211	no_trump	trump	no_trump	ocr
van_taylor_bgimg 20220211	no_trump	trump	no_trump	ocr
velez_for_wisconsin bgimg_20220211	no_trump	trump	no_trump	ocr
vicky_I_cook_for fork_supervisor bgimg_20220211	no_trump	trump	no_trump	ocr
vinny_panico_for assembly_bgimg 20220211	no_trump	trump	no_trump	ocr
west_virginia republican_party bgimg_20220211	no_trump	trump	no_trump	ocr
180515-h-ni589-449 extra_20220219	no_trump	trump	trump	cnn
190829-d-bn624-0276 extra_20220219	no_trump	trump	trump	cnn
african_american history_month reception_extra 20220219	no_trump	trump	trump	cnn
barry_moore_bgimg 20220211	no_trump	trump	trump	cnn
billy_long_bgimg 20220211	no_trump	trump	trump	cnn
cbp_tours_san_diego border_wall prototypes_with potus_extra 20220219	no_trump	trump	trump	cnn
donald_trump_and henry_extra 20220219	no_trump	trump	trump	cnn
donald_trump_and mike_extra_20220219	no_trump	trump	trump	cnn

Image	OCR Prediction	CNN Prediction	True Class	Which Correct
donald_trump_and_peter_extra_20220219	no_trump	trump	trump	cnn
donald_trump_rally_in_youngstown_july_extra_20220219	no_trump	trump	trump	cnn
doug_burgum_for_north_dakota_bgimg_20220211	no_trump	trump	trump	cnn
elise_stefanik_bgimg_20220211	no_trump	trump	trump	cnn
g7biarritz_extra_20220331	no_trump	trump	trump	cnn
gabriel_whitley_bgimg_20220211	no_trump	trump	trump	cnn
harmeet_dhillon-speaks_at_the_white_house-27s_social_media_summit_extra_20220219	no_trump	trump	trump	cnn
lagop_gotvr_dec2016_176_extra_20220219	no_trump	trump	trump	cnn
mo_brooks_bgimg_20220211	no_trump	trump	trump	cnn
president_donald_j_trump,_senator_tom_cotton,_and_senator_david_perdue,_august_2,_2017_extra_20220219	no_trump	trump	trump	cnn
president_donald_j._trump_delivers_remarks_on_5g_deployment_technology_in_the_united_states_extra_20220219	no_trump	trump	trump	cnn
president_trump_and_the_indy_500_winner_extra_20220331	no_trump	trump	trump	cnn
president_trump_announces_1.8_billion_in_opioids_grants_extra_20220219	no_trump	trump	trump	cnn
president_trump_at_dday75_extra_20220331	no_trump	trump	trump	cnn
president_trump_at_the_white_house_summit_on_human_trafficking_extra_20220219	no_trump	trump	trump	cnn

Image	OCR Prediction	CNN Prediction	True Class	Which Correct
president_trump_in_iowa_extra_20220219	no_trump	trump	trump	cnn
president_trump_signs_an_executive_order_extra_20220331	no_trump	trump	trump	cnn
president_trump_welcomes_the_prime_minister_of_the_slovak_republic_to_the_white_house_extra_20220219	no_trump	trump	trump	cnn
roast_and_ride_extra_20220219	no_trump	trump	trump	cnn
roger_penske_receives_the_medal_of_freedom_extra_20220331	no_trump	trump	trump	cnn
save_america_bgimg_20220211	no_trump	trump	trump	cnn
steven_palazzo_bgimg_20220211	trump	no_trump	trump	ocr
the_2019_prison_reform_summit_and_first_step_act_celebration_extra_20220219	no_trump	no_trump	trump	neither
the_congressional_picnic_extra_20220331	no_trump	trump	trump	cnn
the_hispanic_heritage_month_reception_extra_20220219	no_trump	trump	trump	cnn
the_north_dakota_state_bison_visit_the_white_house_extra_20220331	no_trump	trump	trump	cnn
trump_nashville_extra_20220219	no_trump	trump	trump	cnn
trump_welcoming_kenya-e2-80-99s_president_uhuru_extra_20220219	no_trump	trump	trump	cnn
UNGA_extra_20220331	no_trump	trump	trump	cnn
white_house_press_briefing_extra_20220219	no_trump	trump	trump	cnn

## Raw OCR Results

The OCR characters for the classification disagreements largely look like gibberish, though some valid words are recognizable. Here, I randomly sample 10 sets of words and display them beside the name of the associated image. I suspect `tesseract` is struggling with images that contain text.

```
texts <- sapply(combined[disagreements, "ocr_output"], '[[', "word") %>%
  sapply(paste, collapse = " ") %>%
  sapply(strwrap, width = 20) %>%
  lapply(paste, collapse = "\n") %>%
  mapply(FUN = c, markdown_url(
    url = combined[disagreements, "url"],
    display = basename(combined[disagreements, "image"]))
  ), ..) %>%
  unname()

size <- 10
set.seed(1)
indices <- sample(ncol(texts), size)
do.call(what = cat, append(gsub("\\\\(?!\ href)", "\\\\\\ ", paste(
  rep(c("Image:", "Text:"), times = size),
  texts[, indices]
), perl = TRUE), list(sep = "\n\n")))
```

Image: [the\\_congressional\\_picnic\\_extra\\_20220331.jpg](#)

Text: y ell TS Fe ems ali ee Ra 1S ae >» Be . . 7 pe - “a ‘hs - a 5 Cs a . << . oe Ps he > aac F iss xe a = ge Os ae ee . . iff ‘) b | “a - eo “ia? es a 7 a J , > ‘aeare + yj \ iy mnt \ 4 S) —— a 2 z AG! o y” # a’ nm, f \ Vd 1 | 1) ‘a } 4 Bar a a .— ta P y of a q 4 Z) i) y, # i =} ee. wl . ‘ - ¥ 2 & a = [ Z i) nd Si Py: mas < me ‘ . , > . “ ‘ Gea hae x Z 3 r Rit SA ¥ = va <““<— i 5” “ Jt fn eS | Vy , i y//) ; | gape me Me 4 se id, Se an fesse aie os \* i - bs Pate AD ‘ y , NX ; M6 | | = ’ Ret rer” ie, NEA L Ss a rr es (i Fc Saat ky mm \* 4 | 4 RR errr ilins CL one 2. Ls er A Gener onus N x Vp é: . ) RR a é ‘ x dass ht . Se sal Mn i =~ mt ee ae a p’s aa . HORE sr ay, ) x ‘ ‘@ N —— Giser rr casera |). Y i N ANY Serene y y i > “ , »\ pee, RRR > = Gocoacea i Ig \ nt a. / \ RRS eae yO 4 “ ees Rees Peete vy Vy” el y — IRR tare V4 Y |. ave eet i

Image: [african\\_american\\_history\\_month\\_reception\\_extra\\_20220219.jpg](#)

Text: “ OE ae ee a” “ wry “i \_ Sy Ba oe . & \ eae “e . a 4 ae , 7 é 7 “7 \_ \* ‘ ‘ / “a Ce, tistCt” au (i Pies Wh ? . Be ee ae NS % , ee / i | i, ee aa oe . Vis .. i } sii? A. sgn sane si — . es = / < x , i de a 3 ) Pye RA I> @n - | = < a Ne “od o £6 meet, — SS pe ; d aN f — 3 ie ae ; \ y ia, . ‘ef i : 4 0 — a e id es th ~ hey 4 ” , & fis VAR O eeeeeerererreemermenmenmrrermenss 4 so By Te a - J s a : . . : he ao ‘ i 4 y : re i’ , ; A seciead iV ‘ / eer AOE P15 Se ( > ar Ache y 0 eee eo / ~ ee , / ea —— bed y Pee ¥ : . y 4 . : 4 = / a ee me ” a a ‘ TM o~ jt st gg = ws ors — 2 gd ers. Py ne j r i) ie Pe Oa a . ga nas /2 econ 4 Sea “a ow, Pee ee ee ea th OO > es oe —) Cie” es 2 . = . . PK se only - = on aw = 5 a em \_ ll st i, og — eS Pe ne en a? ond : pane ; ee” AB = 3 ey S eee e f Sima - ae BJ ie = PED ali a ae Be i \_ Mice q \* a : § sa, ay > “Se, NS eee] e& ft mo aie ve — = i” te ee > < = [a es ge = ms z SN Se . eo =o s r > a Lf A a oth ae Bee Ba a os ra => Se. a — <i Se wi 4 eo et eG 2. Ne ene ae ee gi eg a tee OP ee a Scar ae “ C Pay) ae Ng es 5 pee is ’ a ail wc s és =i a @ ae A 7” eae Rf Sy Cen) ae re ie ” NN S& ULE << ee — cea

i - i Sa Le? ae a & ; vr: E Ess! ee ~Ses Le ths AS > a 4 ado > wy, — 7 me Phot & f gor 5 we prct P Be o oy I irctiy icc, o Fanelli ess cae — yas : a r Ae S hea Ce Se ee Oho Gs ; » ee eS cf. 5 Be Ee EGBG ean eer : EE SSeS a a

Image: [area\\_25\\_republican\\_committee\\_bgimg\\_20220211.png](#)

Text:

Image: [vicky\\_l\\_cook\\_for\\_fork\\_supervisor\\_bgimg\\_20220211.jpg](#)

Text: ya ey Soe Noe i i | yn N aS" \ gay . \ a

Image: [donald\\_trump\\_and\\_henry\\_extra\\_20220219.jpg](#)

Text: q ee at oda ye a \* ne. zi a ; ie v = KA j XY Taw iM : F p |, i A Pp Z Ce THT f Tn CELLET TEE eeerccer Be NN : ee er, ae f a P al vV: SAM r ~ameniiicasesaa ELE | . Pa fe, 4 - ' a CG ' ee ORS ae - Om oF ve | (ar) ° a wes eee on, ane I " J oe oN. g : A ae o ee he fab . , ; | @r hs i Z Ke é 1 S ( \_ a iy C2 . — Yj, fl. 3 x \ B j} Yo at gy) dl % I ie ~ \_ Pas V 4 i} / é 5 s & ' \_ Wy RL j <<] Be rN << e eae A oe . , rl are — are " es am SS GE ay \_ A Wwe . = —

Image: [jasinski\\_for\\_senate\\_bgimg\\_20220211.jpg](#)

Text: ) eS el “ee DT — ss: ‘ 2 es = = ; 3 ‘ . Pe>) Tee Pons aN i gd is J 4 Bt a me < — - ~ © ane =. , ~ 4 te se , : r & o : — 2a A — \* . ‘ ; “ as a / o : i. ] a “s | Z ’ # / / os a — Fhe

Image: [president\\_trump\\_at\\_the\\_white\\_house\\_summit\\_on\\_human\\_trafficking\\_extra\\_20220219.jpg](#)

Text: 1 : ies, : " a ieee | . i + \_ c., & F a e ? i , “ip , p a i > 4 st Sak Rt a) Jee NS ee EE ' at 5 ME Zi gy Oe. OY B2ete Seat ; 3 “i R : a2 te é ' . 4 oe ; ‘ Rg ” ie 4 se 5 a 7 : \ ~ fe a = (a, # — 4 \_ , i 7 a 44 ee ar Je ' a yey - ' / : . > Chl ats : gee a es a 3 e Z 5 P ” ie Fi j % . = ' / e Ne 2 ete Sh: Ha ; > | : i iz wm rf - ee ae mo ” . — x a : ra ye Pe 2a BS a, ‘ .., i a er 3 oe re a 3 5 iz ie bs Ww ie A : f ye ee B & 3 RS es 5 fl : A 77 f 6 Ke w ee a? 4 .. 4 P ... re ie A 7°? 7 Y Ie L a a . a © i > . : . F ye . A b . 2 : ae ae J a : a : . eo a rn : - ae Bee Be Sb ee - % i if z ys : “i ge om ie os Ss 4 : a f' J q a3 4 , \* oS ie A a # Le 2 , RNIN J ; a: { . 2 | q me : : oS ' ie fd ; et fo pS : : ny : | ue he hanes \ | Ss = . ‘ i A ly iN 7 i ; | a — o 2 yn ba Po oe : 1 \*% 3 ‘ c < A We Pe ee “ ~ - / = . | c wee f A Hi4 ; : 2 a . A a a ‘ , . : - uf & ve a e 4 Bes Lt ae ~ - E i oe Bm: E / RS oe — Nae e ee. CCl Ba wee( ih)! 7m Ne) wi = , : Abe ty, Rm ae Vs BW e : es PY) | ame : “Se : 4 A i , . : : ‘ “= Fo E \ : | | A = Ar : a | a s 7 : ‘ 0 i : ae Ba. : Lj | , Ft . Mam \ KD yo. : | aX hice a aoa mn GA a sg P - 4 SS “? « as \_ tS = Se P 2 : ' i ie ie , . > a Poy oa.

Image: [harmeet\\_dhillon-speaks\\_at\\_the\\_white\\_house-27s\\_social\\_media\\_summit\\_extra\\_20220219.png](#)

Text: tt Ke Oe ag bel / wy oa. i . a | a oe \*x iG ae nee ae a ae Fens, ” + “ ied P= orig ee ) +a ee x x J : A - ov ae ar LI . c VR Tf hi CO ‘Y Ae or iY

Image: [kevin\\_cramer\\_bgimg\\_20220211.jpg](#)

Text: ~ = " wis = ‘ -\_ . (ae 7 a AQ We . —

Image: president\_donald\_j\_trump,\_senator\_tom\_cotton,\_and\_senator\_david\_perdue,\_august\_2,\_2017\_extra\_20

Text: Bi c Ja f | : if ae ¥ | ) | : lee aC ne f : A | | . ae o> &§ oo oe / | i i 4 pa ete a” . } Se x > ot , 7 dah » ’ 3 ej | N/ Bg: M ii i | , MW, : : i | pe: / | ae fl 7 ; \ I Tio Nae a \* , \* Jf ons Ly A - o 1% ’ = . > > ae ey

tesseract also returns its level of confidence (a percentage) in each identified word.

```
combined[["average_confidence"]] <-  
  sapply(combined[["ocr_output"]], function(x) {  
    if (nrow(x) > 0) {  
      mean(x[["confidence"]])  
    } else {  
      NA_real_  
    }  
  })
```

The average is 59.354.

The mean average confidence is quite different for observations where the classifiers disagreed.

```
result <- broom::tidy(t.test(average_confidence ~ disagreements, data = combined))[, 1:5]  
colnames(result)[2:3] <- c("No Disagreement", "Disagreement")  
present_table(result)
```

estimate	No Disagreement	Disagreement	statistic	p.value
30.3	62.5	32.3	14.9	0

Interestingly, it is not much higher for cases where only the OCR prediction was correct than for those where only the CNN prediction was correct.

```
tapply(combined[["average_confidence"]], combined[["prediction_category"]],  
  mean,  
  na.rm = TRUE  
) %>%  
  present_table()
```

	both	cnn	neither	ocr
	62.5	28.3	24.2	36.6

I also check the number of recognizable words tesseract identified in each observation by lemmatizing the words and confirming their presence in a reference of English words

```

library(quanteda)
words <- union(tidytext::stop_words[["word"]], tidytext::parts_of_speech[["word"]])
path <- udpipe::udpipe_download_model(language = "english")[["file_model"]]
English <- udpipe::udpipe_load_model(file = path)

# See https://stackoverflow.com/questions/46731429/quanteda-fastest-way-to-replace-tokens-with-lemma-fr
combined[["valid_words"]] <- sapply(
  combined[["ocr_output"]],
  function(x) {
    if (nrow(x) > 0) {
      tokenized <- tokens(x[["word"]]),
      remove_numbers = TRUE, remove_punct = TRUE,
      remove_url = TRUE, remove_symbols = TRUE
    } %>%
      tokens_tolower() # %>%
      # tokens_wordstem() %>%
      # tokens_replace(pattern = lexicon::hash_lemmas$token, replacement = lexicon::hash_lemmas$lemma)
      raw <- unlist(tokenized)
      if (length(raw) == 0) {
        out <- 0
      } else {
        # Some lemma sequences longer than original word sequence - odd, but likely not important
        lemmas <- udpipe::udpipe_annotate(English,
          unlist(tokenized),
          tagger = "none"
        )[["x"]]
        out <- sum(unname(lemmas) %in% words)
      }
    } else {
      out <- NA
    }
    out
  }
)
# Clean up unneeded file
if (file.exists(path)) invisible(file.remove(path))

```

Most observations seemed to contain at least some valid words.

```

summary(combined[["valid_words"]]) %>%
  present_table()

```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
0	1	3	11.2	6	862	24

Lastly, overall accuracies are almost identical.

```
colMeans(combined[, c("pred_class", "ocr_pred_class")] == combined[["actual_class"]]) %>%  
  present_table()
```

pred_class	ocr_pred_class
0.946	0.948

## Conclusions

OCR seems capable of accurately reading text in at least some cases, but often outputs gibberish. It could be improved by experimenting with parameters, although there is no guarantee that would prove worthwhile. The convolutional network is surprisingly effective, given the mixed training data. I should create more examples of images with Trump to obtain a more accurate view of performance. While further fine-tuning is clearly necessary, I think the strategy of combining a neural network with OCR remains workable. One issue I did not address is how to combine outputs from both methods into a single prediction.