

# FREQUENCY, LENGTH, AND DIVERSITY EFFECTS ON WORD RECOGNITION

A replication of the Howes and Solomon experiment was done through an updated digital approach, analyzed and compared through the lenses of contemporary work, based on updated databases, contextual and semantic information, and initially ignored confound variables such as word length.

## Background

### Howes and Solomon (1951)

In 1951 Howes and Solomon conducted a study focused on the relationship between word-frequency and the speed in which people recognize words.

They use a tachistoscope to measure the visual duration threshold necessary for a participant to correctly recognize a word, the word's frequency in the English language had been previously ranked in a word-frequency count from the work of Thorndike-Lorge Teacher's Word book of 30, 000 words.

Three cases were used to reduce the idiosyncrasies of any individual word count: the Lorge Magazine Count, the Thorndike-Lorge Semantic Count and finally the geometric mean of the frequencies for both previous counts. The Lorge Magazine Count was based on four and a half million words from the issues of current magazines, while Thorndike-Lorge Semantic Count was based on the Encyclopedia Britannica, Bartlett's Familiar Questions, The Library Digest and other books.

In their experiment, Howes and Solomon used a modified tachistoscope to expose a word to a participant for a brief amount of time, and subsequently extending the period by intervals of 10 milliseconds until the subject could correctly recognize the word.

Howes and Solomon recognized a strong relationship between the duration needed for recognition of a word and the word's relative word frequency. Their findings suggested that words with high frequency require shorter exposure before being recognized, as opposed to words with lower frequency.

Correlations found between the two variables ranged from  $-.68$  to  $-.75$  in Howes and Solomon's main experiment, and after some corrections applied, these were raised to  $-.76$  to  $-.83$  after their second experiment, and even when these numbers suggest a strong correlation However Howes and Solomon express the relationship could benefit from a more appropriate word frequency count than the available for the experiment at the moment "For a few words, the frequencies based on the two counts are so disparate that even their ranks in the hierarchy of relative frequencies may be called into question".

Howes and Solomon suggested that future work should focus experiments in which variables like word-frequency and practice are controlled and the physical characteristics of words are varied systematically can determine their precise effects on duration thresholds.

### McGinnies et al. (1952)

In 1952 McGinnies, Comer and Lacey observed word length to correct for Howes and Solomon's original experiment. A clear absence of this relationship considering Howes and Solomon effort to determine the

effect of other parameters like type, patterning of characters and syllable length on duration thresholds. McGinnies et al. (1952) determined a relationship between the word length and word-frequency on recognition thresholds.

McGinnies' experiment selected words of 5, 7, 9 and 11 characters and in each of these categories words with frequencies of 10, 100, 200, 300, and 400 per million. Each category had twenty words randomized for position. A trial before the experiment had four words with different frequencies and lengths to reduce the practice effect.

Participants were students from the University of Alabama, twenty students were recruited for the study, to minimize differences in visual acuity, participants whose threshold duration exceeded 20 milliseconds were excluded from the experiment.

After using partial correlation McGinnies et al. (1952) statistically removed the effect of word length in the replication of the experiment of Howes and Solomon (1951). Still a positive relationship was observed between the word length and the word frequency, increased frequency lowered the thresholds of long words significantly more than those of short words, McGinnies states that the relationship between length and frequency should be considered before stating its individual effects on threshold duration, ultimately contradicting Howes and Solomon's findings on the lack of relationship between word length and thresholds.

## Brysbaert and New (2009)

Kucera and Francis (1967) had been for many years the preferred word count database for word frequency, in 2008 it had been cited 215 times despite being outdated and with a lack of contextual significance. In 2009 Brysbaert and New, believed that Kucera and Francis database, built on adult reading material, was constructed on a relatively small sample of 1.04 million words, as it consistently underperformed in predicting reaction times.

To correct this and other problems with Kucera and Francis popular database, Brysbaert and New developed the SUBTLEXus a database based on subtitles from four different sources: U.S films from 1900-1990 (2,046 files), U.S. films from 1990-2007 (3,218 files), and U.S. television series (4,575 files). For a total of 51.0 million words.

Ultimately the new SUBTLEXus improved the prediction of response time from 58% to 63%, and while Kucera and Francis account for 18% of the variance in the accuracy data and 32% in the latency data, SUBTLEXus improves this percentages to 30% and 44% accordingly.

## Yap & Balota (2012)

Yap & Balota examined differences between individuals who contributed to the English Lexicon Project, their argument that variation in reading skill modulates word recognition performance. Following the methodological aspects of the English Lexicon Project available in Balota et al. (2007). Yap & Balota used an online behavioral database containing nearly four-million-word recognition trials from over 1,289 participants with 470 providing data for the speeded pronunciation task and 819 for the lexical decision task.

Yap & Balota observed considerable within- and between-session reliability across distinct sets of items, in terms of overall mean response time, and sensitivity to underlying lexical dimensions.

In addition, higher vocabulary knowledge was associated with faster, more accurate word recognition performance, attenuated sensitivity to stimuli characteristics, and more efficient accumulation of information. participants who showed more influence of one variable also showed more influence of other variables

## Jones, Johns and Reccia (2012)

Jones, Johns and Reccia observe that counting contexts gives a better quantitative fit to human lexical decision and naming data than counting raw occurrences of words. However, this approach ignores the information redundancy of the contexts in which the word occurs, a factor referred to as semantic diversity.

Jones, Johns and Reccia demonstrate the importance of contextual redundancy in lexical access, suggesting that contextual repetitions in language only increase a word's memory strength if the repetitions are accompanied by a modulation in semantic context. They introduce a cognitive process mechanism to explain the pattern of behavior by encoding the word's context relative to the information redundancy between the current context and the word's current memory representation.

Jones, Johns and Reccia computed word frequency, document count, and semantical diversity count from three corpora: (a) the Touchstone Applied Science Associates (TASA) corpus (Landauer & Dumais, 1997), (b) a Wikipedia corpus (Reccia & Jones, 2009), and (c) a New York Times (NYT) corpus (Jones & Mewhort, 2004). The model gives a better account of identification latency data than models based on either raw frequency or document count, and produces a better-organized space to simulate semantic similarity.

They observed that when words with equivalent document counts are considered, those that occur in more semantically distinct contexts see a larger latency savings when compared to those that occur in redundant contexts.

## Experimental Design

Using a small sample of nine students from the University of Carleton, a mix of native and nonnative English speakers, a replica of the Howes and Solomon's (1951) test was designed, using the PsychoPy program to emulate the functionality of a modern day tachistoscope.

To create a comprehensive analysis the experiment would use the databases of Kucera and Francis (KF), SUBTLEXus and Jones et al. (2012) missing words from any of the databases would be assigned a word frequency (WF) of 1. We expect to find the SUBTLEXus a more reliable database based on the significantly bigger sample, its account for variance and contextual and semantic significance.

A database was created based on Howes and Solomon's initial experiment design, it contained 60 words and 60 made-up non-words, once a word was exposed the participants were expected to choose between word or non-word by pressing one of two keys.

A trial was run beforehand with 10 words presented to eliminate the practice effect, subsequently each participant was presented with 60 words randomized for position, a word would be presented for 100

milliseconds and immediately afterwards a dotted image for one second, the program would then wait for the participant to respond if they thought they read either a word or a non-word.

A second trial was run immediately afterwards; the background and font color was inverted to uncover insights of its effect on word processing. The complete duration of the experiment took an approximate of ~15 minutes for 120 trials plus 10 practice words.

The program would record their accuracy, either a zero or a one for mislabeling or labeling between word or non-word. Response time was also measured and recorded.

## Experiment 1

We compare Kucera & Francis word frequency values with SUBTLEX database. Since every value for correctly or incorrect label of a word or non-word created values of zero or one, we calculated the average number for each word group and used this mean as the basis for our correlation.

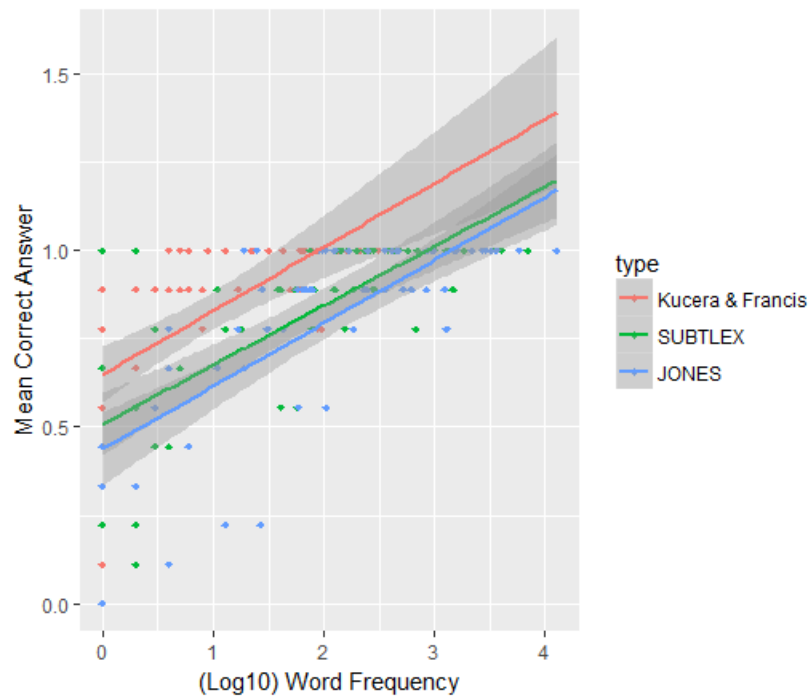


Figure 1 Mean Accuracy against Logged Word Frequency

We hypothesized both databases to perform similarly in accuracy, In *Figure 1* we observe both databases line parallel to each other suggesting a similar correlation. When performed, a correlation test we discover a stronger positive correlation in accuracy to the SUBTLEXus of 0.7259 as than the one obtained by Kucera & Francis of 0.5956, confirming the effort from Brysbaert and New of creating a more reliable updated database.

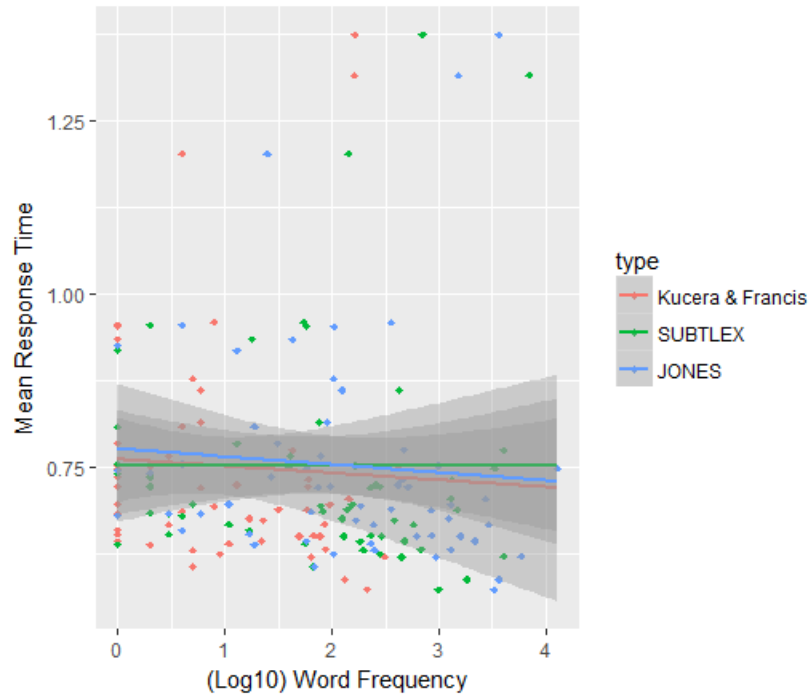


Figure 2 Mean response time against Logged word frequency.

We calculated the mean response time and compared to the logged frequency of both databases, we expected a negative correlation, Contrary to our predictions Kucera and Francis' word frequency *Figure 2* showed a stronger negative correlation while SUBTLEXus showed less correlation than expected. SUBTLEXus correlation coefficient measured -0.1868 against the -0.2633 of Kucera and Francis.

To further analyze these results, we look at the effect of word length on response time by removing it using partial correlation. However, the lack of a statistically significant p values makes the results negligible. Kucera & Francis obtained a partial correlation value of -0.0561 with a p value of 0.6728; SUBTLEX us obtained a partial correlation value of -0.00474 with a p value of 0.9715.

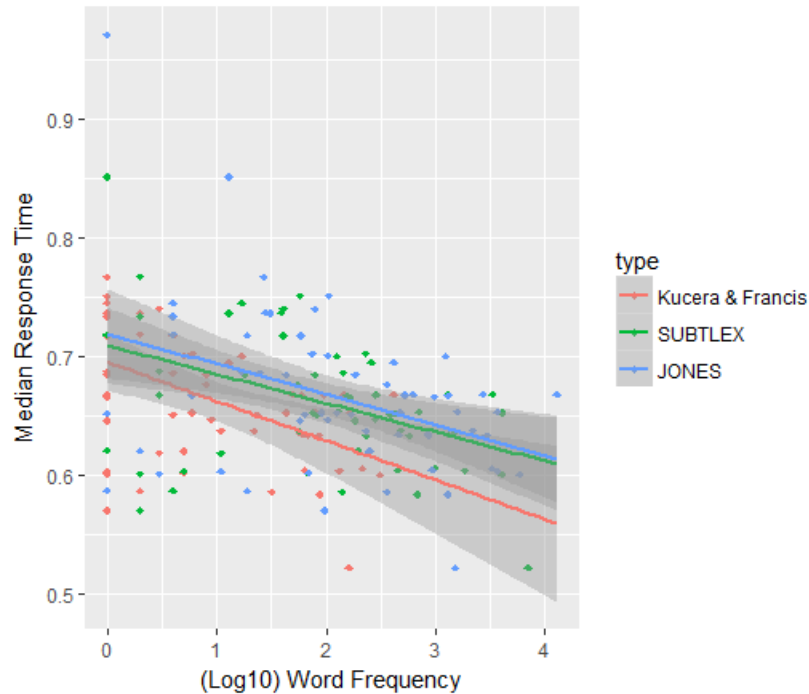


Figure 3 Median Response time against logged word frequency

We felt with such variable responses and with the possibility of outliers a median per group could give us a more insightful result than the calculated mean. In *Figure 3* we compared the median of each word group to the logged frequency of both databases, in an effort to correct for the variability in the data presented in *Figure 2*. Again, *Figure 3* shows a stronger negative correlation between the variables when using the Kucera and Francis database with a correlation of  $-0.4022$  against the SUBTLEXus of  $-0.3878$ .

Again, we look at the effect of word length on response time by removing it using partial correlation. However, in this instance a study of the median creates a much stronger p value for both databases. Kucera & Francis obtained a partial correlation value of  $-0.3991$  with a p value of  $0.00173$ ; SUBTLEX us obtained a partial correlation value of  $-0.3820$  with a p value of  $0.00283$ . The similarity in these results leads us to believe that word length has an approximate effect of  $-.39$  in either correlation.

## Experiment 1 Discussion

Contrary to our initial predictions, the observed data suggests Kucera and Francis word frequency database to be the most reliable than the favorited SUBTLEXus database, presenting stronger correlations in all three comparisons. Mean response time seems to be an unreliable representation of the data when such variance exists in the data. This could also be the product of a significantly small sample size obtained for the experiment.

## Experiment 2

To expand our understanding of all the variables involved in the effect of recognition thresholds we look at the effects of contextual diversity, semantic diversity and word length and its effects on the relationship between word frequency and either mean correct answer, mean response time and median response time.

Using Jones et al. (2012) Wikipedia's database we compare word frequency against different variables, finding some insightful relationships:

*Table 1 – Jones Word Frequency against Mean Response Time*

<b>REMOVING</b>	<b>PARTIAL CORRELATION</b>	<b>P VALUE</b>
<b>CONTEXTUAL DIVERSITY</b>	-0.09966345	0.4526389
<b>SEMANTIC DIVERSITY</b>	-0.1107342	0.4037519
<b>LENGTH</b>	-0.08165402	0.5386836

*Table 2 – Jones Word Frequency against Mean Response Time*

<b>REMOVING</b>	<b>PARTIAL CORRELATION</b>	<b>P VALUE</b>
<b>CONTEXTUAL DIVERSITY</b>	-0.3251564	0.01197684
<b>SEMANTIC DIVERSITY</b>	-0.1607092	0.2240075
<b>LENGTH</b>	-0.379556	0.00302821

We replicated Jones et al. experiment, and even though we find a stronger relationship in semantic diversity than contextual diversity or length, the p value obtained for the mean response time makes the finding negligible.

## DISCUSSION

Overall the results from the experiments were contradictory to our predictions, Kucera and Francis obtained stronger correlations in most comparisons, suggesting the outdated database to be better at predicting word processing than the more contemporary and vastly bigger SUBTLEXus database.

Most of these incorrect findings could be the product of a small sample size, we should also consider words that have been mislabeled and their effect on response time, instead of accepting response times regardless of their correct or incorrect labeling. The imputation of missing words from the SUBTLEXus database could also have a confound effect in the results.

Such inconsistent findings suggest a replication should be performed under more rigorous experiment design and with a much larger participant sample. Other factors remain unaccounted for such as lowercase versus uppercase, and the effect of bright or dark background and font colors.

## CONCLUSION

The results presented in this study contradict recent findings, however the lack of rigorosity during experimentation and small sample make these results interesting at best but inconsequential unless a secondary study is performed.

## Appendix 1 – Code

```
#~~~~~#
##### ADD LIBRARIES
#~~~~~#
library(readr)
library(pastecs)
library(ggplot2)
library(lsr)
library(haven)
library(readxl)
library(reshape2)
library(psych)
library(MASS)
library(car)
library(tidyverse)
library(ppcor)
library(data.table)
library(plyr)

#~~~~~#
##### STEP 1 - IMPORTING DATA
#~~~~~#

#import FINAL
#import SUBTLEX dataframe either full RAW or only the words used
#import Kucera and Francis from CuLearn

FINAL <- read_csv("C:/Users/Jerbo/Google Drive/Carleton/HCIN 5400 Stats/Final Project/Final.csv")
SUBTLEX_RAW <- read_csv("C:/Users/Jerbo/Google Drive/Carleton/HCIN 5400 Stats/Final Project/SUBTLEX_RAW.csv")
KF_RAW <- read_csv("C:/Users/Jerbo/Google Drive/Carleton/HCIN 5400 Stats/Final Project/KF_RAW.csv")
JONES_RAW <- read_csv("C:/Users/Jerbo/Google Drive/Carleton/HCIN 5400 Stats/Final Project/WF_SD_CD.csv")

#~~~~~#
##### STEP 2 - CLEANUP
#~~~~~#

# Remove the rows of all the NON-Words

FINAL<-FINAL[!(FINAL$word_nonword=="nonword"), ]

# Remove some useless columns
FINAL$word_nonword <- NULL
FINAL$corrans <- NULL
FINAL$trials.thisIndex <- NULL
FINAL$X1 <- NULL
FINAL$gender..m.f. <- NULL

#~~~~~#
##### STEP 3 - ADD COLUMNS
#~~~~~#

# Add NEW columns to a dataset
# Counts the length of the word
FINAL["Length"] <- nchar(FINAL$strings)
FINAL["MeanCA"] <- ave(FINAL$resp.corr, FINAL$strings)
FINAL["MeanRT"] <- ave(FINAL$resp.rt, FINAL$strings)
FINAL["MedianRT"] <- NA

# ADD KUCERA AND FRANCIS COLUMNS WITH NAMES
FINAL["KF_WF"] <- NA
FINAL["KF_WF_LOG"] <- NA

# ADD SUBTLEX COLUMNS WITH NAMES
FINAL["SUBTLEX_WF_RAW"] <- NA
FINAL["SUBTLEX_WF_RAW_LOG"] <- NA
FINAL["SUBTLEX_WF_100K"] <- NA
FINAL["SUBTLEX_WF_100K_LOG"] <- NA
FINAL["SUBTLEX_CD"] <- NA
```



```

FINAL["SUBTLEX_CD_LOG"] <- NA

# ADD SUBTLEX COLUMNS WITH NAMES
FINAL["JONES_WF"] <- NA
FINAL["JONES_WF_LOG"] <- NA
FINAL["JONES_CD"] <- NA
FINAL["JONES_CD_LOG"] <- NA
FINAL["JONES_SD"] <- NA
FINAL["JONES_SD_LOG"] <- NA

#~~~~~#
# STEP 4 - IMPORT VALUES
#~~~~~#

# Populate the newly created columns
# With values from the KF and the SUBTLEX files

# Kushera and Francis Word Frequency
FINAL$KF_WF <- KF_RAW$KF_Freq [match (FINAL$strings, KF_RAW$KF_Strings)]
# Subtlex Word Frequency
FINAL$SUBTLEX_WF_RAW <- SUBTLEX_RAW$FREQcount [match (FINAL$strings, SUBTLEX_RAW$Word)]
# Subtlex Contextual Diversity
FINAL$SUBTLEX_CD <- SUBTLEX_RAW$SUBTLCD [match (FINAL$strings, SUBTLEX_RAW$Word)]
# JONES Word frequency
FINAL$JONES_WF <- JONES_RAW$WF [match (FINAL$strings, JONES_RAW$STRING)]
# JONES Contextual Diversity
FINAL$JONES_CD <- JONES_RAW$CD [match (FINAL$strings, JONES_RAW$STRING)]
# JONES Semantic Diversity
FINAL$JONES_SD <- JONES_RAW$SD [match (FINAL$strings, JONES_RAW$STRING)]

# MEDIAN
MEDIAN_RT <- aggregate(FINAL$resp.rt, list(FINAL$strings), median)
FINAL$MedianRT <- MEDIAN_RT$x [match (FINAL$strings, MEDIAN_RT$Group.1)]

#~~~~~#
# STEP 5 - IMPUTATION
#~~~~~#

### FILL MISSING NA VALUES ###

# REPLACE ALL NA VALUES NOT FOUND IN THE SUBTLEX FOR A VALUE OF 1
FINAL$SUBTLEX_WF_RAW[is.na(FINAL$SUBTLEX_WF_RAW)] <- 1
FINAL$SUBTLEX_CD[is.na(FINAL$SUBTLEX_CD)] <- (.01)
FINAL$JONES_WF[is.na(FINAL$JONES_CD)] <- 1
FINAL$JONES_CD[is.na(FINAL$JONES_CD)] <- 1
FINAL$JONES_SD[is.na(FINAL$JONES_SD)] <- 1

# CALCULATE THE 100K WF
FINAL$SUBTLEX_WF_100K <- ( (1000000 * FINAL$SUBTLEX_WF_RAW) / (51000000) )

### LOG COLUMNS #####

# PASTE HERE <- The log 10 of this doc/column:
FINAL$KF_WF_LOG <- log10(FINAL$KF_WF)
FINAL$SUBTLEX_WF_RAW_LOG <- log10(FINAL$SUBTLEX_WF_RAW)
FINAL$SUBTLEX_WF_100K_LOG <- log10(FINAL$SUBTLEX_WF_100K)
FINAL$SUBTLEX_CD_LOG <- log10(FINAL$SUBTLEX_CD)

# FINAL$SUBTLEX_CD_LOG <- log10(FINAL$SUBTLEX_CD) # REVISE
FINAL$JONES_WF_LOG <- log10(FINAL$JONES_WF)
FINAL$JONES_CD_LOG <- log10(FINAL$JONES_CD)
FINAL$JONES_SD_LOG <- log10(FINAL$JONES_SD)

#~~~~~#
# STEP 6 - ACCURACY IN GGLOT
#~~~~~#

df <- data.frame(x=FINAL$KF_WF_LOG,
                 y=FINAL$MeanCA,
                 type='Kucera & Francis')
df <- rbind(df, data.frame(x=FINAL$SUBTLEX_WF_RAW_LOG,
                           y=FINAL$MeanCA,

```

```

        type='SUBTLEX')
df <- rbind(df, data.frame(x=FINAL$JONES_WF_LOG,
                          y=FINAL$MeanCA,
                          type='JONES'))

ggplot(df, aes(x, y, group=type, col=type)) +
  xlab("(Log10) Word Frequency") +
  ylab("Mean Correct Answer") +
  geom_point(shape=18) +
  geom_smooth(method=lm, se=FALSE, fullrange=FALSE)

#~~~~~#
#  MEAN RT IN GGLOT
#~~~~~#

df <- data.frame(x=FINAL$KF_WF_LOG,
                 y=FINAL$MeanRT,
                 type='Kucera & Francis')
df <- rbind(df, data.frame(x=FINAL$SUBTLEX_WF_RAW_LOG,
                          y=FINAL$MeanRT,
                          type='SUBTLEX'))
df <- rbind(df, data.frame(x=FINAL$JONES_WF_LOG,
                          y=FINAL$MeanRT,
                          type='JONES'))

ggplot(df, aes(x, y, group=type, col=type)) +
  xlab("(Log10) Word Frequency") +
  ylab("Mean Response Time") +
  geom_point(shape=18) +
  geom_smooth(method=lm, se=TRUE, fullrange=TRUE)

#~~~~~#
#  MEDIAN RT IN GGLOT
#~~~~~#

df <- data.frame(x=FINAL$KF_WF_LOG,
                 y=FINAL$MedianRT,
                 type='Kucera & Francis')
df <- rbind(df, data.frame(x=FINAL$SUBTLEX_WF_RAW_LOG,
                          y=FINAL$MedianRT,
                          type='SUBTLEX'))
df <- rbind(df, data.frame(x=FINAL$JONES_WF_LOG,
                          y=FINAL$MedianRT,
                          type='JONES'))

ggplot(df, aes(x, y, group=type, col=type)) +
  xlab("(Log10) Word Frequency") +
  ylab("Median Response Time") +
  geom_point(shape=18) +
  geom_smooth(method=lm, se=TRUE, fullrange=TRUE)

#~~~~~#
##### CORRELATION COEFICIENT
#~~~~~#

# Create a subset of FINAL and remove duplicates
# to analyse against means and medians
LeanSubset <- FINAL
LeanSubset <- subset(FINAL,!duplicated(LeanSubset$strings))

## LOGGED CORRELATION COEFICIENT FOR ACCURACY
cor(LeanSubset$KF_WF_LOG,LeanSubset$MeanCA)
cor(LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$MeanCA)
cor(LeanSubset$JONES_WF_LOG,LeanSubset$MeanCA)

## LOGGED CORRELATION COEFICIENT FOR MEAN RESPONSE TIME
cor(LeanSubset$KF_WF_LOG,LeanSubset$resp.rt)
cor(LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$resp.rt)
cor(LeanSubset$JONES_WF_LOG,LeanSubset$resp.rt)

## LOGGED CORRELATION COEFICIENT FOR MEDIAN RESPONSE TIME

```

```

cor(LeanSubset$KF_WF_LOG,LeanSubset$MedianRT)
cor(LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$MedianRT)
cor(LeanSubset$JONES_WF_LOG,LeanSubset$MedianRT)

##### PARTIAL CORRELATION
# MEAN CORRECT ANSWER accounting for LENGTH
pcor.test(LeanSubset$MeanCA,LeanSubset$KF_WF_LOG,LeanSubset$Length)
pcor.test(LeanSubset$MeanCA,LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$Length)

# MEAN RESPONSE TIME accounting for LENGTH
pcor.test(LeanSubset$MeanRT,LeanSubset$KF_WF_LOG,LeanSubset$Length)
pcor.test(LeanSubset$MeanRT,LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$Length)

# MEDIAN RESPONSE TIME accounting for LENGTH
pcor.test(LeanSubset$MedianRT,LeanSubset$KF_WF_LOG,LeanSubset$Length)
pcor.test(LeanSubset$MedianRT,LeanSubset$SUBTLEX_WF_RAW_LOG,LeanSubset$Length)

#~~~~~#
##### PARTIAL CORRELATION COEFICIENT
#~~~~~#

# MEAN CORRECT ANSWER
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanCA,LeanSubset[,c("JONES_CD","JONES_SD","Length")])
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanCA,LeanSubset$JONES_CD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanCA,LeanSubset$JONES_SD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanCA,LeanSubset$Length)
# MEAN RESPONSE TIME
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanRT,LeanSubset[,c("JONES_CD","JONES_SD","Length")])
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanRT,LeanSubset$JONES_CD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanRT,LeanSubset$JONES_SD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MeanRT,LeanSubset$Length)
# MEDIAN RESPONSE TIME
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MedianRT,LeanSubset[,c("JONES_CD","JONES_SD","Length")])
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MedianRT,LeanSubset$JONES_CD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MedianRT,LeanSubset$JONES_SD)
pcor.test(LeanSubset$JONES_WF_LOG,LeanSubset$MedianRT,LeanSubset$Length)

```