

# Measuring Readability in Financial Disclosures

*Journal of Finance*

Tim Loughran

Bill McDonald

University of Notre Dame

---

*“Just as the Black-Scholes model is a commonplace when it comes to compliance with the stock option compensation rules, we may soon be looking to the Gunning-Fog and Flesch-Kincaid models to judge the level of compliance with the plain English rules.”*

- SEC Chairman Christopher Cox in 2007

---

[Overview](#) | [Literature](#) | [Fog et al.](#) | [Data/Parsing](#) | [Readability](#) | [Results](#) | [Conclusion](#)

## The comparison is obvious:

### Black-Scholes:

$$C(s, t) = N\left(\frac{\ln\left(\frac{S}{K}\right) + \left(r + \frac{\sigma^2}{2}\right)(T - t)}{\sigma\sqrt{T - t}}\right) S - N\left(\frac{\ln\left(\frac{S}{K}\right) + \left(r - \frac{\sigma^2}{2}\right)(T - t)}{\sigma\sqrt{T - t}}\right) Ke^{-r(T-t)}$$

### Gunning-Fog (Fog Index):

$$0.4(w + p)$$

where:  $w$  = average # of words per sentence  
 $p$  = percent of “complex” words (> 2 syllables)

# What are we doing?

- We examine the relation between document readability and measures of the information environment (post-filing market model RMSE, earnings surprise, and analyst dispersion) using a large 10-K sample during 1994-2011.
- We define 10-K readability as how effectively management communicates valuation relevant information.

# Why is this important?

- Increasing popularity of textual analysis
- Researchers in this area often need a “readability” measure
- SEC continues to emphasize its Plain English initiative
- Most financial/accounting researchers use the Fog Index because it is the most widely recognized measure of readability from other disciplines.

# What do we find?

- We show that the Fog Index is an inappropriate measure of readability in the context of financial disclosures.
- We propose using the document file size (i.e., the number of megabytes required to store the document as reported on the SEC website) as a simple, and admittedly imperfect, proxy for readability.
- We show that *file size* relates to post-filing return volatility and other measures of the information environment in a manner consistent with the notion of readability.

# What do we find?

- *File size* relates positively and significantly to post-filing date return volatility (market model RMSE),  $|SUE|$ , and analyst dispersion after controlling for other variables.
- Fog does not consistently link to these variables and of its two components one is difficult to measure in 10-Ks, while the other is clearly misspecified in applications to financial disclosures.

# Literature

- Li (2008)
- Biddle, Hilary, and Verdi (2009)
- Lawrence (2013)
- Miller (2010)
- Dougal, Engelberg, Garcia, and Parsons (2012)
- Lehavy, Li, and Merkley (2011)



# The Fog Index

*Fog Index* = 0.4 ( average # of words per sentence +  
percent of complex words)

Complex words = words > 2 syllables

Note: Text with Fog Index values greater than 18 is generally considered unreadable.

---

Overview | Literature | **Fog et al.** | Data/Parsing | Readability | Results | Conclusion

# The Fog Index

Most readability measures, including the Fog Index were derived in the context of grade leveling text books. Consider:

## *The Cat in the Hat (Dr. Seuss):*

*“The sun did not shine. It was too wet to play. So we sat in the house all that cold, cold, wet day.”*

## *The Jungle (Upton Sinclair):*

*All day long the blazing midsummer sun beat down upon that square mile of abominations: upon tens of thousands of cattle crowded into pens whose wooden floors stank and steamed contagion; upon bare, blistering, cinder-strewn railroad tracks and huge blocks of dingy meat factories, whose labyrinthine passages defied a breath of fresh air to penetrate them; and there are not merely rivers of hot blood and carloads of moist flesh, and rendering-vats and soup cauldrons, glue-factories and fertilizer tanks, that smelt like the craters of hell-there are also tons of garbage festering in the sun, and the greasy laundry of the workers hung out to dry and dining rooms littered with food black with flies, and toilet rooms that are open sewers.*

## *Ulysses (James Joyce):*

*Molly Bloom’s soliloquy – two sentences: (1) 11,281 words (2) 12,931 words.*

---

Overview | Literature | **Fog et al.** | Data/Parsing | Readability | Results | Conclusion

# Parsing the 10-Ks

- All 10-Ks/10-K405s, excluding amended, from 1994-2011
- Eliminate
  - all ASCII-encoded text (binaries of pictures, pdfs, etc.)
  - all HTML
  - all tables (where >10% of nonblank characters are numbers)
- Parse into words

# Parsing for *Fog*

- Two components:

1. Average sentence length = # words / # sentences

- # of words from document parse
- # sentences
  - Remove abbreviations
  - Attempt to identify lists and punctuate
  - Divide number of words by the number of sentence termination characters

# Parsing for *Fog*

- Two components:

2. % complex words (words > 2 syllables)

- Syllabification

- “There is no completely agreed-upon definition of a syllable.” Jurafsky and Martin (2009, p. 223)
- Based on 15,000 words manually syllabified, the standard PERL Fathom package is only about 75% accurate.
- Our algorithm, based on Talburt (1986) and tuning, is > 90% accurate.

# What is readability?

- Prior readability research does not provide a precise definition of readability, but acknowledges that its meaning is context dependent.
- Does the SEC intend to make 10-Ks readable for 8<sup>th</sup> graders?
- Do all those SAT-like words in 10-Ks send you running to the dictionary?

# What is readability?

- The definition of Dale and Chall (1948), referred to by Tekfi (1987) as the classic definition and by DuBay (2007) as the most comprehensive definition, specifies “In the broadest sense, readability is the sum total (including interactions) of all the elements with a given piece of printed material that affects the success which a group of readers have with it.”

# What is readability?

- Tekfi (1987) concludes that readability is “ensuring that a given piece of writing reaches and affects its audience in the way the author intends.”



# What is readability?

- We define readability as the ability of individual investors and analysts to assimilate valuation relevant information from a financial disclosure.

# The information environment: Alternative measures

- Post-filing date return volatility
- Absolute value of earnings forecast error (a.k.a. SUE)
- Analyst dispersion

Table I  
Sample  
Creation

	<b>Dropped</b>	<b>Sample Size</b>
SEC 10-K files 1994 to 2011		188,413
Eliminate duplicates within year/CIK	2,930	185,483
Drop if file date < 180 days from prior filing	585	184,898
Drop if number of words < 2,000	8,298	176,600
CRSP PERMNO match	88,800	87,800
Reported on CRSP as ordinary common equity	4,376	83,424
Price on filing date day minus one $\geq$ \$3	13,338	70,086
Book-to-market COMPUSTAT data available and book value > 0	2,874	67,212
Post-filing date market model RMSE for days [6,28]	346	66,866
At least 60 days data available for market model estimates from event days [-252,-6]	147	66,719
Returns for days 0-1 in event period	12	66,707

Table II  
Variable Means  
by Time Period  
1994-2011

Variable	(1) 1994 to 2002	(2) 2003 to 2011	(3) 1994 to 2011
<u>Readability measures:</u>			
<i>Fog Index</i>	18.44	18.94	18.68
<i>Average words per sentence</i>	22.82	23.27	23.04
<i>Complex words</i>	23.28%	24.09%	23.67%
<i>File size (in megabytes)</i>	0.42	2.51	1.43
<u>Dependent variables:</u>			
<i>Post-filing RMSE</i>	3.45	2.26	2.87
<i>Abs(SUE)</i>	0.27	0.39	0.34
<i>Analyst dispersion</i>	0.14	0.21	0.19
<u>Control variables:</u>			
<i>Pre-filing alpha</i>	0.08	0.05	0.06
<i>Pre-filing RMSE</i>	3.54	2.65	3.11
<i>Abs(filing period abnormal return)</i>	0.04	0.03	0.03
<i>Size (market capitalization) in \$ millions</i>	\$2,257.56	\$3,680.13	\$2,946.42
<i>Book-to-market</i>	0.66	0.67	0.66
<i>NASDAQ dummy</i>	0.60	0.58	0.59
<i># of analysts</i>	4.19	5.60	5.05
Number of observations	34,405	32,302	66,707

Table III  
**Dependent variable:**  
**RMSE**  
**Independent variable:**  
**Fog Index**  
**and its**  
**components**

	(1)	(2)	(3)	(4)
<u>Readability measures:</u>				
<i>Fog Index</i>		0.017 (2.04)		
<i>Average words per sentence</i>			0.005 (4.02)	
<i>Complex words</i>				-0.006 (-0.77)
<u>Control variables:</u>				
<i>Pre-filing alpha</i>	-0.913 (-4.12)	-0.908 (-4.09)	-0.908 (-4.10)	-0.912 (-4.11)
<i>Pre-filing RMSE</i>	0.539 (12.07)	0.539 (12.01)	0.539 (12.08)	0.539 (12.18)
<i>Abs(filing period abnormal return)</i>	5.057 (17.52)	5.052 (17.57)	5.051 (17.57)	5.056 (17.53)
<i>Log(size in \$ millions)</i>	-0.105 (-5.45)	-0.105 (-5.45)	-0.105 (-5.52)	-0.105 (-5.50)
<i>Log(book-to-market)</i>	-0.133 (-2.41)	-0.133 (-2.41)	-0.133 (-2.41)	-0.133 (-2.40)
<i>NASDAQ dummy</i>	0.262 (3.37)	0.262 (3.38)	0.263 (3.38)	0.263 (3.45)
$R^2$	46.92%	46.93%	46.93%	46.92%

Table IV  
 First Quartile  
 of Most  
 Frequently  
 Occurring  
 Complex  
 Words in  
 10-Ks

Word	% of Total Complex Words	Cumulative %	Word	% of Total Complex Words	Cumulative %
FINANCIAL	1.51%	1.51%	ACCOUNTING	0.38%	16.76%
COMPANY	1.44%	2.95%	INCORPORATED	0.37%	17.13%
INTEREST	0.99%	3.94%	INCLUDED	0.37%	17.49%
AGREEMENT	0.78%	4.73%	COMPENSATION	0.36%	17.85%
INCLUDING	0.77%	5.50%	APPLICABLE	0.36%	18.21%
OPERATIONS	0.71%	6.21%	PRIMARILY	0.35%	18.56%
PERIOD	0.71%	6.92%	ACCORDANCE	0.35%	18.91%
RELATED	0.60%	7.52%	SIGNIFICANT	0.34%	19.26%
MANAGEMENT	0.60%	8.12%	SUBSIDIARIES	0.34%	19.60%
CONSOLIDATED	0.58%	8.70%	CUSTOMERS	0.34%	19.94%
INFORMATION	0.58%	9.28%	RESPECTIVELY	0.34%	20.28%
SERVICES	0.55%	9.83%	REGISTRANT	0.34%	20.62%
PROVIDED	0.55%	10.38%	OBLIGATIONS	0.33%	20.95%
PURSUANT	0.55%	10.93%	PROVISIONS	0.33%	21.28%
FOLLOWING	0.54%	11.47%	LIABILITIES	0.32%	21.60%
SECURITIES	0.54%	12.01%	ADDITION	0.32%	21.92%
APPROXIMATELY	0.52%	12.54%	OTHERWISE	0.32%	22.24%
REFERENCE	0.49%	13.03%	PROPERTY	0.32%	22.56%
OPERATING	0.47%	13.50%	EMPLOYEES	0.32%	22.87%
MATERIAL	0.46%	13.96%	BENEFIT	0.32%	23.19%
CAPITAL	0.43%	14.39%	REPORTING	0.32%	23.51%
EXPENSES	0.42%	14.81%	PRINCIPAL	0.31%	23.82%
CORPORATION	0.40%	15.21%	DEVELOPMENT	0.31%	24.13%
OUTSTANDING	0.40%	15.61%	REVENUE	0.30%	24.43%
ADDITIONAL	0.39%	16.00%	EQUITY	0.30%	24.73%
EFFECTIVE	0.38%	16.38%	INSURANCE	0.30%	25.04%

Table V  
Correlations  
of Alternative  
Readability  
Measures

	<i>Log(file size)</i>	<i>Fog Index</i>	<i>Average words per sentence</i>	<i>Complex words</i>	<i>Common words</i>	<i>Financial terminology</i>	<i>Vocabulary</i>
<i>Fog Index</i>	0.367						
<i>Average words per sentence</i>	0.316	0.885					
<i>Complex words</i>	-0.015	-0.089	-0.542				
<i>Common words</i>	-0.619	-0.465	-0.572	0.385			
<i>Financial terminology</i>	-0.407	-0.301	-0.372	0.254	0.781		
<i>Vocabulary</i>	0.668	0.497	0.596	-0.377	-0.970	-0.724	
<i>Log(# of words)</i>	0.712	0.560	0.652	-0.384	-0.916	-0.615	0.946

Table VI  
 Dependent  
 variable:  
 RMSE

A Comparison  
 of Log(file  
 size), Fog  
 Index, and the  
 Components  
 of Fog Index

	(1)	(2)	(3)
<u>Readability measures:</u>			
<i>Log(file size)</i>	0.073 (4.60)	0.069 (4.25)	0.076 (3.36)
<i>Fog Index</i>		0.006 (0.73)	
<i>Average words per sentence</i>			0.003 (0.75)
<i>Complex words</i>			0.010 (0.67)
<i>R</i> <sup>2</sup>	46.96%	46.97%	46.97%



**Table VII**  
Robustness

	(1)	(2)	(3)	(4)	(5)	(6)
	Dependent variables: <i>Abs(SUE)</i>			<i>Analyst Dispersion</i>	<i>Analyst Dispersion</i>	<i>Analyst Dispersion</i>
<u>Readability measures:</u>						
<i>Log(file size)</i>		0.046 (5.53)			0.023 (3.51)	
<i>Fog Index</i>			-0.003 (-0.82)			-0.000 (-0.02)
<u>Control variables:</u>						
<i>Pre-filing alpha</i>	-0.365 (-4.68)	-0.361 (-4.68)	-0.366 (-4.69)	-0.270 (-4.99)	-0.268 (-4.99)	-0.270 (-4.99)
<i>Pre-filing RMSE</i>	0.117 (5.47)	0.115 (5.36)	0.117 (5.48)	0.088 (4.84)	0.087 (4.80)	0.088 (4.84)
<i>Abs(filing period abnormal return)</i>	0.960 (4.84)	0.958 (4.84)	0.962 (4.85)	0.514 (4.36)	0.510 (4.33)	0.514 (4.33)
<i>Log(size in \$ millions)</i>	-0.054 (-6.21)	-0.061 (-6.71)	-0.054 (-6.22)	-0.022 (-4.15)	-0.026 (-4.78)	-0.022 (-4.16)
<i>Log(book-to-market)</i>	0.104 (4.85)	0.099 (4.64)	0.104 (4.85)	0.065 (4.09)	0.062 (3.96)	0.065 (4.10)
<i>NASDAQ dummy</i>	-0.089 (-6.56)	-0.086 (-6.58)	-0.089 (-6.53)	-0.053 (-3.81)	-0.051 (-3.87)	-0.053 (-3.81)
<i># of analysts</i>	-0.002 (-1.11)	-0.002 (-1.37)	-0.002 (-1.11)	0.006 (3.35)	0.006 (3.34)	0.006 (3.35)
Sample size	28,434	28,434	28,434	17,960	17,960	17,960
R <sup>2</sup>	23.30%	23.54%	23.30%	25.34%	25.52%	25.34%

**Table VIII**  
Complexity

	(1)	(2)	(3)
Dependent variables:	<i>Post-Filing RMSE</i>	<i>Abs(SUE)</i>	<i>Analyst Dispersion</i>
<i>Log(file size)</i>	0.084 (4.84)	0.044 (4.93)	0.022 (3.49)
<i>Business segment index</i>	0.156 (3.80)	-0.045 (-2.14)	-0.009 (-0.87)
Observations	50,739	22,783	14,516
$R^2$	47.27%	21.83%	23.98%

**Table IX**

Alternative  
Readability  
Measures – 24  
separate  
regressions

	Dependent Variable		
	(1)	(2)	(3)
<u>Readability Measures</u>	<i>Post-Filing RMSE</i>	<i>Abs(SUE)</i>	<i>Analyst Dispersion</i>
<i>Log(file size)</i>	0.073 (4.60)	0.046 (5.53)	0.023 (3.51)
<i>Fog Index</i>	0.017 (2.04)	-0.003 (-0.82)	-0.000 (-0.02)
<i>Average words per sentence</i>	0.005 (4.02)	0.002 (2.23)	0.002 (2.34)
<i>Complex words</i>	-0.006 (-0.77)	-0.014 (-5.75)	-0.009 (-4.08)
<i>Common words</i>	-1.295 (-4.56)	-0.614 (-5.49)	-0.437 (-4.47)
<i>Financial terminology</i>	-8.601 (-4.34)	-1.460 (-2.68)	-0.906 (-2.51)
<i>Vocabulary</i>	7.826 (4.72)	4.094 (6.31)	2.835 (5.68)
<i>Log(# of words)</i>	0.086 (4.27)	0.062 (6.55)	0.041 (4.79)
Number of observations	66,707	28,434	17,960

# Conclusion

- The *Fog Index* doesn't work for financial disclosures.
- We show why.
- We propose an alternative measure for readability –  $\text{Log}(\text{file size})$ .
- *File Size* is imperfect, but it is easy to calculate and relates to all of the information variables we test in a way consistent with the notion of readability.

# Conclusion

- As a simple proxy for readability use file size. If you are focusing on other specific attributes of readability consider some of the alternative measures we examine.

# Conclusion

- SEC – The SEC should focus less on style—which is undifferentiated in 10-Ks—and instead encourage managers to write more concisely.