

The Annual Report Algorithm:
Retrieval of Financial Statements and Extraction of Textual Information*

Jörg Hering¹

Friedrich-Alexander-Universität Erlangen-Nürnberg

First Version: October 4, 2016

Current Version: November 28, 2016

¹ Jörg Hering, Friedrich-Alexander-Universität Erlangen-Nürnberg, School of Business and Economics, Department of Accounting and Auditing, Lange Gasse 20, 90403 Nuremberg, Germany, phone: +49 911 5302 341, e-mail: joerg.hering@fau.de.

* I appreciate the valuable comments by Klaus Henselmann and Daniel Büchs.

The Annual Report Algorithm:

Retrieval of Financial Statements and Extraction of Textual Information

Abstract

U.S. corporations are obligated to file financial statements with the U.S. Securities and Exchange Commission (SEC). The SEC's Electronic Data Gathering, Analysis, and Retrieval (EDGAR) system containing millions of financial statements is one of the most important sources of corporate information available. The paper illustrates which financial statements are publicly available by analyzing the entire SEC EDGAR database since its implementation in 1993. It shows how to retrieve financial statements in a fast and efficient way from EDGAR. The key contribution however is a platform-independent algorithm for business and research purposes designed to extract textual information embedded in financial statements. The dynamic extraction algorithm capable of identifying structural changes within financial statements is applied to more than 180,000 annual reports on Form 10-K filed with the SEC for descriptive statistics and validation purposes.

Keywords: Textual analysis, Textual sentiment, 10-K parsing rules, Information extraction, EDGAR search engine

1. Introduction

Information Extraction (IE) can be defined as the process of “*finding and extracting useful information in unstructured text*” (Grant and Conlon 2006, 119). In contrast to Information Retrieval (IR), a technology that selects a relevant subset of documents from a larger set, IE extracts information from the actual text of documents (Wilks 1997, 1). Important sources for IE are unstructured natural language documents or structured databases (Mooney and Bunescu 2005, 3; Gaizauskas et al. 1997, 28). Since U.S. corporations are obligated by law to file financial statements on a regular basis with the U.S. Securities and Exchange Commission (SEC), the SEC’s Electronic Data Gathering, Analysis, and Retrieval (EDGAR) system containing millions of financial statements is one of the most important sources of corporate information available (Garcia and Norli 2012, 2; Grant and Conlon 2006, 120).

Unfortunately, most of the available textual data in the SEC EDGAR database is weakly structured in technical terms (Stümpert 2008, 357-364; Bovee et al. 2005, 25; O’Riain 2012, 35) especially prior to 2002 when the use of markup languages was less common (Loughran and McDonald 2014, 1650). A limited number of tagged items, formatting errors and other inconsistencies lead to difficulties in accurately identifying and parsing common textual subjects across multiple filings (Gerdes 2003, 16-17; Kambil and Ginsburg 1998, 92-93; Bovee et al. 2005, 20). These issues directly affect the ability to automate the extraction of textual information from SEC submissions (Gerdes 2003, 16; Davis and Tama-Sweet 2012, 811; Loughran and McDonald 2016, 1191-1192). Business data providers are offering expensive commercial products (e.g. AcademicEDGAR+, Edgar Pro, Intelligize). As research in the context of textual analysis is growing (e.g. Tetlock 2007; Loughran and McDonald 2011a; Jegadeesh and Wu 2013) the question occurs which particular financial statements and disclosures are publicly available for free, how to retrieve these corporate documents and how to decode the embedded textual information in order to be incorporated into investment decisions, trading strategies and research studies in financial economics (Garcia and Norli 2012, 1).

Today only a very limited amount of specific literature for extracting textual information from financial statements filed with the SEC and its EDGAR system is available (except Gerdes 2003; Stümpert et al. 2004; Grant and Conlon 2006; Engelberg and Sankaraguruswamy 2007; Cong, Kogan and Vasarhelyi 2007; Thai et al. 2008; Chakraborty and Vasarhelyi 2010; Hernandez et al. 2010; Garcia and Norli 2012; Srivastava 2016). This paper is based on neither of these because first, non-specialist technology is used to retrieve financial statements in an efficient way and secondly, the algorithm designed to extract textual information is platform-independent. The suggested method can compensate for expensive commercial products and help to replicate empirical research results. The paper shall serve as a technical guide on how to retrieve financial statements filed with the SEC and

how to decode the embedded textual information provided by the EDGAR system for business and research purposes.

The remainder of the paper proceeds as follows. Section 2 presents the amount and variety of corporate documents distributed by the SEC's electronic disclosure system. Section 3 demonstrates how to retrieve these documents from the EDGAR database. Section 4 describes the fundamentals of HyperText Markup Language and examines the electronic data provided by the SEC. Section 5 describes the fundamentals of regular expressions and specifies an algorithm to extract textual information embedded in financial statements. Section 6 validates the capabilities of the extraction algorithm. Section 7 presents descriptive statistics of annual reports filed with the EDGAR database. The last section concludes.

2. SEC's EDGAR database

Publicly owned companies, their officers and directors as well as major investors are obligated by law (Securities Exchange Act 1934, Section 2) to file various disclosures (forms) with the SEC (Gerdes 2003, 7). The main purpose of making certain types of corporate information publicly available is to improve the efficiency of security markets and to protect capital market participants (Garcia and Norli 2012, 2). *"The laws and rules that govern the securities industry in the United States derive from a simple and straightforward concept: all investors, whether large institutions or private individuals, should have access to certain basic facts about an investment prior to buying it, and so long as they hold it. To achieve this, the SEC requires public companies to disclose meaningful financial and other information to the public. This provides a common pool of knowledge for all investors to use to judge for themselves whether to buy, sell, or hold a particular security"* (SEC 2013). In order to protect investors, to maintain efficient capital markets and to improve access to publicly available corporate disclosures, the SEC developed the EDGAR database (Gerdes 2003, 9) and describes it as a system which *"performs automated collection, validation, indexing, acceptance, and forwarding of submissions by companies and others who are required by law to file forms with the U.S. Securities and Exchange Commission"* (SEC 2010).

Originally the EDGAR system was developed by the SEC as a pilot system for electronic disclosure in 1983. In order to test and evaluate EDGAR's performance the SEC requested electronic filings in 1994 after completing the phase-in of a mandated test group in December 1993 (the phase-in began on April 26, 1993) (SEC 2006; Kambil and Ginsburg 1998, 91; Pagell 1995, 56). As of May 6, 1996 the SEC obligated all public domestic U.S. companies (issuers) to file submissions electronically through the EDGAR system (SEC Release 34-36997 1996; Kambil and Ginsburg 1998, 91; Pagell 1995, 56; Grant and Conlon 2006, 121) except for certain filings made in paper because of a hardship exemption under Regulation S-T (SEC Regulation S-T 2016, Section 232.201; SEC,

2010). Filing for foreign private issuers (companies organized outside of the U.S.) and foreign governments via EDGAR (SEC 2006) became mandatory on May 14, 2002 (SEC Release 33-8099 2002). The Securities Exchange Act of 1934 (Securities Exchange Act 1934, Section 13(a), (b), Section 15(d)) empowers the SEC to require (periodic) reporting of information from publicly held companies (SEC 2013). In general, all public domestic companies with assets exceeding \$10 million and at least 500 shareholders become subject to Exchange Act reporting requirements (Securities Exchange Act 1934, Section 12(g)) alongside certain individuals (Gerdes 2003, 9). Among other disclosures, corporations with publicly traded securities are required (Securities Exchange Act 1934, Section 13(a), (b), Section 15(d)) to file annual and quarterly reports (Form 10-K, Form 10-Q) as well as current reports (Form 8-K) on an ongoing basis with the SEC and its EDGAR system (SEC 2013). Since by law these public corporate disclosures have to be accurate (Securities Exchange Act 1934, Section 13(i)) and represent a company's operations, they themselves represent a treasure trove of valuable information for investors and researchers (Gerdes 2003, 9; Engelberg and Sankaraguruswamy 2007, 3).

2.1 Underlying data in SEC's EDGAR database

In order to understand the amount and variety of corporate information (e.g. financial statements) distributed by the SEC, I retrieve and analyze all form index files since the implementation of the EDGAR system in 1993. The SEC EDGAR form index files list all publicly available disclosures made through the system in a certain quarter and sort the submissions by their particular filing form type. Table 1 reports the total number of submissions that have been made with the EDGAR system for each quarter and year since the introduction of the EDGAR database.

Table 1. Statistics on EDGAR submissions.

Year	Filings (Number)				Filings (Number)	Filings (%)
	Quarter 1	Quarter 2	Quarter 3	Quarter 4		
2016	307,416	239,528	---	---	546,944	3.42
2015	318,519	259,852	206,628	209,216	994,215	6.21
2014	311,679	252,333	212,352	220,328	996,692	6.23
2013	303,568	257,597	213,031	216,266	990,462	6.19
2012	309,453	246,776	203,723	214,985	974,937	6.09
2011	307,644	262,218	207,142	202,628	979,632	6.12
2010	300,538	255,180	203,920	220,070	979,708	6.12
2009	300,080	229,347	200,688	208,396	938,511	5.87
2008	328,709	267,722	220,732	219,669	1,036,832	6.48
2007	339,872	289,082	252,071	256,460	1,137,485	7.11
2006	335,577	278,960	232,131	249,956	1,096,624	6.85
2005	317,761	271,632	242,173	240,725	1,072,291	6.70
2004	312,029	253,021	217,726	241,435	1,024,211	6.40
2003	183,595	167,119	212,258	227,800	790,772	4.94
2002	125,189	108,013	97,533	118,149	448,884	2.81
2001	111,740	90,283	74,313	75,107	351,443	2.20
2000	116,209	81,129	72,571	72,053	341,962	2.14
1999	105,531	78,272	68,631	68,828	321,262	2.01
1998	106,666	73,830	67,234	65,570	313,300	1.96
1997	91,096	65,470	60,142	63,422	280,130	1.75
1996	49,925	47,659	50,641	54,389	202,614	1.27
1995	31,875	26,104	26,699	28,973	113,651	0.71
1994	20,879	16,500	13,066	15,016	65,461	0.41
1993	4	4	7	20	35	0.00
Filings (Number)	5,035,554	4,117,631	3,355,412	3,489,461	15,998,058	100.00
Filings (%)	31.48	25.74	20.97	21.81	100.00	

Notes: The table presents the total number of filings made on EDGAR for each year between 1993 and 2016. Each individual filing in a particular quarter is listed in an associated EDGAR form index file on the SEC server.

A tremendous amount of publicly available disclosures was filed with the SEC between 1993 and 2016. In total 15,998,058 filings were submitted to the EDGAR system in order to be publicly distributed. On average 31.48 percent (5,035,554) of these filings became available in the first, 25.74 percent (4,117,631) in the second, 20.97 percent (3,355,412) in the third and 21.81 percent (3,489,461) in the last quarter of each year since 1993. Most noticeable is the overall increase in total submissions through the EDGAR system reaching its peak in 2007 with more than 1.1 million disclosures for that particular year. By analyzing the index files more precisely, investors and researchers can gain an insight into the specific type of information the SEC is making publicly available through its EDGAR system (Garcia and Norli 2012, 2). Table 2 describes the most common filing (form) types filed with the EDGAR system.

Table 2. Statistics on EDGAR form types.

Rank	Form/Description	Submission Type	Filings (Number)	Filings (%)
1	Changes in ownership	4	5,850,937	36.57
2	Current report filing	8-K	1,376,248	8.60
3	5% passive ownership triggers amendments	SC 13G/A	587,711	3.67
4	Initial ownership report	3	538,228	3.36
5	Quarterly report	10-Q	522,906	3.27
6	Definitive materials	497	365,987	2.29
7	5% passive ownership triggers	SC 13G	344,030	2.15
8	Current report of foreign issuer	6-K	326,751	2.04
9	Change on a prospectus	424B3	254,046	1.59
10	5% active ownership triggers amendments	SC 13D/A	201,938	1.26
11	Changes in ownership amendments	4/A	197,612	1.24
12	Quarterly holdings, institutional managers	13F-HR	193,463	1.21
13	Annual report on ownership changes	5	186,884	1.17
14	Annual report	10-K	167,599	1.05
15	SEC-originated letters to filers	UPLOAD	159,065	0.99
16	Filer response letters	CORRESP	153,987	0.96
17	Proxy statements	DEF 14A	152,216	0.95
18	Registration management investment companies	485BPOS	151,903	0.95
19	Registration of securities, investment companies	24F-2NT	149,385	0.93
20	Offering of securities without registration	D	147,355	0.92
...
Total			15,998,058	100.00

Notes: The table presents the most frequent form types filed with the EDGAR system between 1993 and 2016. The first column ranks each filing type in descending order of total submissions. The second column gives a short description of each filing form type (Garcia and Norli 2012, 2). The third column lists the form codes used on EDGAR to identify a particular filing type made with the database. The next column contains the number of total submissions of a particular filing form type. The last column shows the amount of total submissions for each filing type in relation to all submissions made with the SEC EDGAR database.

The submission type most often filed with the EDGAR system since its implementation is Form 4. Between 1993 and 2016 5,850,937 filings report purchases or sales of securities by persons who are the beneficial owner of more than ten percent of any class of any equity security, or who are directors or officers of the issuer of the security (Garcia and Norli 2012, 2). The second most frequent submission type filed with the SEC is Form 8-K. 1,376,248 filings of this submission type are listed in the EDGAR index files. The current report filing is required by companies in order to inform shareholders about certain corporate events. These events of material importance for a company include information on significant agreements, impairments, changes in management etc. (Garcia and Norli 2012, 2). Important submission types for investors and researchers such as the annual report on Form 10-K have been submitted 167,599 times. Quarterly reports on Form 10-Q have been filed 522,906 times in total between 1993 and 2016. Another important submission type is Schedule 13G (SC 13G). Investors who are not seeking control over a firm (passive investors) must file this submission type as required by the SEC when crossing the five percent ownership threshold of a company (Garcia and Norli 2012, 2). In total 344,030 filings of this particular submission type alone are reported on EDGAR. Appendix B shows a detailed overview of all publicly available submissions made with the SEC between 1993 and 2016.

The SEC assigns to each filer a Central Index Key (CIK) which is a unique identifier used on the EDGAR database in order to label and identify each individual filer in the system (Gerdes 2003, 16). Since 1993 in total 580,225

unique CIK numbers were assigned and stored in the SEC's electronic disclosure system. The majority of these CIKs were not assigned to publicly traded companies but to private firms, hedge funds and mutual funds as well as to private individuals who receive a CIK when filing with the SEC (Garcia and Norli 2012, 2). Table 3 reports the number of unique CIKs (unique filers) filing a certain submission type with the SEC and its EDGAR system.

Table 3. Statistics on EDGAR filers.

Rank	Form/ Description	Submission Type	Unique CIKs	Mean	Med.	Max.
1	Changes in ownership	4	206,652	28.3	7	12,170
2	Initial ownership report	3	187,366	2.9	1	550
3	Offering of securities without registration	D	104,853	1.4	1	375
4	Regulation D exemption filing (paper submission)	REGDEX	87,285	1.5	1	150
5	Changes in ownership amendments	4/A	62,099	3.2	1	338
6	Annual report on ownership changes	5	47,466	3.9	1	473
7	Change on a prospectus	424B3	45,204	5.6	2	9,911
8	5% active ownership triggers	SC 13D	43,381	2.3	1	730
9	5% passive ownership triggers	SC 13G	41,629	8.3	2	7,726
10	Notification of effectiveness for Securities Act registration statement	EFFECT	40,485	2.4	1	86
11	Registration of securities issued in business combination transactions	S-4	40,139	2.0	1	70
12	Current report filing	8-K	38,713	35.6	10	1,484
13	Offering of securities without registration amendments	D/A	35,673	2.8	2	1,601
14	Registration of securities issued in business combination transactions amendments	S-4/A	35,158	2.8	2	63
15	Annual report	10-K	33,968	4.9	3	107
16	5% passive ownership triggers amendments	SC 13G/A	33,339	17.6	4	25,447
17	SEC-originated letters to filers	UPLOAD	31,720	5.0	3	91
18	Filer response letters	CORRESP	30,031	5.1	3	157
19	5% active ownership triggers amendments	SC 13D/A	29,742	6.8	3	5,528
20	Quarterly report	10-Q	26,271	19.9	14	189

Notes: The table presents the most frequent submission types made on EDGAR in descending order of unique SEC registrants filing a particular submission type. The time period is 1993-2016. The fourth column contains the total number of unique filers submitting a particular form type. Columns 5-7 present the means, medians and maxima of particular filing form types submitted by unique SEC filers.

Submission type Form 4 (Form 3) was submitted by 206,652 (187,366) different filers between 1993 and 2016. Annual reports on Form 10-K were submitted to the SEC by 33,968 filers. Quarterly reports on Form 10-Q can be associated with 26,271 unique filers whereas the number of CIKs assigned to current reports on Form 8-K is 38,713. On average each registrant filed 4.9 annual reports on Form 10-K and 19.9 quarterly reports on Form 10-Q with the EDGAR system in addition to 35.6 current reports on Form 8-K since 1993. AFS SenSub Corp. (CIK 1347185), an issuer of asset-backed securities, filed 107 annual reports on Form 10-K (56 on 10-K/A). PowerShares DB Multi-Sector Commodity Trust (CIK 1367306), an investment company offering several investment funds, filed 189 quarterly reports on Form 10-Q (7 on 10-Q/A). Chase Bank USA, National Association (CIK 869090) filed 1,484 Form 8-K statements (12 on 8-K/A). 730 Schedule 13D Forms were filed by Gamco Investors, INC. (CIK 807249), an investment advisory and brokerage service firm, (5,528 on SC 13D/A) whereas

FMR LLC (CIK 315066), the financial services conglomerate known as Fidelity Investments, filed 7,726 Schedule 13G Forms (25,447 on SC 13G/A).

3. SEC EDGAR Data Gathering

Researchers in the field of finance and accounting often rely on programming languages (Perl, Python, R, SAS, and SPSS) to retrieve financial statements filed with the SEC. The use of a programming language as a tool is problematic for several reasons. First, many people analyzing financial reports are not familiar with these programming languages. For them it is time-consuming to apply a specific and complex coding language to obtain the corporate filings from EDGAR. Secondly, due to downloading only one filing at a time the procedure is very slow especially when obtaining massive data from the database. Thirdly, since incremental changes have to be made to the algorithm to retrieve another filing form type or filings from another company this particular method is very error-prone.

In contrast, widely used internet browsers (e.g. Mozilla-Firefox, Google-Chrome) can be easily equipped with powerful applications (e.g. DownThemAll, GetThemAll) which offer advanced download capabilities. These fully integrated browser extensions are able to identify links contained in a webpage or file and download the desired document parts simultaneously. To feed these applications only a standard MS Excel spreadsheet is necessary.

Every filing made through the EDGAR system in a particular quarter between 1993 and 2016 is stored in an associated index file (file extension *.idx) (Garcia and Norli 2012, 3). The EDGAR index files therefore represent a helpful resource in retrieving massive data from the database. They list important information for each filing such as the name of the filer, the particular central index key, the date and the type of the submission as well as the particular name of the document on the SEC server. In general, four different types of index files are available sorting the filings made on EDGAR by company name, form type, central index key or by submissions containing financial statements formatted in eXtensible Business Reporting Language (XBRL)¹ (SEC 2015). Appendix A lists all index files sorting the entire EDGAR database by filing form type for every quarter between 1993 and 2016 (SEC Index Files 2016). When examining the form index files more precisely one can see that the index files do not only contain the name of any filing made on EDGAR but rather the (entire) server path. Table 4 illustrates an excerpt of information stated in the SEC EDGAR form index file from the first quarter of 2016.

¹ eXtensible Business Reporting Language (XBRL) is a meta language for creating markup languages. XBRL was created for the electronic exchange of business data (Ditter, Henselmann and Scherr 2011, 22). XBRL is a variant of the Extensible Markup Language (XML) and related to the HyperText Markup Language (HTML). XBRL provides semantic context for data reported in SEC EDGAR submissions (Form 10-K) (Bodnaruk, Loughran and McDonald 2015, 643).

Table 4. SEC EDGAR 2016 form index file.

Form	Company Name	CIK	Date Filed	File Name (including partial directory)
10-K	AMAZON COM INC	1018724	2016-01-29	edgar/data/1018724/0001018724-16-000172.txt
10-K	Alphabet Inc.	1652044	2016-02-11	edgar/data/1652044/0001652044-16-000012.txt
10-K	COCA COLA CO	21344	2016-02-25	edgar/data/21344/0000021344-16-000050.txt
...
10-Q	STARBUCKS CORP	829224	2016-01-26	edgar/data/829224/0000829224-16-000049.txt
10-Q	APPLE INC	320193	2016-01-27	edgar/data/320193/0001193125-16-439878.txt
10-Q	VISA INC	1403161	2016-01-28	edgar/data/1403161/0001403161-16-000015.txt
...
8-K	GOLDMAN SACHS INC	886982	2016-01-20	edgar/data/886982/0001193125-16-433035.txt
8-K	EBAY INC	1065088	2016-01-27	edgar/data/1065088/0001065088-16-000262.txt
8-K	MICROSOFT CORP	789019	2016-01-28	edgar/data/789019/0001193125-16-441813.txt
...

Notes: The table presents an excerpt of the EDGAR form index file from the first quarter of 2016. The first column shows the specific submission form type of each filing listed in the quarterly SEC EDGAR form index file. The second column shows the name of each EDGAR registrant submitting a filing. The Central Index Key (CIK) of each filer is shown in the next column. The last two columns contain the submission date and the document name of each filing on the EDGAR database.

By opening the index files for example with a simple MS Excel spreadsheet (file extension *.xlsx) a Uniform Resource Locator (URL) can be created for each financial statement which is listed in a particular index file since the name of the filing and its (partial) server path (directory) is stated. To do so the protocol (https://), the hostname (www.sec.gov/) and a link to the archives directory (Archives/) have to be added to the file name from the index file. Table 5 illustrates the URL components of Coca Cola's 2015 annual report on Form 10-K filed with the SEC on February 25, 2016.

Table 5. Components of a Uniform Resource Locator (URL).

URL Component	String
Protocol	https://
Hostname	www.sec.gov
Directory	/Archives/edgar/data/21344/
File Name	0000021344-16-000050.txt
Composed URL	https://www.sec.gov/Archives/edgar/data/21344/0000021344-16-000050.txt

Notes: The table presents the different components of a Uniform Resource Locator. The first column shows the different component types within the URL. The second column shows the particular value of a certain component type in the URL.

Figure 1 illustrates the composition of URLs in MS Excel using a simple formula command.

Figure 1. Composition of Uniform Resource Locators (URLs) in MS Excel.

A	B	C	D	E	F
Form Type	Company Name	CIK	Date Filed	File Name	Composed Uniform Resource Locator (URL)
10-K	COCA COLA CO	21344	25.02.2016	edgar/data/21344/0000021344-16-000050.txt	https://www.sec.gov/Archives/edgar/data/21344/0000021344-16-000050.txt
10-K	MCDONALDS CORP	63908	25.02.2016	edgar/data/63908/0000063908-16-000103.txt	https://www.sec.gov/Archives/edgar/data/63908/0000063908-16-000103.txt
10-K	Alphabet Inc.	1652044	11.02.2016	edgar/data/1652044/0001652044-16-000012.txt	https://www.sec.gov/Archives/edgar/data/1652044/0001652044-16-000012.txt
10-K	EXXON MOBIL CORP	34088	24.02.2016	edgar/data/34088/0000034088-16-000065.txt	
10-K	AMERICAN EXPRESS CO	4962	19.02.2016	edgar/data/4962/0001193125-16-469798.txt	
10-K	MASTERCARD INC	1141391	12.02.2016	edgar/data/1141391/0001141391-16-000085.txt	
10-K	BlackRock Inc.	1364742	26.02.2016	edgar/data/1364742/0001564590-16-013511.txt	

Notes: The table demonstrates the composition of URLs in MS Excel based on the SEC EDGAR form index files. The last column presents the URL for each disclosure which has been composed using the filename listed in the form index file and a simple MS Excel formula command (=\"https://www.sec.gov/Archives/\"&File Name).

These URLs which have been composed based on the EDGAR index files can be copied into a plain text file (file extension *.txt). By opening it with the browser extensive data (financial statements) can be retrieved from the SEC and its EDGAR system in a fast and efficient way using a browser extension (however, the composed URLs can also be implemented in any other data gathering method).

This method offers various significant advantages. First, for many people composing URLs with commonly used and easy accessible computer software like MS Excel is simpler and faster than relying on complex coding languages to identify and retrieve the documents in question. Secondly, since multiple documents can be retrieved at the same time using browser extensions, the described method is again a lot faster especially when obtaining massive data from EDGAR. Thirdly, by sorting or filtering the different index files in MS Excel the proposed method can easily be adjusted to retrieve another filing form type or data from another company. The result of this procedure is validated through obtaining exactly the same financial statements investors and researchers would retrieve using a complex, slow and error-prone alternative.

4. HyperText Markup Language in SEC filings

Because financial statements filed with the SEC are formatted in HyperText Markup Language (HTML) the fundamentals of HTML are illustrated first, followed by an examination of the data formatted in HTML provided by the SEC and its EDGAR system.

4.1 Fundamentals of HyperText Markup Language

HyperText Markup Language (HTML) is a universally understood digital language which is used to publish and distribute information globally. HTML is the publishing language of the World Wide Web (W3C Recommendation 1999, Section 2.2). HTML is used to create HyperText documents that are portable from one platform to another (Filer Manual 2016, Section 5-12) due to their generic semantics as a Standard Generalized Markup Language (SGML)² application (W3C Recommendation 1999, Section 3.1). HTML enables authors to publish documents online, assign a specific look or layout to document content (tagging) (Ditter, Henselmann and Scherr 2011, 20; Chakraborty and Vasarhelyi 2010, 4) or to retrieve information online via HyperText links (W3C Recommendation 1999, Section 2.2). The World Wide Web Consortium (W3C) is maintaining and specifying the

² Standard Generalized Markup Language (SGML) is a meta language or system for defining markup languages. A markup language defined in SGML is a “*SGML application*” (HTML) (W3C Recommendation 1999, Section 3.1).

vocabulary (applicable markups) and grammar (logical structure) of HTML documents (Ditter, Henselmann and Scherr 2011, 20).

A valid HTML document is composed of three different parts (W3C Recommendation 1999, Section 7.1). First, it declares which version of HTML is used in the document through the document type declaration (<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN" "http://www.w3.org/TR/html4/strict.dtd">). The document type declaration names the document type definition (DTD) specifying which elements and attributes can be implemented into a document formatted in HTML (W3C Recommendation 1999, Section 7.2). HTML 4.01 specifies three different DTDs: HTML 4.01 Strict DTD; HTML 4.01 Transitional DTD and HTML 4.01 Frameset DTD (W3C Recommendation 1999, Section 7.2). The W3C recommends to use HTML 4.01 Strict DTD which excludes presentation attributes since these elements are supposed to be replaced by style sheets (W3C Strict DTD, 1999). The second part of a HTML document is the document head (<HEAD>). This section contains information about the current document such as the title and relevant keywords for search engines. In general, the elements appearing in the head section are not presented by a document formatted in HTML (W3C Recommendation 1999, Section 7.4). The third and most important part of a HTML document is the body (<BODY>). This section contains the actual content of the document such as text paragraphs, images, graphics, tables, links, etc. (W3C Recommendation 1999, Section 7.5). The content in the document body can be structured in many different ways using various HTML elements (tags) to accomplish a certain look or layout to present the embedded information. Figure 2 illustrates a simplified excerpt of a document formatted in HTML (W3C Recommendation 1999, Section 7.1).

Figure 2. Document formatted in HTML 4.01.

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN" "http://www.w3.org/TR/html4/strict.dtd">

<HTML>

  <HEAD>
    <TITLE>"Annual Report on Form 10-K"</TITLE>
  </HEAD>

  <BODY>"United States Securities and Exchange Commission ..."</BODY>

</HTML>
```

Notes: The figure shows a simplified excerpt of a well-formed document formatted in HTML 4.01.

4.2 SEC EDGAR HTML Data

“Official” financial statements filed with the SEC have to be formatted either in American Standard Code for Information Interchange (ASCII) or in HyperText Markup Language (HTML 3.2/4.0).³ Financial statements formatted in Portable Document Format (PDF) or XBRL are considered “unofficial” documents (submissions formatted in PDF and XBRL may qualify as official documents as well when specific criteria are met) (Filer Manual 2016, Section 5-1, Section 5-26). Due to a limited support of HTML in order to reduce the number of inconsistencies caused by HTML 4.0 implementation variances (SEC 2000), the EDGAR system only accepts a subset of HTML 3.2 semantics (tags) and several HTML 4.0 attributes (Filer Manual 2016, Section 5-17) therefore enforcing several restrictions (no active content, no external references etc.) of HTML formatting in financial statement submissions (Filer Manual 2016, Section 5-12).

The “*Complete Submission Text File*” (file extension *.txt) provided by the EDGAR system represents an aggregation of all information in a particular financial statement filed with the SEC. The text version of the filings on the SEC server contains the 10-K document formatted in HTML, XBRL, exhibits and ASCII-encoded graphics (“*binary-to-text*” encoding or “*uuencoding*” converts binary data files to plain ASCII-printable characters to facilitate transfer across various hardware platforms) (Bodnaruk, Loughran and McDonald 2015, 643; Loughran and McDonald 2011b, 1). Besides the “*Complete Submission Text File*” several submission parts (documents) are also provided in HTML (file extension *.htm) such as the core 10-K document and the exhibits which have been submitted (Bodnaruk, Loughran and McDonald 2015, 643). For example, Coca Cola’s 10-K filing on February 25, 2016 lists the core 10-K filing in HTML format, ten exhibits, eight graphic files (file extension *.jpg), six XBRL files and a single “*Complete Submission Text File*” containing all of these documents (SEC EDGAR Archives 2016). Figure 3 illustrates a simplified excerpt of Coca Cola’s 2015 annual report on Form 10-K formatted in HTML embedded in the “*Complete Submission Text File*”.

³ American Standard Code for Information Interchange (ASCII) was created to standardize the storing of alphanumeric characters by defining a unique binary 7-bits number for each of the 128 (2^7) storable characters. ASCII included the Roman/Latin alphabet and essential characters for writing English. Using 8-bits (single 8-bit byte) the character encoding of HTML 2.0/4.01 is ISO-8859-1, an extension to ASCII which allows to present 256 (2^8) characters in HTML documents (W3 Schools, 2016; Palmer 2010, 11).

Figure 3. Coca Cola’s 2015 annual report on Form 10-K formatted in HTML 4.01.

```

<SEC-HEADER>...</SEC-HEADER>

<DOCUMENT>
<TYPE>10-K
<TEXT>
<!DOCTYPE html PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN" "http://www.w3.org/TR/html4/loose.dtd">

<html>
  <head>
    <title>10-K</title>
  </head>

  <body style="font-family:Times New Roman;font-size:10pt;">...

  <div style...><font style...>UNITED STATES</font></div>...
  <div style...><font style...>SECURITIES AND EXCHANGE COMMISSION</font></div>...

  <div style...><font style...>FORM 10-K</font></div>...
  <div style...><font style...>For the fiscal year ended December 31, 2015</font></div>...

  <div style...><font style...>ITEM 1. BUSINESS</font></div>...
  <div style...><font style...>The Coca-Cola Company is the world's largest beverage company...

  <div style...><font style...>ITEM 1A. RISK FACTORS</font></div>...
  <div style...><font style...>In addition to the other information set forth in this report, you should...

  <div style...><font style...>ITEM 7. MANAGEMENT'S DISCUSSION AND ANALYSIS OF FINANCIAL
    CONDITION AND RESULTS OF OPERATIONS</font></div>...
  <div style...><font style...>The following Management's Discussion and Analysis of Financial Condition...

  </body>

</html>

</TEXT>
</DOCUMENT>

```

Notes: The figure presents a simplified excerpt of Coca Cola’s 2015 annual report on Form 10-K formatted in HTML 4.01 contained in the “*Complete Submission Text File*” available on the SEC server.

5. Textual Information in Financial Statements

This section describes how regular expressions are used to extract textual information from financial statements filed with the SEC. First, I illustrate the fundamentals of regular expressions. Then I discuss the algorithm to extract textual information from financial statements using only regular expressions before presenting the actual text embedded in financial statements as a result of the designed algorithm. Due to their high relevance for investors and researchers an actual annual report on Form 10-K from the Coca Cola Company serves as basis for the illustration.

5.1 Fundamentals of Regular Expressions

Regular expressions or regular sets were first used as an algebra by mathematicians to describe models developed by physiologists of how the nervous system would work at the neuron level. The first published computational use of regular expressions was in 1968 by Ken Thompson (Friedl 2006, 85) who describes regular expressions as “a

method for locating specific character strings embedded in character text” (Thompson 1968, 419). They are implemented not only in modern programming languages, but also in application programs that can be used for text analysis without special programming skills (e.g. RapidMiner).

Regular expressions (“*RegEx*”; “*RegExp*”; “*RegExes*”) with a general pattern notation (pattern language) allow to process all kinds of text and data in a flexible and efficient way (Friedl 2006, 1; Loughran and McDonald 2016, 1219-1220). In particular RegExes can be used to modify textual elements or to identify and extract certain information from different documents (Goyvaerts and Levithan 2012, 1-2). The two types (full) regular expressions are composed of special characters (metacharacters) and normal (literal) text characters acting as the grammar and the words of the regular expression language (Friedl 2006, 5; Goyvaerts and Levithan 2012, 28). For example, RegEx: “[0-9]” identifies all digits, RegEx: “[a-zA-Z]” isolates all upper and lower-case letters (character classes) and RegEx: “.” matches all of these elements (metacharacter) embedded in an underlying text document (Friedl 2006, 9-12; Goyvaerts and Levithan 2012, 33-39). Another metacharacter and counting element (quantifier) within the regular expression language is a star or an asterisk (*) which quantifies the immediately preceding item within the defined expression (match any number of the preceding element including none) (Friedl 2006, 18; Goyvaerts and Levithan 2012, 74). Counting elements or quantifiers are used to specify the search pattern of regular expressions in more detail. “*Greedy*” quantifiers like “*” match as much as possible whereas “*lazy*” quantifiers such as “*?” match as little as possible to satisfy the search pattern of a composed regular expression (Friedl 2006, 141; Goyvaerts and Levithan 2012, 75-78).⁴

In addition, regular expressions can be modified in the way they are interpreted and applied using different regular expression modes (modifiers). These modifiers allow to change the search pattern of a particular regular expression (matching mode) in modern programming languages or in application programs. Regular expressions equipped with “*case-insensitive match mode*” ((?i)) ignore the letter case of the input (textual elements) during the matching process allowing the search pattern to match both upper and lower case letters (Friedl 2006, 110-111; Goyvaerts and Levithan 2012, 29). Since modern applications work with multiple (coding) lines regular expressions need to be modified in order to match a string across different lines. “*Dot-matches-all match mode*” also known as “*single-line mode*” ((?s)) modifies the search pattern of a regular expression in a way that it matches a character string across multiple lines (Friedl 2006, 110-113; Goyvaerts and Levithan 2012, 39-40). By designing regular

⁴ “*Nondeterministic Finite Automation (NFA)*” or “*regex-directed*” and “*Deterministic Finite Automation (DFA)*” or “*text-directed*” engines are the two basic kinds of regular expression engines (Friedl 2006, 153-156). Lazy quantifiers can only be implemented in “*regex-directed*” regular expression engines (Goyvaerts 2007, 7). The DFA engine type is used by awk, egrep, flex, lex, MySQL and Procmail. The NFA engine type is used by GNU Emacs, Java, grep, less, more, .Net languages, PCRE library, Perl, PHP, Python, Ruby, sed and vi (Friedl 2006, 145).

expressions and implementing them into modern computer software the results of various search patterns (textual information) can be highlighted and changed or even removed from the underlying text at all (Friedl 2006, 2; Goyvaerts and Levithan 2012, 1).

5.2 Extraction of Textual Information

Researchers in the field of finance and accounting (as well as business data providers) use the “*Complete Submission Text Files*” (file extension *.txt) provided by the SEC and its EDGAR system to extract textual information from financial statements. In order to delete all non-textual elements (HTML tags and their corresponding attributes) most often special text-processing programs and their predefined applications (HTML-Parser) are used. This again is problematic for several reasons. First, using predefined text-processing operators to delete non-textual elements makes one platform-dependent since a specific HTML-Parser can not be (easily) implemented into any other text-processing program in use. Secondly, since the extraction algorithm of the HTML-Parser is complex or not presented at all its extraction results can hardly be validated. Thirdly, because of these drawbacks empirical research results are challenging to replicate for a particular or any other data sample. Regular expressions can in fact overcome these problems in extracting textual information embedded in financial statements filed with the SEC. They offer platform-independent (research) results which can be validated and replicated for any data sample at any given time.

The proposed extraction algorithm (“*Annual Report Algorithm*”) first decomposes the “*Complete Submission Text File*” (file extension *.txt) into its components (RegEx 1). In the end, the entire algorithm is validated through obtaining exactly one core (Form 10-K) document and the number of exhibits which have been embedded in the “*Complete Submission Text File*” for every financial statement in the data sample. Next, the “*Annual Report Algorithm*” identifies all other file types contained in the submission since these additional documents are not either a core document or an exhibit within the text version of the filing (RegEx 2). Table 6 illustrates the regular expressions needed to decompose the “*Complete Submission Text File*” of a financial statement filed with the SEC and to identify the embedded document (file) types.

Table 6. Regular expressions contained in the “*Annual Report Algorithm*”.

ID	Description	Regular Expression
1	Decomposition of “ <i>Complete Submission Text File</i> ”	(?s)<DOCUMENT>.*?</DOCUMENT>
2	Identification of document (file) types	<TYPE>.*

Notes: The table presents the regular expressions contained in the “*Annual Report Algorithm*” for extracting documents and identifying document (file) types.

In addition to the filing components described earlier (10-K section, exhibits, XBRL, graphics), several other document (file) types might be embedded in financial statements such as MS Excel files (file extension *.xlsx), ZIP files (file extension *.zip) and encoded PDF files (file extension *.pdf). By applying additional rules in the “*Annual Report Algorithm*” (RegExes 3-22) these documents are deleted to be able to extract textual information only from the core document and the various exhibits contained in the “*Complete Submission Text File*”. The additional SEC-header is not supposed to be removed separately since it has already been deleted by the algorithm.⁵ Table 7 illustrates the regular expressions applied to delete document (file) types other than the core document and the corresponding exhibits.

Table 7. Regular expressions contained in the “*Annual Report Algorithm*”.

ID	Description	Regular Expression
3	Removal of graphic files	(?s)<TYPE>GRAPHIC.*?</TEXT>
4	Removal of MS Excel files	(?s)<TYPE>EXCEL.*?</TEXT>
5	Removal of PDF files	(?s)<TYPE>PDF.*?</TEXT>
6	Removal of ZIP files	(?s)<TYPE>ZIP.*?</TEXT>
7	Removal of cover letter	(?s)<TYPE>COVER.*?</TEXT>
8	Removal of correspondence between SEC staff and EDGAR participant	(?s)<TYPE>CORRESP.*?</TEXT>
9	Removal of XBRL instance document	(?s)<TYPE>EX-10[01].INS.*?</TEXT>
10	Removal of XBRL instance document	(?s)<TYPE>EX-99.SDR [KL].INS.*?</TEXT>
11	Removal of XBRL taxonomy extension schema document	(?s)<TYPE>EX-10[01].SCH.*?</TEXT>
12	Removal of XBRL taxonomy extension schema document	(?s)<TYPE>EX-99.SDR [KL].SCH.*?</TEXT>
13	Removal of XBRL taxonomy extension linkbase document	(?s)<TYPE>EX-10[01].CAL.*?</TEXT>
14	Removal of XBRL taxonomy extension linkbase document	(?s)<TYPE>EX-99.SDR [KL].CAL.*?</TEXT>
15	Removal of XBRL taxonomy extension definition linkbase document	(?s)<TYPE>EX-10[01].DEF.*?</TEXT>
16	Removal of XBRL taxonomy extension definition linkbase document	(?s)<TYPE>EX-99.SDR [KL].LAB.*?</TEXT>
17	Removal of XBRL taxonomy extension labels linkbase document	(?s)<TYPE>EX-10[01].LAB.*?</TEXT>
18	Removal of XBRL taxonomy extension labels linkbase document	(?s)<TYPE>EX-99.SDR [KL].LAB.*?</TEXT>
19	Removal of XBRL taxonomy extension presentation linkbase document	(?s)<TYPE>EX-10[01].PRE.*?</TEXT>
20	Removal of XBRL taxonomy extension presentation linkbase document	(?s)<TYPE>EX-99.SDR [KL].PRE.*?</TEXT>
21	Removal of XBRL taxonomy extension reference linkbase document	(?s)<TYPE>EX-10[01].REF.*?</TEXT>
22	Removal of XBRL documents	(?s)<TYPE>XML.*?</TEXT>

Notes: The table presents the regular expressions contained in the “*Annual Report Algorithm*” for deleting nonrelevant document (file) types.

Next, the “*Annual Report Algorithm*” deletes all metadata included in the core document and the exhibits (RegExes 23-27). Table 8 illustrates the regular expressions for deleting metadata in SEC EDGAR documents.

⁵ RegEx 1 extracts filing parts declared as documents. The SEC header disclosed at the top of every “*Complete Submission Text File*” is declared as “<SEC-HEADER>”, therefore the head section is not a separate document. The SEC header contains additional information about a particular filer such as the industry classification and the address. The phrases “BEGIN PRIVACY-ENHANCED MESSAGE” and “END PRIVACY-ENHANCED MESSAGE” in early SEC EDGAR complete submission text filings are not deleted separately for the same reason.

Table 8. Regular expressions contained in the “*Annual Report Algorithm*”.

ID	Description	Regular Expression
23	Removal of document type information	<TYPE>.*
24	Removal of sequence information	<SEQUENCE>.*
25	Removal of filename	<FILENAME>.*
26	Removal of description	<DESCRIPTION>.*
27	Removal of head section (including document title)	(?s)<HEAD>.*?</HEAD>

Notes: The table presents the regular expressions contained in the “*Annual Report Algorithm*” for deleting nonrelevant document metadata.

Before deleting all HTML elements and their corresponding attributes (RegEx 29) the algorithm deletes tables since they contain non-textual (quantitative) information (RegEx 28).⁶ Table 9 illustrates the set of regular expressions applied to delete tables and HTML elements embedded in financial statements filed with the SEC.

Table 9. Regular expressions contained in the “*Annual Report Algorithm*”.

ID	Description	Regular Expression
28	Removal of table content	(?s)(?i)<Table.*?</Table>
29	Removal of HTML tags and attributes	(?s)<[<^>]*>

Notes: The table presents the regular expressions contained in the “*Annual Report Algorithm*” for deleting tables and HTML elements.

After extracting the core document and the exhibits as well as deleting all HTML elements, the “*Annual Report Algorithm*” adjusts the content embedded in the body section of each HTML-formatted document in order to extract textual elements from financial statements on the EDGAR database. According to the SEC filer manual the EDGAR system suspends financial statements which contain extended ASCII characters. However, it supports submissions with extended character references. By using ISO-8859-1/Latin-1 decimal character references or entity-names (either technique is allowed within SEC submissions) extended ASCII characters can be embedded in financial statement submissions. These extended character sets within HTML documents included in the “*Complete Submission Text File*” need to be decoded to be able to extract human-readable textual information from financial statements (Filer Manual 2016, Section 5-19). The “*Annual Report Algorithm*” finally decodes all extended character sets (RegExes 30-680) most likely embedded in financial statements filed with the SEC and its EDGAR system formatted in HTML 4.01 (ASCII, ANSI/Windows-1252, ISO-8859-1/Latin-1, mathematical, Greek, symbolic and special characters).⁷

⁶ SEC EDGAR filers might use HTML table tags to structure textual content in electronic submissions.

⁷ Instead of applying RegExes 30-680 investors might want to use a predefined text processing operator since they are less dependent on extraction results that can be replicated for other data samples.

5.3 Extraction Results

By applying the “*Annual Report Algorithm*” investors and researchers are able to extract textual information from financial statements filed with the SEC for thousands of companies in a fully automated process. Based on the “*Complete Submission Text File*” provided by the EDGAR system the algorithm extracts the core (Form 10-K) document and the exhibits which have been embedded in the text version of a company’s financial statement. For example for Coca Cola’s 2015 annual report on Form 10-K filed on February 25, 2016 via EDGAR the algorithm extracts one core document in addition to ten different exhibits. Figure 4 illustrates partial extraction results for the 10-K section of the annual report as well as for two exhibits.

Figure 4. Examples of the extraction result of the “*Annual Report Algorithm*”.

UNITED STATES SECURITIES AND EXCHANGE COMMISSION Washington, D.C. 20549 FORM 10-K For the fiscal year ended December 31, 2015 OR For the transition period from to Commission File No. 001-02217 (Exact name of Registrant as specified in its charter) Registrant's telephone number, including area code: (404) 676-2121 Securities registered pursuant to Section 12(b) of the Act: Securities registered pursuant to Section 12(g) of the Act: None...
Exhibit 23.1 CONSENT OF INDEPENDENT REGISTERED PUBLIC ACCOUNTING FIRM We consent to the incorporation by reference in the registration statements and related prospectuses of The Coca-Cola Company listed below of our reports dated February 25, 2016, with respect to the consolidated financial statements of The Coca-Cola Company and subsidiaries, and the effectiveness of internal control over financial reporting of The Coca-Cola Company and subsidiaries, included in this Annual Report (Form 10-K) for the year ended December 31, 2015. /s/ ERNST & YOUNG LLP Atlanta, Georgia February 25, 2016...
EXHIBIT 31.1 CERTIFICATIONS I, Muhtar Kent, Chairman of the Board of Directors and Chief Executive Officer of The Coca-Cola Company, certify that: 1. I have reviewed this annual report on Form 10-K of The Coca-Cola Company; 2. Based on my knowledge, this report does not contain any untrue statement of a material fact or omit to state a material fact necessary to make the statements made, in light of the circumstances under which such statements were made, not misleading with respect to the period covered by this report...

Notes: The figure presents extraction results from Coca Cola’s 2015 annual report on Form 10-K filed with the SEC. The first part of the figure displays the actual 10-K section embedded in text version of the submission. The second part shows the statement of the auditing firm. The certification of the annual report by the CEO is presented in the last part of the figure.

Besides from textual content of entire documents (10-K section and exhibits) contained in the “*Complete Submission Text File*” investors and researchers might be interested in extracting textual information from particular sections (Items) within the core 10-K section of an annual report (like Item 1A - Risk Factors; Item 3 - Legal Proceedings; Item 7 - Management’s Discussion and Analysis of Financial Condition and Results of Operations etc.). In order to extract textual information from particular 10-K items the “*Annual Report Algorithm*” is modified to the “*Items Algorithm*”. Excluding all exhibits, the modified “*Items Algorithm*” isolates only the 10-K section within the SEC submission. After deleting nonrelevant information and decoding reserved characters

within the document investors and researchers can extract textual information from specific 10-K items.⁸ Table 10 specifies the modified “*Items Algorithm*” applied to extract textual information from particular items of the annual report on Form 10-K filed with the SEC.

Table 10. Regular expressions contained in the “*Items Algorithm*”.

ID	Description	Regular Expression
1.1	Extraction of 10-K section	(?s)<TYPE>10-K.*?</TEXT>
2.1	Removal of document metadata	RegExes 23-28
3.1	Removal of table content	(?s)(?i)<Table.*?</Table>
4.1	Decoding of reserved characters	See RegExes 30-680
5.1	Identification and renaming of item headings (“>°Item”)	(?s)(?i)(?m)> +Item >Item ^Item
6.1	Removal of multiple empty spaces	(?s) +
7.1	Extraction of Item 1.	- Business
7.2	Extraction of Item 1A.	- Risk Factors
7.3	Extraction of Item 1B.	- Unresolved Staff Comments
7.4	Extraction of Item 2.	- Properties
7.5	Extraction of Item 3.	- Legal Proceedings
7.6	Extraction of Item 4.	- Mine Safety Disclosures
7.7	Extraction of Item X.	- Executive Officers of the Company
7.8	Extraction of Item 5.	- Market for Registrant’s Common Equity, Related Stockholder Matters and Issuer Purchases of Equity Securities
7.9	Extraction of Item 6.	- Selected Financial Data
7.10	Extraction of Item 7.	- Management’s Discussion and Analysis of Financial Condition and Results of Operations
7.11	Extraction of Item 7A.	- Quantitative and Qualitative Disclosures About Market Risk
7.12	Extraction of Item 8.	- Financial Statements and Supplementary Data
7.13	Extraction of Item 9.	- Changes in and Disagreements with Accountants on Accounting and Financial Disclosure
7.14	Extraction of Item 9A.	- Controls and Procedures
7.15	Extraction of Item 9B.	- Other Information
7.16	Extraction of Item 10.	- Directors, Executive Officers and Corporate Governance
7.17	Extraction of Item 11.	- Executive Compensation
7.18	Extraction of Item 12.	- Security Ownership of Certain Beneficial Owners and Management and Related Stockholder Matters
7.19	Extraction of Item 13.	- Certain Relationships and Related Transactions, and Director Independence
7.20	Extraction of Item 14.	- Principal Accounting Fees and Services
7.21	Extraction of Item 15.	- Exhibits, Financial Statement Schedules
8.1	Removal of HTML tags and attributes	(?s)<[>]*>

Notes: The table presents the regular expressions contained in the modified “*Items Algorithm*” for extracting particular items from the annual report on Form 10-K. RegExes 1.1-6.1 modify the text version of a financial statement to be able to extract (clear) textual information from particular items. RegExes 7.1-7.21 represent the actual regular expressions designed to extract particular sections from the text version of the annual report.

⁸ HTML-formatting inconsistencies might influence the capabilities of the “*Items Algorithm*” in extracting and separating Form 10-K items from multiple filings. Parsing errors due to formatting inconsistencies across documents (SEC filings) have the potential to produce “extraordinary results” in form of misspecified 10-K sections (Loughran and McDonald 2016, 1218). Despite several SEC EDGAR filers using HTML table tags to structure textual content (section headings and section content) tables are deleted due to a majority of research studies in the field of textual analysis excluding table content in general. By deleting tables just before removing HTML tags (RegEx 8.1) the capability of the “*Items Algorithm*” in extracting certain items from the entire 10-K section can be enhanced. At the same time forfeiting the ability to quantify this capability in extracting different sections in relation to the overall text length (several sections would not appear in the overall text length of the entire 10-K section due to being embedded in table tags).

Using only regular expressions to extract textual information from financial statements investors and researchers can implement the designed extraction algorithms in any modern application and computer program available today. By applying either the “*Annual Report Algorithm*” or the “*Items Algorithm*” entire documents (10-K section and exhibits) or particular items from the core 10-K section can be extracted from the annual SEC submissions in order to be analyzed. More importantly, while compensating for expensive commercial products the algorithms and their extraction results can be validated and replicated for any data sample at any given time. Figure 5 finally illustrates several extraction results of the “*Items Algorithm*” from the annual report on Form 10-K highly relevant to investors and researchers alike.

Figure 5. Examples of the extraction result of the “*Items Algorithm*”.

<p>Item 1A. RISK FACTORS In addition to the other information set forth in this report, you should carefully consider the following factors, which could materially affect our business, financial condition or results of operations in future periods. The risks described below are not the only risks facing our Company. Additional risks not currently known to us or that we currently deem to be immaterial also may materially adversely affect our business, financial condition or results of operations in future periods...</p>
<p>Item 3. LEGAL PROCEEDINGS The Company is involved in various legal proceedings, including the proceedings specifically discussed below. Management believes that the total liabilities to the Company that may arise as a result of currently pending legal proceedings will not have a material adverse effect on the Company taken as a whole. Aqua-Chem Litigation On December 20, 2002, the Company filed a lawsuit (The Coca-Cola Company v. Aqua-Chem, Inc., Civil Action No. 2002CV631-50) in the Superior Court of Fulton County, Georgia...</p>
<p>Item 7. MANAGEMENT'S DISCUSSION AND ANALYSIS OF FINANCIAL CONDITION AND RESULTS OF OPERATIONS Overview The following Management's Discussion and Analysis of Financial Condition and Results of Operations ("MD&A") is intended to help the reader understand The Coca-Cola Company, our operations and our present business environment. MD&A is provided as a supplement to - and should be read in conjunction with - our consolidated financial statements and the accompanying notes thereto contained in "Item 8. Financial Statements and Supplementary Data" of this report. This overview summarizes the MD&A, which includes the following sections...</p>

Notes: The figure presents extraction results from Coca Cola's 2015 annual report on Form 10-K filed with the SEC. The first part of the figure displays Item 1A (Risk Factors) embedded in the overall 10-K section. The last two parts of the figure show Item 3 (Legal Proceedings) and Item 7 (Management's Discussion and Analysis of Financial Condition and Results of Operations) contained in the 10-K section of the “*Complete Submission Text File*”.

6. Validation of Extraction Algorithms

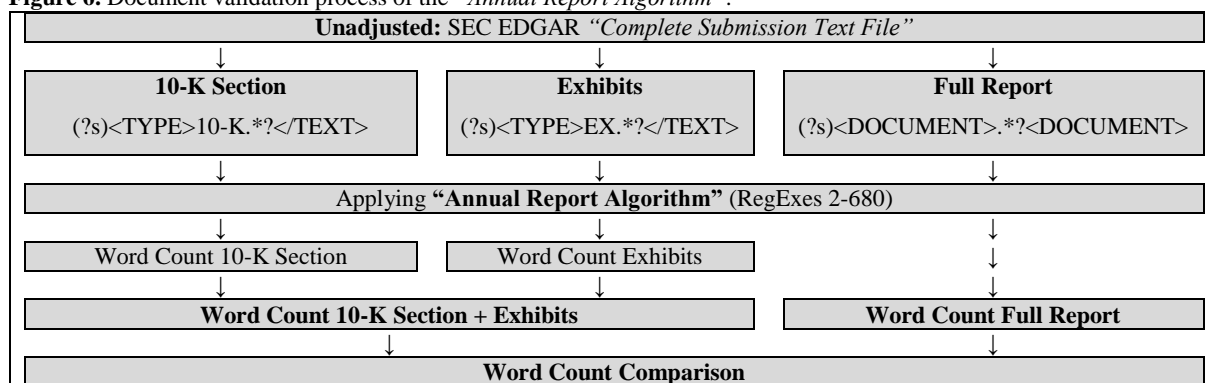
In order to validate the proposed extraction algorithms and to test their capabilities, I retrieve all Form 10-K filings listed in the SEC EDGAR form index files. Using the data gathering method as described in Section 3 in total 188,875 annual reports (167,599 on Form 10-K and 21,276 on Form 10-K405⁹) filed between 1993 and 2016 are retrieved from the EDGAR database (SEC EDGAR Form 10-K types as used in Loughran and McDonald 2011a).¹⁰ The “*Annual Report Algorithm*” is applied to all submissions to derive different word counts for each filing made

⁹ Before eliminated by the SEC in 2003 due to confusion and inconsistency in its application an annual report on Form 10-K405 (405 classification) was used to indicate that a disclosure of delinquent filers pursuant to Item 405 was not included in the current 10-K filing (Loughran and McDonald 2011a, 39).

¹⁰ One out of 188,876 annual reports (Form: 10-K405/ CIK: 884219/ Year: 1996) could not be retrieved (<https://www.sec.gov/Archives/edgar/data/884219/0000884219-96-000033.txt>).

with the SEC. In addition to the overall word count of an annual report, for each core document (10-K section) and the exhibits embedded in a “*Complete Submission Text File*” an individual word count is retrieved in order to be compared (XBRL files declared as exhibits are deleted). Figure 6 illustrates how word counts for each filing and its components are obtained from the “*Complete Submission Text File*” for the document validation process of the “*Annual Report Algorithm*”.

Figure 6. Document validation process of the “*Annual Report Algorithm*”.



Notes: The figure presents the document validation process of the “*Annual Report Algorithm*”. The “*Complete Submission Text File*” of each financial statement as provided on the SEC server is used to extract all relevant components (documents). The “*Annual Report Algorithm*” is applied to each filing in order to retrieve word counts for all relevant documents embedded in the submission. The word count of all relevant documents is compared with the overall length of the submission. A mismatch between the word counts would indicate that the entire report contains nonrelevant document (file) types after applying the “*Annual Report Algorithm*”.

This word count comparison between the overall report on full length and its different components cannot be a validation of the “*Annual Report Algorithm*” since the same algorithm is simply applied to different sets of textual information (10-K section, exhibits, full report). However, if the entire report would still contain document (file) types or elements which are not a part of the core 10-K section or a corresponding exhibit the word count of a certain financial statement would be artificially increased (Word Count Full Report). Table 11 presents the word count comparison of all annual reports and their embedded documents.

Table 11. Results of the document validation process.

Year	Filings	Word Count Comparison				
		10-K Sections + Exhibits	Full Reports	Extraction Error		
				%	Filings	Type
2016	6,467	298,196,698	298,196,698	0.0000	0	-
2015	7,985	350,610,193	350,610,193	0.0000	0	-
2014	8,084	351,653,600	351,664,987	0.0032	1	- Included Form 10-Q
2013	8,105	355,676,863	355,676,863	0.0000	0	-
2012	8,393	347,050,833	347,082,370	0.0091	2	- Declaration error - Included Form 10-Q
2011	8,840	363,208,197	363,208,197	0.0000	0	-
2010	9,165	371,951,148	371,951,148	0.0000	0	-
2009	9,839	397,555,819	397,555,819	0.0000	0	-
2008	8,746	342,695,246	342,695,246	0.0000	0	-
2007	8,574	340,907,681	340,907,681	0.0000	0	-
2006	8,852	326,730,970	326,730,970	0.0000	0	-
2005	9,017	326,083,277	326,111,171	0.0086	1	- Declaration error
2004	8,567	330,972,848	330,972,848	0.0000	0	-
2003	8,468	331,839,061	331,889,877	0.0153	2	- Declaration errors
2002	8,927	332,483,818	332,574,108	0.0272	14	- Declaration errors
2001	9,248	325,089,007	325,099,248	0.0032	1	- Declaration error
2000	9,869	335,180,436	335,242,190	0.0184	5	- Declaration errors - Syntax error
1999	10,122	340,328,009	340,444,185	0.0341	10	- Declaration errors - Syntax errors
1998	10,287	363,400,331	363,480,640	0.0221	8	- Declaration errors - Syntax errors
1997	9,899	319,367,142	319,428,271	0.0191	5	- Declaration errors - Syntax errors
1996	6,258	181,885,167	181,911,885	0.0147	6	- Declaration errors - Syntax errors
1995	3,236	112,604,017	112,622,922	0.0168	3	- Declaration errors - Syntax error
1994	1,923	76,399,982	76,399,982	0.0000	0	-
1993	4	82,283	82,283	0.0000	0	-
Total	188,875	7,221,952,626	7,222,539,782	0.0081	58	-

Notes: The table presents the document validation results of the “*Annual Report Algorithm*”. The third column displays the total word count of all reports in a particular year if only relevant sections have been included. The next column shows the actual word count of all reports in a particular year retrieved from the “*Complete Submission Text Files*” by applying the “*Annual Report Algorithm*”. Mismatches between the word counts for every year are shown in the next column indicating that the algorithm is not able to delete all nonrelevant document (file) types contained in the submissions.

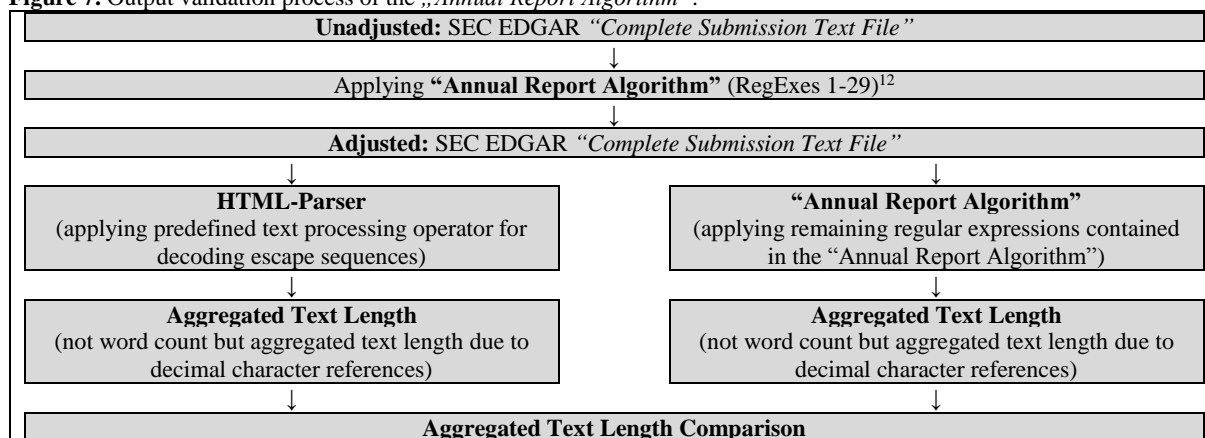
In fact, the ability to validate the entire extraction procedure by applying an alternative to the “*Annual Report Algorithm*” (e.g. HTML-Parser) is limited since to a certain extent the same regular expressions have to be used to create the input for both extraction methods in the first place (extracting core 10-K document and exhibits, deleting nonrelevant document (file) types etc.). Due to this disability in validating the entire extraction process from the beginning by applying an HTML-Parser one has to validate the input the proposed algorithm is creating and its extraction results separately, therefore validating the entire information extraction process. The validation of the textual input created by the “*Annual Report Algorithm*” is represented by the extraction algorithm itself since it uses only regular expressions combined with the electronic filing requirements introduced by the SEC

(precisely not the SEC but Attain, LLC)¹¹. According to the SEC, all documents embedded in a “*Complete Submission Text File*” must be equipped with a “<TYPE>” tag representing the conformed document type of that particular submission part within the text version of the filing (<TYPE>10-K, <TYPE>10-Q, <TYPE>8-K, <TYPE>EX-1, <TYPE>EX-2 etc.) (SEC EDGAR 2015, 12-13, 29). The “*Annual Report Algorithm*” (RegExes 1-29) uses these requirements in order to extract the core document and the corresponding exhibits from annual reports while deleting all documents associated with XBRL and other document (file) types. The search patterns of the “*Annual Report Algorithm*” which have been designed accordingly to the filing requirements of the SEC can be validated due to the general pattern notation of the regular expression language.

An output comparison between the “*Annual Report Algorithm*” and a common HTML-Parser shall serve as an additional validation for the remaining extraction procedure. Therefore, I modify the “*Complete Submission Text Files*” as provided by the SEC (unadjusted filings) and apply the first part of the “*Annual Report Algorithm*” (RegExes 1-29) in order to make the text version of the financial statements readable for the predefined HTML-Parser (adjusted filings). Since this part of the overall validation process focuses on how well the “*Annual Report Algorithm*” is capable of decoding escape sequences embedded in a “*Complete Submission Text File*” the aggregated text length of both procedures are compared rather than the word counts due to decimal character encodings (a simple word count comparison would not fully capture the disability of the “*Annual Report Algorithm*” in decoding these character references in relation to the HTML-Parser). Figure 7 illustrates the output validation process of the “*Annual Report Algorithm*”.

¹¹ The EDGAR Public Dissemination Service (PDS) System is a privatized PDS System managed by Attain, LLC. The system is the primary source to receive all accepted and valid EDGAR filings. The system became operational on November 1, 1998 (SEC EDGAR 2016).

Figure 7. Output validation process of the „Annual Report Algorithm”.



Notes: The figure presents the output validation process of the “Annual Report Algorithm”. The “Complete Submission Text File” of each financial statement as provided on the SEC server is adjusted in order to compare the output of the algorithm with the output a common HTML-Parser would produce. RegExes 1-29 modify the unadjusted document as provided on the EDGAR database before applying a predefined text processing operator (HTML-Parser). The aggregated text length for all filings of both procedures is compared in order to validate the capability of the „Annual Report Algorithm” in decoding escape sequences. The aggregated text length includes each individual element in an underlying text document (text, digits, spaces, special characters etc.).

Table 12 presents the validation results for the “Annual Report Algorithm”.

¹² HTML tags and their corresponding attributes are removed (RegEx 29) before applying the HTML-Parser in order to decode HTML escape sequences since unescaping might interfere HTML structure.

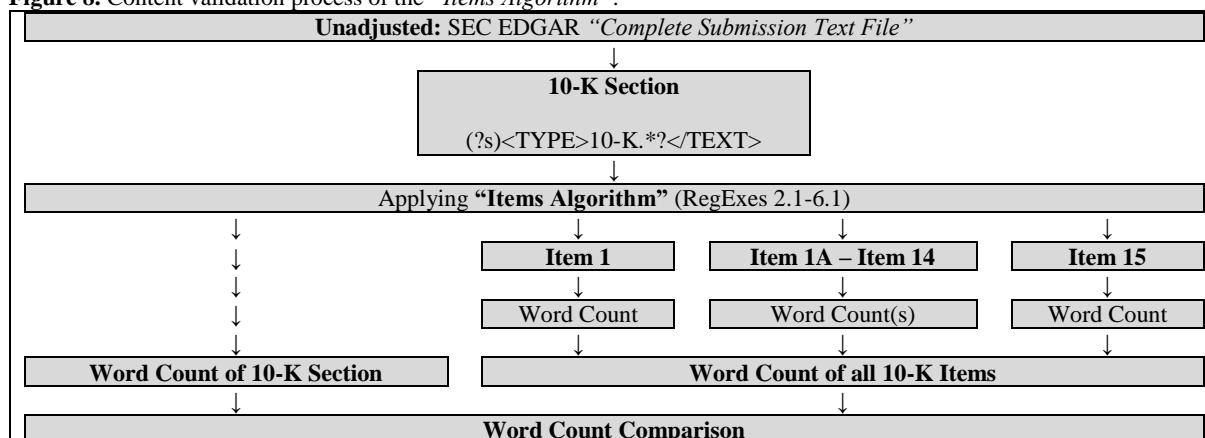
Table 12. Validation results of the “*Annual Report Algorithm*”.

Year	Filings	“Annual Report Algorithm”					
		No Parsing	HTML-Parser	Annual Report Algorithm	Error (%)		
					Mean	Med.	Max.
2016	6,467	101.84	100.00	100.08	0.08	0.01	5.93
2015	7,985	102.04	100.00	100.06	0.06	0.01	5.42
2014	8,084	102.14	100.00	100.05	0.05	0.01	3.54
2013	8,105	102.14	100.00	100.04	0.05	0.00	6.10
2012	8,393	102.40	100.00	100.05	0.05	0.00	9.62
2011	8,840	102.63	100.00	100.05	0.05	0.01	19.07
2010	9,165	102.77	100.00	100.04	0.05	0.01	8.80
2009	9,839	102.92	100.00	100.04	0.04	0.00	2.89
2008	8,746	102.95	100.00	100.05	0.05	0.00	5.39
2007	8,574	102.41	100.00	100.04	0.05	0.00	2.72
2006	8,852	102.59	100.00	100.04	0.05	0.00	3.02
2005	9,017	102.80	100.00	100.04	0.04	0.00	2.87
2004	8,567	102.49	100.00	100.05	0.05	0.00	40.63
2003	8,468	102.28	100.00	100.03	0.04	0.00	2.23
2002	8,927	101.46	100.00	100.02	0.02	0.00	5.30
2001	9,248	100.53	100.00	100.01	0.01	0.00	4.34
2000	9,869	100.17	100.00	100.01	0.01	0.00	4.67
1999	10,122	100.01	100.00	100.00	0.00	0.00	1.08
1998	10,287	100.00	100.00	100.00	0.00	0.00	0.03
1997	9,899	100.00	100.00	100.00	0.00	0.00	0.16
1996	6,258	100.00	100.00	100.00	0.00	0.00	0.90
1995	3,236	100.00	100.00	100.00	0.00	0.00	0.00
1994	1,923	100.00	100.00	100.00	0.00	0.00	0.00
1993	4	100.00	100.00	100.00	0.00	0.00	0.00
Total	188,875	101.95	100.00	100.04	0.03	0.00	40.63

Notes: The table presents the validation results of the “*Annual Report Algorithm*”. The third column presents the aggregated text length including all elements in a filing (text, digits, spaces, special characters etc.) as provided on the SEC server for all submissions in a particular year. The next column represents the aggregated text length of all submissions in a certain year if an HTML-Parser is used to decode character references within the submissions. The fifth column shows the aggregated text length of all filings when using the “*Annual Report Algorithm*” instead of a predefined text processing operator (HTML-Parser) to decode the embedded character references. The inaccuracy in decoding escape sequences when using the “*Annual Report Algorithm*” rather than applying an HTML-Parser is shown in the next column.

In contrast to the “*Annual Report Algorithm*” the modified “*Items Algorithm*” is validated by its ability to distribute the extracted information to the individual items an annual report filed with the SEC is composed of. In order to test and validate the capabilities of the “*Items Algorithm*” I again use the “*Complete Submission Text Files*” as provided by the SEC and extract only the 10-K section of each filing. For each submission, I retrieve separate word counts for the 10-K section and for all individual items extracted by the “*Items Algorithm*”. Despite textual information embedded in the 10-K section not contained in a particular item (introduction) a word count comparison between the overall 10-K section and all items represents an attempt to validate the capabilities of the “*Items Algorithm*” in extracting certain sections from the core document of an annual report filed with the SEC and its EDGAR system. Figure 8 illustrates the content validation process of the “*Items Algorithm*”.

Figure 8. Content validation process of the “Items Algorithm”.



Notes: The figure presents the content validation process of the “Items Algorithm”. First, the entire 10-K section of each filing from the “Complete Submission Text File” as provided on the SEC server is extracted. Word counts for the entire 10-K section as well as for all individual items are retrieved by applying the “Items Algorithm” in order to be compared. Due to structural changes of the annual report on Form 10-K over time (different number of items) the relation of text length between the overall 10-K section and all individual items shall represent the ability of the algorithm to extract particular items from the 10-K section.

Table 13 presents the validation results for the “Items Algorithm”.

Table 13. Validation results of the “*Items Algorithm*”.

Year	Filings		“Items Algorithm”								
			Word Count Comparison		Precision, Recall, and F-measure						
	Number	%	Σ of Items (%)	Rest/ Error (%)	Filings	Items					
					Tested	Exists	Extracted	Correct	Precision	Recall	F-measure
2016	2,886	44.63	97.72	2.28	10	195	191	185	96.86	94.87	95.85
2015	3,714	46.51	97.60	2.40	10	200	200	195	97.50	97.50	97.50
2014	4,004	49.53	97.52	2.48	10	198	197	193	97.97	97.47	97.72
2013	3,962	48.88	97.53	2.47	10	192	188	188	100.00	97.92	98.95
2012	3,938	46.92	97.39	2.61	10	198	192	183	95.31	92.42	93.85
2011	4,104	46.43	97.43	2.57	10	193	191	189	98.95	97.93	98.44
2010	2,805	30.61	97.10	2.90	10	197	168	155	92.26	78.68	84.93
2009	2,719	27.63	97.15	2.85	10	196	181	164	90.61	83.67	87.00
2008	2,077	23.75	97.28	2.72	10	196	184	170	92.39	86.73	89.47
2007	2,065	24.08	97.38	2.62	10	195	184	173	94.02	88.72	91.29
2006	2,662	30.07	97.49	2.51	10	198	184	165	89.67	83.33	86.39
2005	3,122	34.62	97.70	2.30	10	181	175	163	93.14	90.06	91.57
2004	3,496	40.81	97.60	2.40	10	173	170	163	95.88	94.22	95.04
2003	3,903	46.09	97.33	2.67	10	161	161	154	95.65	95.65	95.65
2002	4,961	55.57	97.65	2.35	10	150	150	150	100.00	100.00	100.00
2001	5,799	62.71	97.61	2.39	10	146	144	135	93.75	92.47	93.10
2000	6,268	63.51	97.55	2.45	10	150	149	138	92.62	92.00	92.31
1999	6,302	62.26	97.55	2.45	10	146	145	143	98.62	97.95	98.28
1998	6,492	63.11	97.56	2.44	10	140	140	128	91.43	91.43	91.43
1997	6,397	64.62	97.43	2.57	10	132	129	125	96.90	94.70	95.79
1996	3,918	62.61	97.27	2.73	10	136	121	112	92.56	82.35	87.16
1995	1,907	58.93	97.07	2.93	10	135	135	127	94.07	94.07	94.07
1994	1,039	54.03	97.23	2.77	10	140	139	133	95.68	95.00	95.34
1993	1	25.00	98.49	1.51	1	14	14	14	100.00	100.00	100.00
Total	88,541	46.88	97.48	2.52	231	3,962	3,832	3,645	95.12	92.00	93.53

Notes: The table presents the validation results of the “*Items Algorithm*”. The second and third columns show the number of filings of which items could be extracted from by applying the “*Items Algorithm*” (filings were not machine-parsable due to lacks of content, inconsistent filing structure, table tags and HTML formatting inconsistencies). Only filings with extracted items length exceeding 90 percent of 10-K section are presented. The next two columns show the average amount of extracted information from each filing in a particular year since 1993. The next columns show the performance evaluation of the “*Items Algorithm*” using precision (=number of correct answers/number of total answers), recall (=number of correct answers/total possible correct answers), and F-measure (=2*precision*recall/precision+recall).

7. Descriptive Statistics on Form 10-K contents

In total, I examine the textual composition of 188,875 annual reports filed with the SEC between 1993 and 2016.

On average, an annual report on Form 10-K submitted to the EDGAR system during the sample period is composed of 38,240 words. The average word count of an annual submission increased from 39,730 in 1994 to 46,111 in 2016. The medians of the word counts increased accordingly. The majority of textual information embedded in an annual report on Form 10-K are contained in the core document (64.95 percent) whereas the disclosed exhibits represent only a minority of the overall textual elements stated in annual submissions (35.04 percent). By examining the EDGAR database and its Form 10-K filings in more detail, investors and researchers can see that the average file size (Megabyte) of an annual report made with the electronic disclosure system increased in recent years due to HTML formatting, ASCII-encodings and XBRL documents. Table 14 presents descriptive statistics of the text length and the file size of 188,875 annual reports on Form 10-K (Form 10-K405) filed with the SEC between 1993 and 2016.

Table 14. Descriptive statistics of SEC EDGAR Form 10-K reports.

Year	Filings	Word Count					File Size		
		Full Report (Number)			10-K Sections (%)	Exhibits (%)	Mean (MB)	Med. (MB)	Max. (MB)
		Mean	Med.	Max.					
2016	6,467	46,111	39,997	1,112,167	79.54	20.46	12.50	9.11	261.90
2015	7,985	43,909	37,262	1,657,009	79.51	20.49	15.12	10.18	414.52
2014	8,084	43,501	35,840	2,884,474	78.38	21.62	14.08	9.72	402.86
2013	8,105	43,884	35,181	6,257,121	77.32	22.68	13.22	9.38	254.18
2012	8,393	41,354	34,135	1,441,676	78.62	21.37	8.68	4.90	139.48
2011	8,840	41,087	33,008	1,031,964	77.33	22.67	4.48	1.71	212.57
2010	9,165	40,584	32,448	957,870	77.65	22.35	2.50	1.49	95.27
2009	9,839	40,406	32,074	3,997,528	74.97	25.03	1.90	1.33	86.21
2008	8,746	39,183	32,501	779,558	72.72	27.28	1.72	1.27	61.97
2007	8,574	39,761	32,206	2,617,579	73.67	26.33	1.81	1.28	91.99
2006	8,852	36,910	30,247	908,916	70.76	29.24	1.42	1.01	61.16
2005	9,017	36,166	28,854	1,442,810	66.13	33.86	1.19	0.82	80.62
2004	8,567	38,633	28,655	1,008,146	60.55	39.45	0.98	0.67	27.82
2003	8,468	39,193	28,738	911,982	58.15	41.83	0.90	0.55	24.01
2002	8,927	37,255	26,201	1,545,636	52.82	47.15	0.59	0.34	26.59
2001	9,248	35,153	24,531	1,308,749	52.03	47.97	0.40	0.28	23.34
2000	9,869	33,969	23,619	1,258,064	51.01	48.97	0.35	0.26	19.91
1999	10,122	33,634	23,290	496,458	49.40	50.56	0.33	0.25	8.29
1998	10,287	35,334	22,206	667,721	44.27	55.71	0.33	0.24	4.82
1997	9,899	32,269	20,496	650,347	44.84	55.14	0.30	0.22	4.82
1996	6,258	29,069	19,082	447,469	45.68	54.31	0.28	0.21	4.25
1995	3,236	34,803	22,570	361,832	38.51	61.48	0.34	0.24	4.03
1994	1,923	39,730	25,510	553,782	37.55	62.45	0.39	0.28	4.27
1993	4	20,571	18,247	31,993	83.01	16.99	0.23	0.26	0.27
Total	188,875	38,240	28,772	6,257,121	64.95	35.04	3.56	0.71	414.52

Notes: The table presents descriptive statistics of the text lengths, document compositions and file sizes for all annual reports filed with the SEC since 1993. Columns 3-5 show the means, medians and maxima of word counts of Form 10-K filings made on EDGAR. The average distribution of textual information between the 10-K sections and exhibits contained in the “*Complete Submission Text Files*” is presented in column 6 and 7. The last three columns of the table present the means, medians and maxima of the file sizes (Megabyte) for all Form 10-K submissions.

The distribution of textual elements among the various 10-K items is unequal. On average 22.65 percent of all textual information are contained in Item 1 (“*Business*”). Describing a company’s business as well as its main products and services, the item may also include information about the competition, regulations and other issues a particular company is faced with (SEC 2011; SEC Regulation S-K 2016, Section 229.101 Item 101). Item 7 (“*Management’s Discussion and Analysis of Financial Condition and Results of Operations – MD&A*”) represents 18.58 percent of the given information within Form 10-K filings made with the SEC. The item states information about a company’s operations and financial results in addition to its liquidity and capital resources. The section may include off-balance sheet arrangements and contractual obligations alongside key business risks (SEC 2011; SEC Regulation S-K 2016, Section 229.303 Item 303). Item 8 (“*Financial Statements and Supplementary Data*”) requires a company to disclose audited financial statements (SEC 2011; SEC Regulation S-X 2016, Section 210.3.01-210.3.02; SEC Regulation S-K 2016, Section 229.302 Item 302). Additional information explaining the financial statements in more detail (“*Notes to Consolidated Financial Statements*”, “*Report of Management*”, “*Report of Independent Registered Accounting Firm*” etc.) represent 15.96 percent of all given information in the

10-K section of an annual report. Item 1A (“*Risk Factors*”) describes significant factors that may adversely affect a filer’s business, financial condition or future financial performance (SEC 2011; SEC Regulation S-K 2016, Section 229.503 Item 503(c)). Since electronic filings became available on average 8.42 percent of all textual information disclosed in annual submissions are contained in this section. Each of the remaining items only represent a fraction of the overall textual information embedded in Form 10-K filings. While the length for most sections in annual reports remained constant over time the amount of textual information contained in Item 1A (“*Risk Factors*”) increased from 12.56 percent in 2006 to 20.10 percent in 2016 indicating that SEC EDGAR filers disclose more information about risks in recent years. Table 15 reports descriptive statistics of the distribution of textual information in Form 10-K reports.

Table 15. Distribution of textual information in SEC EDGAR Form 10-K reports.

Panel A													
Year	Filings	10-K Sections	Item										
			1	1A	1B	2	3	4	X	5	6	7	7A
2016	2,886	100.00	16.71	20.10	0.03	0.74	1.08	0.36	0.00	1.45	1.35	17.82	1.65
2015	3,714	100.00	16.69	18.81	0.04	0.88	0.78	0.40	0.00	1.50	1.42	17.79	1.63
2014	4,004	100.00	17.22	17.40	0.08	0.89	0.92	0.40	0.00	1.53	1.47	17.96	1.60
2013	3,962	100.00	17.28	16.46	0.07	0.95	1.08	0.57	0.00	1.64	1.52	17.87	1.37
2012	3,938	100.00	17.58	15.45	0.07	1.00	0.98	0.53	0.00	1.67	1.53	17.88	1.66
2011	4,104	100.00	17.21	14.54	0.08	1.15	2.10	0.54	0.00	1.67	1.60	17.65	1.34
2010	2,805	100.00	17.53	14.35	0.11	1.03	1.08	0.87	0.00	1.70	1.69	17.82	1.33
2009	2,719	100.00	17.51	14.08	0.03	1.03	1.09	0.82	0.00	1.65	1.81	17.93	1.36
2008	2,077	100.00	18.16	13.95	0.06	1.03	1.31	1.14	0.00	1.27	1.91	18.45	1.56
2007	2,065	100.00	17.69	12.91	0.05	1.11	1.37	2.11	0.00	1.12	2.40	18.46	1.78
2006	2,662	100.00	19.05	12.56	0.22	1.09	1.33	1.26	0.01	1.03	2.04	19.76	1.71
2005	3,122	100.00	24.08	0.30	0.00	1.22	1.59	0.94	0.00	1.25	2.00	24.14	1.93
2004	3,496	100.00	24.33	0.03	0.01	1.50	1.73	1.36	0.00	1.26	2.14	23.48	1.95
2003	3,903	100.00	25.22	0.07	0.00	1.25	1.69	0.83	0.01	1.26	2.02	21.65	1.83
2002	4,961	100.00	26.97	0.07	0.00	1.37	1.62	0.73	0.01	1.34	1.40	20.93	1.96
2001	5,799	100.00	29.78	0.06	0.00	1.41	1.47	0.73	0.01	1.38	1.14	18.74	1.69
2000	6,268	100.00	29.91	0.04	0.00	1.64	1.48	1.02	0.01	1.47	1.10	17.94	1.38
1999	6,302	100.00	30.16	0.04	0.00	1.72	1.51	0.80	0.01	1.40	1.10	19.23	1.24
1998	6,492	100.00	31.27	0.03	0.00	1.80	1.59	0.99	0.01	1.49	1.13	16.87	0.37
1997	6,397	100.00	31.58	0.04	0.00	2.09	1.64	1.00	0.01	1.44	1.10	15.90	0.03
1996	3,918	100.00	29.80	0.03	0.01	2.27	2.24	1.06	0.01	1.45	0.89	15.31	0.00
1995	1,907	100.00	28.77	0.01	0.00	2.23	2.75	1.39	0.01	1.20	0.80	14.07	0.00
1994	1,039	100.00	27.78	0.02	0.00	2.26	3.97	1.54	0.01	1.10	0.72	12.61	0.00
1993	1	100.00	18.45	0.00	0.00	0.92	0.32	1.47	0.00	0.58	0.06	12.96	0.00
Total	88,541	100.00	22.65	8.42	0.04	1.29	1.42	0.83	0.01	1.44	1.51	18.58	1.39
Panel B													
Year	Filings	10-K Sections	Item										
			8	9	9A	9B	10	11	12	13	14	15	Rest/ Error
2016	2,886	100.00	18.58	0.12	1.65	0.26	1.80	1.25	0.38	0.62	0.58	11.18	2.28
2015	3,714	100.00	18.38	0.15	1.80	0.26	2.20	1.58	0.50	0.75	0.65	11.41	2.40
2014	4,004	100.00	17.74	0.13	1.89	0.26	2.40	1.78	0.52	0.79	0.72	11.80	2.48
2013	3,962	100.00	17.26	0.15	2.00	0.25	2.59	1.94	0.57	0.86	0.80	12.29	2.47
2012	3,938	100.00	17.29	0.15	2.09	0.25	2.70	2.04	0.65	0.87	0.80	12.20	2.61
2011	4,104	100.00	16.35	0.19	2.16	0.27	2.89	2.70	0.68	0.95	0.95	12.42	2.57
2010	2,805	100.00	17.21	0.24	2.19	0.31	2.63	1.99	0.66	0.93	1.11	12.31	2.90
2009	2,719	100.00	17.21	0.22	2.31	0.24	2.29	1.79	0.72	0.79	0.91	13.36	2.85
2008	2,077	100.00	18.51	0.18	2.19	0.29	1.71	1.88	0.52	0.62	0.73	11.82	2.72
2007	2,065	100.00	17.33	0.15	1.83	0.34	1.44	2.33	0.50	0.53	0.79	13.13	2.62
2006	2,662	100.00	15.69	0.16	1.97	0.37	1.59	1.25	0.50	0.61	0.68	14.62	2.51
2005	3,122	100.00	17.18	0.24	1.79	0.41	1.73	1.31	0.58	0.64	0.85	15.53	2.30
2004	3,496	100.00	17.37	0.21	1.00	0.01	1.89	1.29	0.67	0.71	1.09	15.57	2.40
2003	3,903	100.00	16.46	0.79	0.11	0.00	1.70	1.43	0.74	0.80	1.80	17.66	2.67
2002	4,961	100.00	14.66	0.61	0.00	0.00	1.79	1.67	0.73	0.98	19.55	1.23	2.35
2001	5,799	100.00	14.17	0.53	0.00	0.00	2.03	1.78	0.81	1.02	20.85	0.02	2.39
2000	6,268	100.00	13.13	0.69	0.00	0.00	2.27	2.00	0.89	1.15	21.40	0.02	2.45
1999	6,302	100.00	12.50	0.72	0.00	0.00	2.33	1.98	0.88	1.12	20.77	0.04	2.45
1998	6,492	100.00	12.57	0.85	0.00	0.00	2.33	2.01	0.88	1.25	22.09	0.02	2.44
1997	6,397	100.00	12.71	0.99	0.00	0.00	2.54	1.94	0.87	1.18	22.35	0.02	2.57
1996	3,918	100.00	13.08	0.96	0.00	0.00	2.72	2.00	0.96	1.16	23.26	0.07	2.73
1995	1,907	100.00	14.87	0.89	0.00	0.00	2.36	1.67	0.74	0.81	24.46	0.03	2.93
1994	1,039	100.00	15.53	1.31	0.00	0.00	1.91	1.84	0.76	0.69	25.17	0.00	2.77
1993	1	100.00	0.09	0.09	0.00	0.00	0.29	0.13	0.27	0.16	62.70	0.00	1.51
Total	88,541	100.00	15.96	0.40	1.16	0.16	2.20	1.81	0.67	0.88	8.12	8.53	2.52

Notes: The table reports the average distribution of textual information among different 10-K items in annual reports filed with the SEC and its EDGAR database since 1993 (only filings with extracted items length exceeding 90 percent of 10-K section are presented).

8. Summary

This paper displays the huge amount and variety of publicly available corporate information filed with the SEC and distributed by its EDGAR database. It shows how massive data can be retrieved from the SEC server in a fast and efficient way using simple and easy accessible software. The second main purpose of this paper is to create standardized procedures (*“Annual Report Algorithm”* and *“Items Algorithm”*) investors and researchers can use to extract any kind of textual information from financial statements filed with the SEC. This is achieved by providing regular expressions for multiple steps of data cleaning and filtering. Using these dynamic and platform-independent extraction algorithms the paper analyses the textual composition of more than 180,000 annual reports filed with the SEC via the EDGAR system between 1993 and 2016. The algorithms are tested for validity in several ways. The tools and algorithms intend to reduce costs and lower technical boundaries for researchers in the field of finance and accounting to engage in textual analysis.

References

- Bodnaruk A, Loughran T, McDonald B (2015) Using 10-K Text to Gauge Financial Constraints, in: Journal of Financial and Quantitative Analysis (JFQA), 50(4)/2015, 623-646
- Bovee M, Kogan A, Nelson K, Srivastava R P, Vasarhelyi M A, (2005) Financial Reporting and Auditing Agent with Net Knowledge (FRAANK) and eXtensible Business Reporting Language (XBRL), in: Journal of Information Systems (JIS), 19(1)/2005, 19-41
- Chakraborty V, Vasarhelyi M A (2010) Automating the process of taxonomy creation and comparison of taxonomy structures, in: 19th Annual Research Workshop on Strategic and Emerging Technologies, American Accounting Association. San Francisco, California, USA
- Cong Y, Kogan A, Vasarhelyi M A (2007) Extraction of Structure and Content from the Edgar Database: A Template-Based Approach, in: Journal of Emerging Technologies in Accounting (JETA) 4(1)/2007, 69-86
- Davis A K, Tama-Sweet I (2012) Managers' Use of Language Across Alternative Disclosure Outlets: Earnings Press Releases versus MD&A, in: Contemporary Accounting Research (CAR), 29(3)/2012, 804-837
- Ditter D, Henselmann K, Scherr E (2011) Using XBRL Technology to Extract Competitive Information from Financial Statements, in: Journal of Intelligence Studies in Business (JISIB), 1/2011, 19-28
- Engelberg J, Sankaraguruswamy S (2007) How to Gather Data Using a Web Crawler: An Application Using SAS to Search Edgar. Working Paper, SSRN
- Filer Manual (2016) Filer Manual – Volume II EDGAR Filing, Available online on URL: <https://www.sec.gov/info/edgar/edgarfm-vol2-v37.pdf>
- Friedl J E F (2006) Mastering Regular Expressions, Third Edition, O'Reilly Media, Inc., Sebastopol, California, USA
- Gaizauskas R, Humphreys K, Azzam S, Wilks Y (1997) Concepticons vs. Lexicons: an Architecture for Multilingual Information Extraction, in: M T Pazienza, Ed. Information Extraction: A Multidisciplinary Approach to an Emerging Information Technology. Springer-Verlag, Berlin Heidelberg, Germany
- Garcia D, Norli O (2012) Crawling EDGAR, in: The Spanish Review of Financial Economics (SRFE), 10/2012, 1-10

Gerdes J Jr (2003) EDGAR-Analyzer: automating the analysis of corporate data contained in the SECs EDGAR database, in: Decision Support Systems 35/2003, 7-29

Goyvaerts J (2007) Regular Expressions: The Complete Tutorial, Available online on URL: <https://www.princeton.edu/~mlovett/reference/Regular-Expressions.pdf>

Goyvaerts J, Levithan S (2012) Regular Expressions Cookbook, Second Edition, O'Reilly Media, Inc., Sebastopol, California, USA

Grant G H, Conlon S J (2006) EDGAR Extraction System: An Automated Approach to Analyze Employee Stock Option Disclosures, in: Journal of Information Systems (JIS), 20(2)/2006, 119-142

Hernandez M A, Ho H, Koutrika G, Krishnamurthy R, Popa L, Stanoi I R, Vaithyanathan S, Das S (2010) Unleashing the Power of Public Data for Financial Risk Measurement, Regulation, and Governance. IBM Technical Report #RJ10475

Jegadeesh N, Wu D (2013) Word power: A new approach for content analysis, in: Journal of Financial Economics (JFE), 110(3)/2013, 712-729

Kambil A, Ginsburg M (1998) Public Access Web Information Systems: Lessons from the Internet EDGAR Project, in: Communications of the ACM (CACM), 41(7)/1998, 91-97

Loughran T, McDonald B (2011a) When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks, in: The Journal of Finance (JoF), 66(1)/2011, 35-65

Loughran T, McDonald B (2011b) Internet Appendix for “When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks”, Available online on URL: <http://www.afajof.org/SpringboardWebApp/userfiles/afa/file/Supplements%20and%20Data%20Sets/Internet%20Appendix%20for%20When%20Is%20a%20Liability%20Not%20a%20Liability%20Textual%20Analysis,%20Dictionaries,%20and%2010-Ks%206989-IA-Feb-2011.pdf>

Loughran T, McDonald B (2014) Measuring Readability in Financial Disclosures, in: The Journal of Finance (JoF), 69(4)/2014, 1643-1671

Loughran T, McDonald B (2016) Textual Analysis in Accounting and Finance: A Survey, in: Journal of Accounting Research (JAR), 54(4)/2016, 1187-1230

Mooney R J, Bunescu R (2005) Mining Knowledge from Text Using Information Extraction, in: SIGKDD Explorations (SIGKDD), 7(1)/2005, 3-10

O'Riain S (2012) Semantic Paths in Business Filings Analysis. Ph.D. thesis, National University of Ireland, Galway, Ireland

Pagell R A (1995) EDGAR: Electronic Data Gathering and Receiving, in: Business Information Review (BIR), 11(3)/1995, 56-68

Palmer D D (2010) Text Preprocessing, in: N Indurkha, F J Damerau, Eds. Handbook of Natural Language Processing, Second Edition, Taylor & Francis Group, Boca Raton, Florida, USA

SEC (1934) Securities Exchange Act of 1934, Available online on URL: <https://www.sec.gov/about/laws/sea34.pdf>

SEC (2000) HTML Specifications for EDGAR Rel. 7.0, Available online on URL: <https://www.sec.gov/info/edgar/ednews/edhtml.htm>

SEC (2006) Electronic Filing and the EDGAR System: A Regulatory Overview, Available online on URL: <https://www.sec.gov/info/edgar/regoverview.htm>

SEC (2010) Important Information about EDGAR, Available online on URL: <https://www.sec.gov/edgar/aboutedgar.htm>

SEC (2011) Fast Answers – How to Read a 10-K, Available online on URL: <https://www.sec.gov/answers/reada10k.htm>

SEC (2013) What We Do, Available online on URL: <https://www.sec.gov/about/whatwedo.shtml>

SEC (2015) Information for FTP Users, Available online on URL: <https://www.sec.gov/edgar/searchedgar/ftpusers.htm>

SEC EDGAR (2016) EDGAR Public Dissemination Service (PDS) System, Available online on URL: <https://www.sec.gov/oit/announcement/public-dissemination-service-system-contact.html>

SEC EDGAR (2015) Public Dissemination Service (PDS) Technical Specification, Available online on URL: https://www.sec.gov/info/edgar/specifications/pds_dissemination_spec.pdf

SEC EDGAR Archives (2016) Coca Cola Company's Financial Statement Submissions on 2016-02-25, Available online on URL: <https://www.sec.gov/Archives/edgar/data/21344/000002134416000050/0000021344-16-000050-index.htm>

SEC Form Glossary (2015) Index to Forms, Available online on URL: <https://www.sec.gov/info/edgar/forms/edgform.pdf>

SEC Index Files (2016) Full Index Files, Available online on URL: <ftp://ftp.sec.gov/edgar/full-index/>

SEC Regulation S-K (2016) Standard Instructions for filing Forms under Securities Act of 1933, Securities Exchange Act of 1934 and Energy Policy and Conservation Act of 1975-Regulation S-K, Available online on URL: http://www.ecfr.gov/cgi-bin/text-idx?SID=8e0ed509ccc65e983f9eca72ceb26753&node=17:3.0.1.1.11&rgn=div5#se17.3.229_1101

SEC Regulation S-T (2016) General Rules and Regulations for electronic Filings, Available online on URL: http://www.ecfr.gov/cgi-bin/text-idx?node=17:3.0.1.1.14&rgn=div5#se17.3.232_1100

SEC Regulation S-X (2016) Form and Content of and Requirements for Financial Statements; Securities Act of 1933, Securities Exchange Act of 1934, Investment Company Act of 1940, Investments Advisers Act of 1940, and Energy Policy and Conservation Act of 1975- Regulation S-X, Available online on URL: http://www.ecfr.gov/cgi-bin/text-idx?SID=8e0ed509ccc65e983f9eca72ceb26753&node=17:3.0.1.1.8&rgn=div5#se17.3.210_11_601

SEC Release 33-8099 (2002) Mandated EDGAR Filing for Foreign Issuers, Available online on URL: <https://www.sec.gov/rules/final/33-8099.htm>

SEC Release 34-36997 (1996) EDGAR Phase-in Complete on May 6, 1996, Available online on URL: <https://www.sec.gov/info/edgar/ednews/34-36997.htm>

Srivastava R P (2016) Textual Analysis and Business Intelligence in Big Data Environment: Search Engine versus XBRL, in: Indian Accounting Review (IAR), 20(1)/2016, 1-20

Stümpert T (2008) Extracting Financial Data from SEC Filings for US GAAP Accountants, in: D Seese, C Weinhardt, F Schlottmann, Eds. Handbook on Information Technology in Finance. Springer-Verlag, Berlin Heidelberg, Germany

Stümpert T, Seese D, Centinkaya Ö, Spöth R (2004) EASE – a software agent that extracts financial data from the SEC's EDGAR database, in: Proceedings of the 4th International ICSC Symposium on Engineering of Intelligent Systems (EIS 2004). Funchal, Portugal

Tetlock P C (2007) Giving Content to Investor Sentiment: The Role of the Media in the Stock Market, in: The Journal of Finance (Jof), 62(3)/2007, 1139-1168

Thai V, Davis B, O'Riain S, O'Sullivan D, Handschuh S (2008) Semantically Enhanced Passage Retrieval for Business Analysis Activity, in: Proceedings of the 16th European Conference on Information Systems (ECIS 2008). Galway, Ireland

Thompson K (1968) Regular Expression Search Algorithm, in: Communications of the ACM (CACM), 11(6)/1968, 419-422

Wilks Y (1997) Information Extraction as a Core Language Technology, in: M T Pazienza, Ed. Information Extraction: A Multidisciplinary Approach to an Emerging Information Technology. Springer-Verlag, Berlin Heidelberg, Germany

W3 Schools (2016), HTML Character Sets, Available online on URL:
<http://www.w3schools.com/charsets/default.asp>

W3C Recommendation (1999) HTML 4.01 Specification, Available online on URL:
<https://www.w3.org/TR/html401/cover.html>

W3C Strict DTD (1999) HTML 4.01 Strict DTD, Available online on URL:
<https://www.w3.org/TR/html4/strict.dtd>

Appendix B: Form Index Files - Descriptive Statistics

Sub_Type	Total	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016		
1	33	0	0	0	0	0	0	0	0	0	6	1	4	2	3	5	1	4	0	2	2	0	2	1	0		
3	538228	0	0	0	1149	2983	4317	5282	6852	5327	8694	27814	47801	47356	46869	52898	37477	31171	33279	32264	30265	31834	36150	32839	15607		
4	5850937	0	0	0	4312	12758	18432	19899	22606	21689	56942	324449	476172	474215	468480	486241	456300	387113	404831	393744	398395	400134	400454	397908	225863		
5	186884	0	0	0	301	2189	3446	4741	5346	5632	6560	16067	19881	16593	15040	13070	11495	10620	9374	9003	8427	8084	8097	6991	5927		
25	3642	0	0	0	0	0	0	0	0	1	0	16	406	468	445	425	335	362	325	205	118	111	108	102	89	76	50
26	1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
144	12111	0	0	0	9	32	127	107	128	99	193	9	775	129	1082	1209	687	672	773	790	839	989	856	530	244	0	
425	81918	0	0	0	0	0	0	0	4303	4849	4271	4176	4878	6601	4971	5834	3736	4535	4339	4932	3020	4295	6732	7304	3152	0	
487	15478	0	346	352	393	394	461	537	552	512	460	473	441	498	564	679	723	874	1042	1074	1101	1063	1130	1187	622	0	
497	365987	0	6285	10675	13227	17378	14624	13905	14883	19491	15718	16738	20578	18592	19068	17563	19155	18238	18325	15496	16586	19076	16923	17153	9950	0	
1-A	413	0	0	0	0	0	0	0	0	0	1	0	0	0	2	1	0	0	0	0	31	94	135	136	3	11	
10-12B	578	0	4	4	37	37	25	29	29	23	19	31	12	15	19	33	39	21	19	27	22	31	38	45	19	0	
10-12B/A	1318	0	1	2	69	50	44	50	69	48	40	30	30	15	44	69	106	38	64	67	66	70	150	138	58	0	
10-12G	2106	0	2	3	49	76	94	148	84	58	57	38	43	58	75	63	271	159	225	143	118	103	99	92	48	0	
10-12G/A	3861	0	4	6	91	146	136	234	235	114	107	80	63	104	116	91	337	274	423	414	262	222	201	134	67	0	
10-C	1303	0	82	396	745	80	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-C/A	25	0	2	12	9	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-D	45492	0	0	0	0	0	0	0	0	0	0	0	0	3	6950	8081	2412	891	934	1775	2983	4571	6035	7097	3760	0	
10-D/A	2040	0	0	0	0	0	0	0	0	0	0	0	0	0	483	720	370	54	23	20	55	124	84	69	38	0	
10-K	167599	4	1912	2218	4315	6698	6930	6761	6652	6248	6759	8468	8567	9017	8852	8574	8746	9839	9165	8840	8393	8105	8084	7985	6467	0	
10-K/A	38888	2	616	933	1495	2152	1943	1798	1530	1578	2010	2021	2096	2180	1510	1470	1801	2320	2213	1995	1840	1765	1557	1258	805	0	
10-K405	21277	0	11	1018	1944	3201	3357	3361	3217	3000	2168	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10-K405/A	3349	0	5	190	317	538	589	610	483	494	123	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-KSB	36912	0	1	95	1052	2013	2001	2173	2856	3199	3419	3410	3557	3458	3399	3513	2686	80	0	0	0	0	0	0	0	0	0
10-KSB	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-KSB/A	11909	0	0	63	494	713	582	518	724	833	858	1024	955	1380	1382	1009	1284	90	0	0	0	0	0	0	0	0	0
10-KSB/A	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-KSB40	3441	0	0	24	213	540	532	547	661	617	307	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-KSB40/A	625	0	0	7	42	103	129	110	98	116	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-KT	418	0	1	6	2	18	14	9	12	18	22	25	38	19	13	11	111	21	27	23	22	41	22	29	0	0	
10-KT/A	86	0	0	2	1	5	1	6	2	6	3	5	6	3	2	5	3	2	9	9	5	17	8	5	5	0	
10KT405	106	0	0	3	14	8	7	10	12	20	12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10KT405/A	17	0	0	3	6	3	1	1	2	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10-Q	522906	7	6628	14120	25749	28996	29239	28687	28282	25997	24084	21866	20872	20670	20038	20047	27040	27817	26550	25676	24186	23250	22876	22166	8063	0	
10-Q/A	41651	0	558	1692	2319	2036	2135	2426	1926	1722	1778	1694	1704	1946	1574	1065	1517	2195	2060	3539	3045	1829	1420	1095	376	0	
10QSB	120120	0	15	1297	5670	7448	7455	8419	11619	11581	11160	10348	10671	10209	10191	10703	3334	0	0	0	0	0	0	0	0	0	0
10-QSB	4	0	0	0	0	0	0	0	0	2	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10QSB/A	17117	0	1	198	884	984	805	883	1325	1428	1464	1527	1355	1840	2051	1585	787	0	0	0	0	0	0	0	0	0	0
10-QSB/A	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10-QT	197	0	3	9	8	18	12	13	17	12	25	7	5	6	10	7	6	5	4	0	5	8	7	7	3	0	
10-QT/A	33	0	1	2	1	10	3	1	2	0	0	3	1	0	1	0	1	1	2	1	1	2	0	1	0	0	
10SB12B	119	0	0	0	9	12	5	22	12	11	7	3	5	12	13	6	2	0	0	0	0	0	0	0	0	0	0
10SB12B/A	124	0	0	0	6	8	8	31	31	25	3	3	0	1	4	4	0	0	0	0	0	0	0	0	0	0	0
10SB12G	4105	0	0	0	71	113	175	1048	1140	348	230	109	138	170	233	274	56	0	0	0	0	0	0	0	0	0	0
10SB12G/A	5523	0	0	0	79	167	167	1037	1672	694	480	282	184	250	235	253	23	0	0	0	0	0	0	0	0	0	0
11-K	36426	0	811	611	809	949	1127	1264	1411	1550	2412	2476	2369	2251	2162	2058	1913	1824	1711	1626	1569	1508	1436	1379	1200	0	
11-K/A	911	0	37	22	21	26	34	38	33	25	47	95	129	45	43	44	42	45	30	30	25	16	13	58	13	0	
11-KT	88	1	4	1	6	14	5	5	1	5	5	10	2	3	8	3	5	0	2	1	2	2	2	1	0	0	
11-KT/A	2	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1	0	0	
12G3-2B	609	0	0	0	0	0	0	0	3	9	66	66	0	92	109	100	78	0	0	0	0	0	0	0	0	0	0
12G33BR	178	0	0	0	0	0	0	0	0	0	0	45	26	32	24	15	4	0	0	0	0	0	0	0	0	0	0
13FCONP	1249	6	5	4	10	13	18	33	92	145	4	143	133	103	85	0	99	72	43	112	20	1	0	0	1	0	
13FCONP/A	99	0	4	10	12	15	20	15	5	12	0	6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
13F-E	1188	0	193	232	222	212	252	77	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13F-E/A	120	0	5	28	7	15	25	40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13F-HR	193463	4	4	8	5	8	34	5079	7428	8180	8212	8416	8821	9550	10418	11301	12299	11996	11687	12630	13246	13829	15239	16562	8507	0	
13F-HR/A	20575	0	0	0	0	3	1	0	474	770	979	1074	1055	1228	1199	1353	1435	1803	1109	1234	1525	1264	1300	949	1278	542	
13F-NT	64392	0	0	0	0	0	0	968	1657	1811	1941	2194	2397	2795	3215	3678	4223	4371	4822	5026	5098	5535	5586	6049	3026	0	
13F-NT/A	1514	0	0	0	0	0	0	27	37	44	100	134	182	100	125	112	115	122	57	74	109	84	16	45	31	0	
144/A	455	0	0	0	0	6	50	14	24	11	8	25	37	42	38	49	13	14	34	18	12	19	24	12	5	0	
15-12B	4302	0	8	23	49</																						

38

39

40

41

42

43

U-57	617	0	2	12	42	92	77	92	96	74	61	23	28	9	9	0	0	0	0	0	0	0	0	0	0	0
U-57/A	240	0	0	0	3	16	14	42	9	50	54	23	19	10	0	0	0	0	0	0	0	0	0	0	0	0
USA	40	0	0	0	0	1	3	1	14	6	7	1	4	2	1	0	0	0	0	0	0	0	0	0	0	0
USA/A	4	0	0	0	0	0	0	0	1	0	0	0	2	1	0	0	0	0	0	0	0	0	0	0	0	0
USB	32	0	0	1	0	1	3	0	6	10	5	0	3	3	0	0	0	0	0	0	0	0	0	0	0	0
USB/A	10	0	0	1	0	0	1	0	0	3	2	2	1	0	0	0	0	0	0	0	0	0	0	0	0	0
USS	248	0	10	15	16	16	17	20	23	20	26	27	29	29	0	0	0	0	0	0	0	0	0	0	0	0
USS/A	68	0	3	2	3	3	7	4	5	1	12	9	10	9	0	0	0	0	0	0	0	0	0	0	0	0
U-6B-2	1291	0	39	49	50	52	115	91	91	128	164	177	170	157	8	0	0	0	0	0	0	0	0	0	0	0
U-6B-2/A	9	0	0	0	0	1	2	1	0	0	1	3	1	0	0	0	0	0	0	0	0	0	0	0	0	0
U-7D	38	0	2	6	0	3	9	0	2	7	3	3	2	1	0	0	0	0	0	0	0	0	0	0	0	0
U-7D/A	66	0	0	4	4	4	26	18	6	4	2	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0
U-9C-3	636	0	0	0	22	55	63	62	75	80	0	82	95	103	0	0	0	0	0	0	0	0	0	0	0	0
U-9C-3/A	16	0	0	0	1	4	2	1	4	3	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
UNDER	41	0	0	0	0	0	0	0	7	4	7	1	2	5	1	3	2	0	0	0	0	0	0	0	0	0
UNDER/A	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
UPLOAD	159065	0	0	0	0	0	0	0	0	0	0	641	8844	13289	12301	13095	15324	17093	17930	15413	16906	12043	13119	3067		
WDL-REQ	55	0	0	0	0	0	0	0	0	0	5	8	2	1	4	2	6	8	3	4	1	4	6	1		
X-17A-5	71022	0	0	0	0	2	0	0	7	5538	5410	5281	5083	5170	5124	4966	4793	4718	4543	4464	4308	4108	4095	3412		
X-17A-5/A	4936	0	0	0	0	1	2	0	0	2	337	284	477	434	441	441	611	279	152	135	92	319	165			
Total	15998058	35	65461	113651	202614	280130	313300	321262	341962	351443	448884	790772	1024211	1072291	1096624	1137485	1036832	938511	979708	979632	974937	990462	996692	994215	546944	

Appendix C: „Annual Report Algorithm“

ID.	Description	Regular Expression	Repl. by
1	Decomposition of “Complete Submission Text File”	(?s)<DOCUMENT>.*?</DOCUMENT>	
2	Identification of document (file) types	<TYPE>.*	
3	Removal of graphic files	(?s)<TYPE>GRAPHIC.*?</TEXT>	“ ”
4	Removal of MS Excel files	(?s)<TYPE>EXCEL.*?</TEXT>	“ ”
5	Removal of PDF files	(?s)<TYPE>PDF.*?</TEXT>	“ ”
6	Removal of ZIP files	(?s)<TYPE>ZIP.*?</TEXT>	“ ”
7	Removal of cover letter	(?s)<TYPE>COVER.*?</TEXT>	“ ”
8	Removal of correspondance between SEC staff and EDGAR participant	(?s)<TYPE>CORRESP.*?</TEXT>	“ ”
9	Removal of XBRL instance document	(?s)<TYPE>EX-10[01].INS.*?</TEXT>	“ ”
10	Removal of XBRL instance document	(?s)<TYPE>EX-99.SDR [KL].INS.*?</TEXT>	“ ”
11	Removal of XBRL taxonomy extension schema document	(?s)<TYPE>EX-10[01].SCH.*?</TEXT>	“ ”
12	Removal of XBRL taxonomy extension schema document	(?s)<TYPE>EX-99.SDR [KL].SCH.*?</TEXT>	“ ”
13	Removal of XBRL taxonomy extension linkbase document	(?s)<TYPE>EX-10[01].CAL.*?</TEXT>	“ ”
14	Removal of XBRL taxonomy extension linkbase document	(?s)<TYPE>EX-99.SDR [KL].CAL.*?</TEXT>	“ ”
15	Removal of XBRL taxonomy extension definition linkbase document	(?s)<TYPE>EX-10[01].DEF.*?</TEXT>	“ ”
16	Removal of XBRL taxonomy extension definition linkbase document	(?s)<TYPE>EX-99.SDR [KL].LAB.*?</TEXT>	“ ”
17	Removal of XBRL taxonomy extension labels linkbase document	(?s)<TYPE>EX-10[01].LAB.*?</TEXT>	“ ”
18	Removal of XBRL taxonomy extension labels linkbase document	(?s)<TYPE>EX-99.SDR [KL].LAB.*?</TEXT>	“ ”
19	Removal of XBRL taxonomy extension presentation linkbase document	(?s)<TYPE>EX-10[01].PRE.*?</TEXT>	“ ”
20	Removal of XBRL taxonomy extension presentation linkbase document	(?s)<TYPE>EX-99.SDR [KL].PRE.*?</TEXT>	“ ”
21	Removal of XBRL taxonomy extension reference linkbase document	(?s)<TYPE>EX-10[01].REF.*?</TEXT>	“ ”
22	Removal of XBRL documents	(?s)<TYPE>XML.*?</TEXT>	“ ”
23	Removal of document type information	<TYPE>.*	“ ”
24	Removal of sequence information	<SEQUENCE>.*	“ ”
25	Removal of filename	<FILENAME>.*	“ ”
26	Removal of description	<DESCRIPTION>.*	“ ”
27	Removal of head section (including document title)	(?s)(?i)<Head>.*?</Head>	“ ”
28	Removal of table content	(?s)(?i)<Table>.*?</Table>	“ ”
29	Removal of HTML tags and attributes	(?s)<[^>]*>	“ ”
<p style="text-align: right;">ASCII (ISO 8859-1/ANSI) character set: 32-126</p> <p style="text-align: right;">[http://www.w3schools.com/charsets/ref_html_ascii.asp] [http://www.w3schools.com/charsets/ref_html_8859.asp] [http://www.w3schools.com/charsets/ref_html_ansi.asp] [https://www.sec.gov/info/edgar/edgarfm-vol2-v37.pdf] [http://www.htmlhelp.com/reference/charset/]</p>			
30	Decoding normal space	 	“ ”
31	Decoding exclamation mark	!	“!”
32	Decoding quotation mark	"	“”
33	Decoding quotation mark	"	“”

34	Decoding number sign	#	“#”
35	Decoding dollar sign	$	“\$”
36	Decoding percent sign	%	“%”
37	Decoding ampersand	&	“&”
38	Decoding ampersand	&	“&”
39	Decoding apostrophe	'	“’”
40	Decoding left parenthesis	(“(“
41	Decoding right parenthesis)	“)”
42	Decoding asterisk	*	“*”
43	Decoding plus sign	+	“+”
44	Decoding comma	,	“,”
45	Decoding hyphen-minus	-	“-”
46	Decoding full stop	.	“.”
47	Decoding solidus	/	“/”
48	Decoding digit zero	0	“0”
49	Decoding digit one	1	“1”
50	Decoding digit two	2	“2”
51	Decoding digit three	3	“3”
52	Decoding digit four	4	“4”
53	Decoding digit five	5	“5”
54	Decoding digit six	6	“6”
55	Decoding digit seven	7	“7”
56	Decoding digit eight	8	“8”
57	Decoding digit nine	9	“9”
58	Decoding colon	:	“.”
59	Decoding semicolon	;	“,”
60	Decoding less-than sign	<	“<”
61	Decoding less-than sign	<	“<”
62	Decoding equals sign	=	“=”
63	Decoding greater-than sign	>	“>”
64	Decoding greater-than sign	>	“>”
65	Decoding question mark	?	“?”
66	Decoding commercial at	@	“@”
67	Decoding latin capital letter A	A	“A”
68	Decoding latin capital letter B	B	“B”
69	Decoding latin capital letter C	C	“C”
70	Decoding latin capital letter D	D	“D”
71	Decoding latin capital letter E	E	“E”
72	Decoding latin capital letter F	F	“F”
73	Decoding latin capital letter G	G	“G”
74	Decoding latin capital letter H	H	“H”
75	Decoding latin capital letter I	I	“I”
76	Decoding latin capital letter J	J	“J”

77	Decoding latin capital letter K	K	“K”
78	Decoding latin capital letter L	L	“L”
79	Decoding latin capital letter M	M	“M”
80	Decoding latin capital letter N	N	“N”
81	Decoding latin capital letter O	O	“O”
82	Decoding latin capital letter P	P	“P”
83	Decoding latin capital letter Q	Q	“Q”
84	Decoding latin capital letter R	R	“R”
85	Decoding latin capital letter S	S	“S”
86	Decoding latin capital letter T	T	“T”
87	Decoding latin capital letter U	U	“U”
88	Decoding latin capital letter V	V	“V”
89	Decoding latin capital letter W	W	“W”
90	Decoding latin capital letter X	X	“X”
91	Decoding latin capital letter Y	Y	“Y”
92	Decoding latin capital letter Z	Z	“Z”
93	Decoding left square bracket	[[“
94	Decoding reverse solidus	\	“\”
95	Decoding right square bracket]	”]
96	Decoding circumflex accent	^	“^”
97	Decoding low line	_	“_”
98	Decoding grave accent	`	“`”
99	Decoding latin small letter a	a	“a”
100	Decoding latin small letter b	b	“b”
101	Decoding latin small letter c	c	“c”
102	Decoding latin small letter d	d	“d”
103	Decoding latin small letter e	e	“e”
104	Decoding latin small letter f	f	“f”
105	Decoding latin small letter g	g	“g”
106	Decoding latin small letter h	h	“h”
107	Decoding latin small letter i	i	“i”
108	Decoding latin small letter j	j	“j”
109	Decoding latin small letter k	k	“k”
110	Decoding latin small letter l	l	“l”
111	Decoding latin small letter m	m	“m”
112	Decoding latin small letter n	n	“n”
113	Decoding latin small letter o	o	“o”
114	Decoding latin small letter p	p	“p”
115	Decoding latin small letter q	q	“q”
116	Decoding latin small letter r	r	“r”
117	Decoding latin small letter s	s	“s”
118	Decoding latin small letter t	t	“t”
119	Decoding latin small letter u	u	“u”

120	Decoding latin small letter v	v	“v”
121	Decoding latin small letter w	w	“w”
122	Decoding latin small letter x	x	“x”
123	Decoding latin small letter y	y	“y”
124	Decoding latin small letter z	z	“z”
125	Decoding left curly bracket	{	“{”
126	Decoding vertical line	|	“ ”
127	Decoding right curly bracket	}	“}”
128	Decoding tilde	~	“~”
<p style="text-align: right;">ANSI (Windows-1252) character set: 128-159</p> <p style="text-align: right;">[http://www.w3schools.com/charsets/ref_html_ansi.asp] [https://www.sec.gov/info/edgar/edgarfm-vol2-v37.pdf] [http://www.htmlhelp.com/reference/charset/]</p>			
129	Decoding euro sign	€	“[Euro]”
130	Decoding euro sign	€	“[Euro]”
131	Decoding single low-9 quotation mark	‚	“’”
132	Decoding single low-9 quotation mark	‚	“’”
133	Decoding latin small letter f with hook	ƒ	“f”
134	Decoding latin small letter f with hook	ƒ	“f”
135	Decoding double low-9 quotation mark	„	“”
136	Decoding double low-9 quotation mark	„	“”
137	Decoding horizontal ellipsis	…	“...”
138	Decoding horizontal ellipsis	…	“...”
139	Decoding dagger	†	“†”
140	Decoding dagger	†	“†”
141	Decoding double dagger	‡	“‡”
142	Decoding double dagger	‡	“‡”
143	Decoding modifier letter circumflex accent	ˆ	“ˆ”
144	Decoding modifier letter circumflex accent	ˆ	“ˆ”
145	Decoding per mille sign	‰	“‰”
146	Decoding per mille sign	‰	“‰”
147	Decoding latin capital letter S with caron	Š	“Š”
148	Decoding latin capital letter S with caron	Š	“Š”
149	Decoding single left-pointing angle quotation mark	‹	“‹”
150	Decoding single left-pointing angle quotation mark	‹	“‹”
151	Decoding latin capital ligature OE	Œ	“Œ”
152	Decoding latin capital ligature OE	Œ	“Œ”
153	Decoding Latin capital letter Z with caron	Ž	“Ž”
154	Decoding Latin capital letter Z with caron	Ž	“Ž”
155	Decoding left single quotation mark	‘	“’”
156	Decoding left single quotation mark	‘	“’”
157	Decoding right single quotation mark	’	“’”

158	Decoding right single quotation mark	’	“ ’ ”
159	Decoding left double quotation mark	“	“ ”
160	Decoding left double quotation mark	“	“ ”
161	Decoding right double quotation mark	”	“ ”
162	Decoding right double quotation mark	”	“ ”
163	Decoding bullet	•	“ ”
164	Decoding bullet	•	“ ”
165	Decoding en dash	–	“ _ ”
166	Decoding en dash	–	“ _ ”
167	Decoding em dash	—	“ _ ”
168	Decoding em dash	—	“ _ ”
169	Decoding small tilde	˜	“ ~ ”
170	Decoding small tilde	˜	“ ~ ”
171	Decoding trade mark sign	™	“ [trade mark sign] ”
172	Decoding trade mark sign	™	“ [trade mark sign] ”
173	Decoding latin small letter s with caron	š	“ Š ”
174	Decoding latin small letter s with caron	š	“ Š ”
175	Decoding single right-pointing angle quotation mark	›	“ ’ ”
176	Decoding single right-pointing angle quotation mark	›	“ ’ ”
177	Decoding latin small ligature oe	œ	“ œ ”
178	Decoding latin small ligature oe	œ	“ œ ”
179	Decoding latin small letter z with caron	ž	“ Ž ”
180	Decoding latin small letter z with caron	ž	“ Ž ”
181	Decoding latin capital letter Y with diaeresis	Ÿ	“ Ÿ ”
182	Decoding latin capital letter Y with diaeresis	Ÿ	“ Ÿ ”
<p style="text-align: right;">ISO 8859-1 (Latin 1) character set: 160-255</p> <p style="text-align: right;">[https://www.w3.org/TR/html4/sgml/entities.html] [http://www.w3schools.com/charsets/ref_html_8859.asp] [http://www.w3schools.com/charsets/ref_html_ansi.asp] [https://www.sec.gov/info/edgar/edgarfm-vol2-v37.pdf] [http://www.htmlhelp.com/reference/charset/]</p>			
183	Decoding non-breaking space	 	“ ”
184	Decoding non-breaking space	 	“ ”
185	Decoding inverted exclamation	¡	“ ¡ ”
186	Decoding inverted exclamation	¡	“ ¡ ”
187	Decoding cent sign	¢	“ [Cent] ”
188	Decoding cent sign	¢	“ [Cent] ”
189	Decoding pound sign	£	“ [Pound Sterling] ”
190	Decoding pound sign	£	“ [Pound Sterling] ”
191	Decoding currency sign	¤	“ [Currency] ”
192	Decoding currency sign	¤	“ [Currency] ”
193	Decoding yen sign	¥	“ [Yen] ”

194	Decoding yen sign	¥	“[Yen]”
195	Decoding broken bar	¦	“ ”
196	Decoding broken bar	¦	“ ”
197	Decoding section sign	§	“§”
198	Decoding section sign	§	“§”
199	Decoding diaeresis	¨	“ ”
200	Decoding diaeresis	¨	“ ”
201	Decoding diaeresis	¨	“ ”
202	Decoding copyright sign	©	“[copyright sign]”
203	Decoding copyright sign	©	“[copyright sign]”
204	Decoding feminine ordinal indicator	ª	“ ”
205	Decoding feminine ordinal indicator	ª	“ ”
206	Decoding left-pointing double angle quotation mark	«	“”
207	Decoding left-pointing double angle quotation mark	«	“”
208	Decoding not sign	¬	“ ”
209	Decoding not sign	¬	“ ”
210	Decoding soft hyphen	­	“ ”
211	Decoding soft hyphen	­	“ ”
212	Decoding registered sign	®	“[registered trademark sign]”
213	Decoding registered sign	®	“[registered trademark sign]”
214	Decoding macron	¯	“ ”
215	Decoding macron	¯	“ ”
216	Decoding degree sign	°	“ ”
217	Decoding degree sign	°	“ ”
218	Decoding plus-minus sign	±	“[+/-]”
219	Decoding plus-minus sign	±	“[+/-]”
220	Decoding superscript two	²	“ ”
221	Decoding superscript two	²	“ ”
222	Decoding superscript three	³	“ ”
223	Decoding superscript three	³	“ ”
224	Decoding acute accent	´	“ ”
225	Decoding acute accent	´	“ ”
226	Decoding micro sign	µ	“ ”
227	Decoding micro sign	µ	“ ”
228	Decoding pilcrow sign	¶	“ ”
229	Decoding pilcrow sign	¶	“ ”
230	Decoding middle dot	·	“ ”
231	Decoding middle dot	·	“ ”
232	Decoding cedilla	¸	“ ”
233	Decoding cedilla	¸	“ ”
234	Decoding superscript one	¹	“ ”
235	Decoding superscript one	¹	“ ”
236	Decoding masculine ordinal	º	“ ”

237	Decoding masculine ordinal	º	“ ¨ ”
238	Decoding right-pointing double angle quotation mark	»	“ ”
239	Decoding right-pointing double angle quotation mark	»	“ ”
240	Decoding vulgar fraction one quarter	¼	“1/4”
241	Decoding vulgar fraction one quarter	¼	“1/4”
242	Decoding vulgar fraction one half	½	“1/2”
243	Decoding vulgar fraction one half	½	“1/2”
244	Decoding vulgar fraction one half	½	“1/2”
245	Decoding vulgar fraction three quarters	¾	“3/4”
246	Decoding vulgar fraction three quarters	¾	“3/4”
247	Decoding inverted question mark	¿	“?”
248	Decoding inverted question mark	¿	“?”
249	Decoding latin capital letter A with grave	À	“A”
250	Decoding latin capital letter A with grave	&#Agrave	“A”
251	Decoding latin capital letter A with acute	Á	“A”
252	Decoding latin capital letter A with acute	Á	“A”
253	Decoding latin capital letter A with circumflex	Â	“A”
254	Decoding latin capital letter A with circumflex	Â	“A”
255	Decoding latin capital letter A with tilde	Ã	“A”
256	Decoding latin capital letter A with tilde	Ã	“A”
257	Decoding latin capital letter A with diaeresis	Ä	“A”
258	Decoding latin capital letter A with diaeresis	Ä	“A”
259	Decoding latin capital letter A with ring above	Å	“A”
260	Decoding latin capital letter A with ring above	Å	“A”
261	Decoding latin capital letter AE	Æ	“AE”
262	Decoding latin capital letter AE	Æ	“AE”
263	Decoding latin capital letter C with cedilla	Ç	“C”
264	Decoding latin capital letter C with cedilla	Ç	“C”
265	Decoding latin capital letter E with grave	È	“E”
266	Decoding latin capital letter E with grave	È	“E”
267	Decoding latin capital letter E with acute	É	“E”
268	Decoding latin capital letter E with acute	É	“E”
269	Decoding latin capital letter E with circumflex	Ê	“E”
270	Decoding latin capital letter E with circumflex	Ê	“E”
271	Decoding latin capital letter E with diaeresis	Ë	“E”
272	Decoding latin capital letter E with diaeresis	Ë	“E”
273	Decoding latin capital letter I with grave	Ì	“I”
274	Decoding latin capital letter I with grave	Ì	“I”
275	Decoding latin capital letter I with acute	Í	“I”
276	Decoding latin capital letter I with acute	Í	“I”
277	Decoding latin capital letter I with circumflex	Î	“I”
278	Decoding latin capital letter I with circumflex	Î	“I”
279	Decoding latin capital letter I with diaeresis	Ï	“I”

280	Decoding latin capital letter I with diaeresis	Ï	“I”
281	Decoding latin capital letter ETH	Ð	“ETH”
282	Decoding latin capital letter ETH	Ð	“ETH”
283	Decoding latin capital letter N with tilde	Ñ	“N”
284	Decoding latin capital letter N with tilde	Ñ	“N”
285	Decoding latin capital letter O with grave	Ò	“O”
286	Decoding latin capital letter O with grave	Ò	“O”
287	Decoding latin capital letter O with acute	Ó	“O”
288	Decoding latin capital letter O with acute	Ó	“O”
289	Decoding latin capital letter O with circumflex	Ô	“O”
290	Decoding latin capital letter O with circumflex	Ô	“O”
291	Decoding latin capital letter O with tilde	Õ	“O”
292	Decoding latin capital letter O with tilde	Õ	“O”
293	Decoding latin capital letter O with diaeresis	Ö	“O”
294	Decoding latin capital letter O with diaeresis	Ö	“O”
295	Decoding multiplication sign	×	“*”
296	Decoding multiplication sign	×	“*”
297	Decoding latin capital letter O with stroke	Ø	“O”
298	Decoding latin capital letter O with stroke	Ø	“O”
299	Decoding latin capital letter U with grave	Ù	“U”
300	Decoding latin capital letter U with grave	Ù	“U”
301	Decoding latin capital letter U with acute	Ú	“U”
302	Decoding latin capital letter U with acute	Ú	“U”
303	Decoding latin capital letter U with circumflex	Û	“U”
304	Decoding latin capital letter U with circumflex	Û	“U”
305	Decoding latin capital letter U with diaeresis	Ü	“U”
306	Decoding latin capital letter U with diaeresis	Ü	“U”
307	Decoding latin capital letter Y with acute	Ý	“Y”
308	Decoding latin capital letter Y with acute	Ý	“Y”
309	Decoding latin capital letter THORN	Þ	“THORN”
310	Decoding latin capital letter THORN	Þ	“THORN”
311	Decoding latin small letter sharp s	ß	“ß”
312	Decoding latin small letter sharp s	ß	“ß”
313	Decoding latin small letter a with grave	à	“a”
314	Decoding latin small letter a with grave	à	“a”
315	Decoding latin small letter a with acute	á	“a”
316	Decoding latin small letter a with acute	á	“a”
317	Decoding latin small letter a with circumflex	â	“a”
318	Decoding latin small letter a with circumflex	â	“a”
319	Decoding latin small letter a with tilde	ã	“a”
320	Decoding latin small letter a with tilde	ã	“a”
321	Decoding latin small letter a with diaeresis	ä	“a”
322	Decoding latin small letter a with diaeresis	ä	“a”

323	Decoding latin small letter a with ring above	å	“a”
324	Decoding latin small letter a with ring above	å	“a”
325	Decoding latin small letter ae	æ	“ae”
326	Decoding latin small letter ae	æ	“ae”
327	Decoding latin small letter c with cedilla	ç	“c”
328	Decoding latin small letter c with cedilla	ç	“c”
329	Decoding latin small letter e with grave	è	“e”
330	Decoding latin small letter e with grave	è	“e”
331	Decoding latin small letter e with acute	é	“e”
332	Decoding latin small letter e with acute	é	“e”
333	Decoding latin small letter e with circumflex	ê	“e”
334	Decoding latin small letter e with circumflex	ê	“e”
335	Decoding latin small letter e with diaeresis	ë	“e”
336	Decoding latin small letter e with diaeresis	ë	“e”
337	Decoding latin small letter i with grave	ì	“i”
338	Decoding latin small letter i with grave	ì	“i”
339	Decoding latin small letter i with acute	í	“i”
340	Decoding latin small letter i with acute	í	“i”
341	Decoding latin small letter i with circumflex	î	“i”
342	Decoding latin small letter i with circumflex	î	“i”
343	Decoding latin small letter i with diaeresis	ï	“i”
344	Decoding latin small letter i with diaeresis	ï	“i”
345	Decoding latin small letter eth	ð	“eth”
346	Decoding latin small letter eth	ð	“eth”
347	Decoding latin small letter n with tilde	ñ	“n”
348	Decoding latin small letter n with tilde	ñ	“n”
349	Decoding latin small letter o with grave	ò	“o”
350	Decoding latin small letter o with grave	ò	“o”
351	Decoding latin small letter o with acute	ó	“o”
352	Decoding latin small letter o with acute	ó	“o”
353	Decoding latin small letter o with circumflex	ô	“o”
354	Decoding latin small letter o with circumflex	ô	“o”
355	Decoding latin small letter o with tilde	õ	“o”
356	Decoding latin small letter o with tilde	õ	“o”
357	Decoding latin small letter o with diaeresis	ö	“o”
358	Decoding latin small letter o with diaeresis	ö	“o”
359	Decoding division sign	÷	“/”
360	Decoding division sign	÷	“/”
361	Decoding latin small letter o with stroke	ø	“o”
362	Decoding latin small letter o with stroke	ø	“o”
363	Decoding latin small letter u with grave	ù	“u”
364	Decoding latin small letter u with grave	ù	“u”
365	Decoding latin small letter u with acute	ú	“u”

366	Decoding latin small letter u with acute	ú	“u”
367	Decoding latin small letter u with circumflex	û	“u”
368	Decoding latin small letter u with circumflex	û	“u”
369	Decoding latin small letter u with diaeresis	ü	“u”
370	Decoding latin small letter u with diaeresis	ü	“u”
371	Decoding latin small letter y with acute	ý	“y”
372	Decoding latin small letter y with acute	ý	“y”
373	Decoding latin small letter thorn	þ	“thorn”
374	Decoding latin small letter thorn	þ	“thorn”
375	Decoding latin small letter y with diaeresis	ÿ	“y”
376	Decoding latin small letter y with diaeresis	ÿ	“y”
Mathematical, Greek and Symbolic characters set:			
https://www.w3.org/TR/html4/sgml/entities.html http://www.w3schools.com/charsets/ref_html_symbols.asp http://www.w3schools.com/charsets/ref_html_entities_4.asp			
377	Decoding latin small f with hook	ƒ	“ “
378	Decoding latin small f with hook	ƒ	“ “
379	Decoding greek capital letter alpha	Α	“ “
380	Decoding greek capital letter alpha	Α	“ “
381	Decoding greek capital letter beta	Β	“ “
382	Decoding greek capital letter beta	Β	“ “
383	Decoding greek capital letter gamma	Γ	“ “
384	Decoding greek capital letter gamma	Γ	“ “
385	Decoding greek capital letter delta	Δ	“ “
386	Decoding greek capital letter delta	Δ	“ “
387	Decoding greek capital letter epsilon	Ε	“ “
388	Decoding greek capital letter epsilon	Ε	“ “
389	Decoding greek capital letter zeta	Ζ	“ “
390	Decoding greek capital letter zeta	Ζ	“ “
391	Decoding greek capital letter eta	Η	“ “
392	Decoding greek capital letter eta	Η	“ “
393	Decoding greek capital letter theta	Θ	“ “
394	Decoding greek capital letter theta	Θ	“ “
395	Decoding greek capital letter iota	Ι	“ “
396	Decoding greek capital letter iota	Ι	“ “
397	Decoding greek capital letter kappa	Κ	“ “
398	Decoding greek capital letter kappa	Κ	“ “
399	Decoding greek capital letter lambda	Λ	“ “
400	Decoding greek capital letter lambda	Λ	“ “
401	Decoding greek capital letter mu	Μ	“ “
402	Decoding greek capital letter mu	Μ	“ “
403	Decoding greek capital letter nu	Ν	“ “

404	Decoding greek capital letter nu	Ν	“ “
405	Decoding greek capital letter xi	Ξ	“ “
406	Decoding greek capital letter xi	Ξ	“ “
407	Decoding greek capital letter omicron	Ο	“ “
408	Decoding greek capital letter omicron	Ο	“ “
409	Decoding greek capital letter pi	Π	“ “
410	Decoding greek capital letter pi	Π	“ “
411	Decoding greek capital letter rho	Ρ	“ “
412	Decoding greek capital letter rho	Ρ	“ “
413	Decoding greek capital letter sigma	Σ	“ “
414	Decoding greek capital letter sigma	Σ	“ “
415	Decoding greek capital letter tau	Τ	“ “
416	Decoding greek capital letter tau	Τ	“ “
417	Decoding greek capital letter upsilon	Υ	“ “
418	Decoding greek capital letter upsilon	Υ	“ “
419	Decoding greek capital letter phi	Φ	“ “
420	Decoding greek capital letter phi	Φ	“ “
421	Decoding greek capital letter chi	Χ	“ “
422	Decoding greek capital letter chi	Χ	“ “
423	Decoding greek capital letter psi	Ψ	“ “
424	Decoding greek capital letter psi	Ψ	“ “
425	Decoding greek capital letter omega	Ω	“ “
426	Decoding greek capital letter omega	Ω	“ “
427	Decoding greek small letter alpha	α	“ “
428	Decoding greek small letter alpha	α	“ “
429	Decoding greek small letter beta	β	“ “
430	Decoding greek small letter beta	β	“ “
431	Decoding greek small letter gamma	γ	“ “
432	Decoding greek small letter gamma	γ	“ “
433	Decoding greek small letter delta	δ	“ “
434	Decoding greek small letter delta	δ	“ “
435	Decoding greek small letter epsilon	ε	“ “
436	Decoding greek small letter epsilon	ε	“ “
437	Decoding greek small letter zeta	ζ	“ “
438	Decoding greek small letter zeta	ζ	“ “
439	Decoding greek small letter eta	η	“ “
440	Decoding greek small letter eta	η	“ “
441	Decoding greek small letter theta	θ	“ “
442	Decoding greek small letter theta	θ	“ “
443	Decoding greek small letter iota	ι	“ “
444	Decoding greek small letter iota	ι	“ “
445	Decoding greek small letter kappa	κ	“ “
446	Decoding greek small letter kappa	κ	“ “

447	Decoding greek small letter lambda	λ	“ ”
448	Decoding greek small letter lambda	λ	“ ”
449	Decoding greek small letter mu	μ	“ ”
450	Decoding greek small letter mu	μ	“ ”
451	Decoding greek small letter nu	ν	“ ”
452	Decoding greek small letter nu	ν	“ ”
453	Decoding greek small letter xi	ξ	“ ”
454	Decoding greek small letter xi	ξ	“ ”
455	Decoding greek small letter omicron	ο	“ ”
456	Decoding greek small letter omicron	ο	“ ”
457	Decoding greek small letter pi	π	“ ”
458	Decoding greek small letter pi	π	“ ”
459	Decoding greek small letter rho	ρ	“ ”
460	Decoding greek small letter rho	ρ	“ ”
461	Decoding greek small letter final sigma	ς	“ ”
462	Decoding greek small letter final sigma	ς	“ ”
463	Decoding greek small letter sigma	σ	“ ”
464	Decoding greek small letter sigma	σ	“ ”
465	Decoding greek small letter tau	τ	“ ”
466	Decoding greek small letter tau	τ	“ ”
467	Decoding greek small letter upsilon	υ	“ ”
468	Decoding greek small letter upsilon	υ	“ ”
469	Decoding greek small letter phi	φ	“ ”
470	Decoding greek small letter phi	φ	“ ”
471	Decoding greek small letter chi	χ	“ ”
472	Decoding greek small letter chi	χ	“ ”
473	Decoding greek small letter psi	ψ	“ ”
474	Decoding greek small letter psi	ψ	“ ”
475	Decoding greek small letter omega	ω	“ ”
476	Decoding greek small letter omega	ω	“ ”
477	Decoding greek small letter theta symbol	ϑ	“ ”
478	Decoding greek small letter theta symbol	ϑ	“ ”
479	Decoding greek upsilon with hook symbol	ϒ	“ ”
480	Decoding greek upsilon with hook symbol	ϒ	“ ”
481	Decoding greek pi symbol	ϖ	“ ”
482	Decoding greek pi symbol	ϖ	“ ”
483	Decoding bullet	•	“ ”
484	Decoding bullet	•	“ ”
485	Decoding horizontal ellipsis	…	“ ” ...
486	Decoding horizontal ellipsis	…	“ ” ...
487	Decoding prime	′	“ ”
488	Decoding prime	′	“ ”
489	Decoding double prime	″	“ ”

490	Decoding double prime	″	“”
491	Decoding overline	‾	“ ”
492	Decoding overline	‾	“ ”
493	Decoding fraction slash	⁄	“/”
494	Decoding fraction slash	⁄	“/”
495	Decoding blackletter capital I	ℑ	“ ”
496	Decoding blackletter capital I	ℑ	“ ”
497	Decoding script capital P	℘	“ ”
498	Decoding script capital P	℘	“ ”
499	Decoding blackletter capital	ℜ	“ ”
500	Decoding blackletter capital	ℜ	“ ”
501	Decoding trade mark sign	™	[trade mark sign]
502	Decoding trade mark sign	™	[trade mark sign]
503	Decoding alef symbol	ℵ	“ ”
504	Decoding alef symbol	ℵ	“ ”
505	Decoding leftwards arrow	←	“ ”
506	Decoding leftwards arrow	←	“ ”
507	Decoding upwards arrow	↑	“ ”
508	Decoding upwards arrow	↑	“ ”
509	Decoding rightwards arrow	→	“ ”
510	Decoding rightwards arrow	→	“ ”
511	Decoding downwards arrow	↓	“ ”
512	Decoding downwards arrow	↓	“ ”
513	Decoding left right arrow	↔	“ ”
514	Decoding left right arrow	↔	“ ”
515	Decoding downwards arrow with corner leftwards	↵	“ ”
516	Decoding downwards arrow with corner leftwards	↵	“ ”
517	Decoding leftwards double arrow	⇐	“ ”
518	Decoding leftwards double arrow	⇐	“ ”
519	Decoding upwards double arrow	⇑	“ ”
520	Decoding upwards double arrow	⇑	“ ”
521	Decoding rightwards double arrow	⇒	“ ”
522	Decoding rightwards double arrow	⇒	“ ”
523	Decoding downwards double arrow	⇓	“ ”
524	Decoding downwards double arrow	⇓	“ ”
525	Decoding left right double arrow	⇔	“ ”
526	Decoding left right double arrow	⇔	“ ”
527	Decoding for all	∀	“ ”
528	Decoding for all	∀	“ ”
529	Decoding partial differential	∂	“ ”
530	Decoding partial differential	∂	“ ”
531	Decoding there exists	∃	“ ”
532	Decoding there exists	∃	“ ”

533	Decoding empty set	∅	“ ”
534	Decoding empty set	∅	“ ”
535	Decoding nabla	∇	“ ”
536	Decoding nabla	∇	“ ”
537	Decoding element of	∈	“ ”
538	Decoding element of	∈	“ ”
539	Decoding not an element of	∉	“ ”
540	Decoding not an element of	∉	“ ”
541	Decoding contains as member	∋	“ ”
542	Decoding contains as member	∋	“ ”
543	Decoding n-ary product	∏	“ ”
544	Decoding n-ary product	∏	“ ”
545	Decoding n-ary sumation	∑	“ ”
546	Decoding n-ary sumation	∑	“ ”
547	Decoding minus sign	−	“_”
548	Decoding minus sign	−	“_”
549	Decoding asterisk operator	∗	“*”
550	Decoding asterisk operator	∗	“*”
551	Decoding square root	√	“ ”
552	Decoding square root	√	“ ”
553	Decoding proportional to	∝	“ ”
554	Decoding proportional to	∝	“ ”
555	Decoding infinity	∞	“ ”
556	Decoding infinity	∞	“ ”
557	Decoding angle	∠	“ ”
558	Decoding angle	∠	“ ”
559	Decoding logical and	∧	“ ”
560	Decoding logical and	∧	“ ”
561	Decoding logical or	∨	“ ”
562	Decoding logical or	∨	“ ”
563	Decoding intersection	∩	“ ”
564	Decoding intersection	∩	“ ”
565	Decoding union	∪	“ ”
566	Decoding union	∪	“ ”
567	Decoding integral	∫	“ ”
568	Decoding integral	∫	“ ”
569	Decoding therefore	∴	“ ”
570	Decoding therefore	∴	“ ”
571	Decoding tilde operator	∼	“~”
572	Decoding tilde operator	∼	“~”
573	Decoding approximately equal to	≅	“ ”
574	Decoding approximately equal to	≅	“ ”
575	Decoding almost equal to	≈	“ ”

576	Decoding almost equal to	≈	“ ”
577	Decoding not equal to	≠	“ ”
578	Decoding not equal to	≠	“ ”
579	Decoding identical to	≡	“ ”
580	Decoding identical to	≡	“ ”
581	Decoding less-than or equal to	≤	“ ”
582	Decoding less-than or equal to	≤	“ ”
583	Decoding greater-than or equal to	≥	“ ”
584	Decoding greater-than or equal to	≥	“ ”
585	Decoding subset of	⊂	“ ”
586	Decoding subset of	⊂	“ ”
587	Decoding superset of	⊃	“ ”
588	Decoding superset of	⊃	“ ”
589	Decoding not a subset of	⊄	“ ”
590	Decoding not a subset of	⊅	“ ”
591	Decoding subset of or equal to	⊆	“ ”
592	Decoding subset of or equal to	⊆	“ ”
593	Decoding superset of or equal to	⊇	“ ”
594	Decoding superset of or equal to	⊇	“ ”
595	Decoding circled plus	⊕	“ ”
596	Decoding circled plus	⊕	“ ”
597	Decoding circled times	⊗	“ ”
598	Decoding circled times	⊗	“ ”
599	Decoding up tack	⊥	“ ”
600	Decoding up tack	⊥	“ ”
601	Decoding dot operator	⋅	“ ”
602	Decoding dot operator	⋅	“ ”
603	Decoding left ceiling	⌈	“ ”
604	Decoding left ceiling	⌈	“ ”
605	Decoding right ceiling	⌉	“ ”
606	Decoding right ceiling	⌉	“ ”
607	Decoding left floor	⌊	“ ”
608	Decoding left floor	⌊	“ ”
609	Decoding right floor	⌋	“ ”
610	Decoding right floor	⌋	“ ”
611	Decoding left-pointing angle bracket	〈	“ ”
612	Decoding left-pointing angle bracket	⟨	“ ”
613	Decoding right-pointing angle bracket	〉	“ ”
614	Decoding right-pointing angle bracket	⟩	“ ”
615	Decoding lozenge	◊	“ ”
616	Decoding lozenge	◊	“ ”
617	Decoding black spade suit	♠	“ ”
618	Decoding black spade suit	♠	“ ”

619	Decoding black club suit	♣	“ ”
620	Decoding black club suit	♣	“ ”
621	Decoding black heart suit	♥	“ ”
622	Decoding black heart suit	♥	“ ”
623	Decoding black diamond suit	♦	“ ”
624	Decoding black diamond suit	♦	“ ”
Special characters set: [https://www.w3.org/TR/html4/sgml/entities.html]			
625	Decoding latin capital ligature OE	Œ	“OE”
626	Decoding latin capital ligature OE	Œ	“OE”
627	Decoding latin small ligature oe	œ	“oe”
628	Decoding latin small ligature oe	œ	“oe”
629	Decoding latin capital letter S with caron	Š	“Š”
630	Decoding latin capital letter S with caron	Š	“Š”
631	Decoding latin small letter s with caron	š	“š”
632	Decoding latin small letter s with caron	š	“š”
633	Decoding latin capital letter Y with diaeresis	Ÿ	“Ÿ”
634	Decoding latin capital letter Y with diaeresis	Ÿ	“Ÿ”
635	Decoding modifier letter circumflex accent	ˆ	“ ”
636	Decoding modifier letter circumflex accent	ˆ	“ ”
637	Decoding small tilde	˜	“~”
638	Decoding small tilde	˜	“~”
639	Decoding en space	 	“ ”
640	Decoding en space	 	“ ”
641	Decoding em space	 	“ ”
642	Decoding em space	 	“ ”
643	Decoding thin space	 	“ ”
644	Decoding thin space	 	“ ”
645	Decoding zero width non-joiner	‌	“ ”
646	Decoding zero width non-joiner	‍	“ ”
647	Decoding zero width joiner	‍	“ ”
648	Decoding zero width joiner	‍	“ ”
649	Decoding left-to-right mark	‎	“ ”
650	Decoding left-to-right mark	‎	“ ”
651	Decoding right-to-left mark	‏	“ ”
652	Decoding right-to-left mark	‏	“ ”
653	Decoding en dash	–	“_”
654	Decoding en dash	–	“_”
655	Decoding em dash	—	“_”
656	Decoding em dash	—	“_”
657	Decoding left single quotation mark	‘	“'”
658	Decoding left single quotation mark	‘	“'”

659	Decoding right single quotation mark	’	“ ’ ”
660	Decoding right single quotation mark	’	“ ’ ”
661	Decoding single low-9 quotation mark	‚	“ ’ ”
662	Decoding single low-9 quotation mark	‚	“ ’ ”
663	Decoding left double quotation mark	“	“ ”
664	Decoding left double quotation mark	“	“ ”
665	Decoding right double quotation mark	”	“ ”
666	Decoding right double quotation mark	”	“ ”
667	Decoding double low-9 quotation mark	„	“ ”
668	Decoding double low-9 quotation mark	„	“ ”
669	Decoding dagger	†	“ ”
670	Decoding dagger	†	“ ”
671	Decoding double dagger	‡	“ ”
672	Decoding double dagger	‡	“ ”
673	Decoding per mille sign	‰	“ ”
674	Decoding per mille sign	‰	“ ”
675	Decoding single left-pointing angle quotation mark	‹	“ ’ ”
676	Decoding single left-pointing angle quotation mark	‹	“ ’ ”
677	Decoding single right-pointing angle quotation mark	›	“ ’ ”
678	Decoding single right-pointing angle quotation mark	›	“ ’ ”
679	Decoding euro sign	€	“[Euro]”
680	Decoding euro sign	€	“[Euro]”