

Reading the Disclosures with New Eyes: Bridging the Gap between Information Security Disclosures and Incidents[†]

Ta-Wei “David” Wang

Krannert Graduate School of Management
Purdue University
West Lafayette, IN 47907
wang131@purdue.edu

Jackie Rees

Krannert Graduate School of Management
Center for Education and Research in Information Assurance and Security (CERIAS)
Purdue University
West Lafayette, IN 47907
jrees@purdue.edu

Karthik Kannan

Krannert Graduate School of Management
Center for Education and Research in Information Assurance and Security (CERIAS)
Purdue University
West Lafayette, IN 47907
kkarthik@purdue.edu

Abstract

This paper investigates how the characteristics of information security incidents and disclosures in financial reports affect the valuation of a firm. Building on theories of disclosures in the accounting literature, we investigate investor reaction to disclosures through both quantitative and qualitative analyses. A cross-sectional analysis is first performed to examine the effect of the number of disclosures on stock price reactions to information security incidents. The results suggest that information security risk factors disclosed in financial reports increases the impact of information security incidents. Such an observation is consistent with investors perceiving the disclosures as a warning of future incidents. In order to provide a richer interpretation of the results, we further explore the contents of the disclosures using text mining techniques. One of the key findings is that breached firms react to information security incidents by disclosing additional and more specific risk factors. We further build a model to link disclosures with different stock price reactions to information security incidents to provide insights into how companies should disclose security concerns and practices. The model suggests that the disclosures associated with non-negative reactions are more generic and include actionable terms which confirms different disclosure patterns from companies with and without breach announcements. Thus, the paper not only contributes to the literature in information security and accounting but also sheds light on how managers can evaluate their information security policies and convey information security practices more effectively to the investors.

Keywords: information security, Sarbanes-Oxley Act (SOX), risk disclosure, text mining

[†] The authors are grateful to the Center for Education and Research in Information Assurance and Security (CERIAS) at Purdue University for funding part of the research.

1. Introduction

Organizations rely heavily on information technology (IT) to enable daily operations. Because of this dependency on IT, there may be a tremendous impact on business when there is an information security related incident. For example, a series of Denial of Service (DoS) attacks in 2000 resulted in online retailers and portals such as Amazon.com and Yahoo! losing service for hours (Sandoval and Wolverton 2000). According to the CSI/FBI computer crime and security report in 2006 (CSI/FBI 2007), the total dollar amount of financial losses resulting from security breaches is approximately \$200,000 US dollars per respondent. The losses of different attacks ranged from \$90,000 to \$15,000,000, accompanied by the fast growing number of reported security incidents (CERT 2007). This evidence highlights organizational concerns regarding information security. One way senior managers convey concerns about such potential disruptions is through financial report disclosures.

Disclosures in general are relevant to issues involving information asymmetry between a firm and its investors. In the accounting literature, two different motivations for disclosures are provided. On one hand, papers such as Verrecchia (1983), Dye (1985), and Verrecchia (2001) demonstrate that a firm only discloses information that is positively correlated to its business value. On the other hand, papers such as Skinner (1994), and Kasznik and Lev (1995) present evidence that a firm discloses in order to reduce its legal and reputation costs from the disappointing information it expects. It is not ex ante clear which specific motivation would be applicable to information security disclosures.

As to the first motivation, one may expect information security disclosures to indicate preparedness for security incidents. Therefore, the disclosures could have a positive impact on the valuation of the firm when an information security incident is observed. On the contrary, as

to the second motivation, disclosure itself implies future litigation or reputation costs, which result in a decreasing of future cash flow and also the valuation of the firm. Understanding which motivation is applicable should aid managers in deciding the extent of information security disclosures provided. If the first motivation holds, the managers should encourage disclosure. However, if the second motivation holds, the managers should be careful about how they convey their security practices to the public.

In light of this opponent conflict, we seek to answer the following research questions: Do information security disclosures in financial reports mitigate or worsen stock price reactions when a firm faces information security incidents?¹ What are the elements within these disclosures that have significant impacts on stock prices and characterize these disclosures? Do companies change their disclosure policies, such as the number or the focus of information security related risk factors, after experiencing information security incidents?

To answer these questions, we associate the information security incidents and stock price reactions to incidents, with the disclosures in financial reports. For the disclosures, we employ two different sources. One is the internal control report, which is mandated by Sarbanes-Oxley Act (SOX)² Section 404, describing the weaknesses of internal controls and financial systems. The other piece of information is the voluntary disclosure of risk factors that firms include regarding their future performance and forward-looking statements. Using the data, we perform cross-sectional analysis on the performance of the firm's stock price to various aspects of disclosures. Since how risk factors are disclosed in financial reports and the readability of financial reports can affect investors' expectations (Katz 2001; Li 2006), we also investigate the

¹ We focus on the disclosures of possible risk factors in financial reports instead of the disclosures of security breaches or sharing vulnerabilities. The latter has been addressed by, for example, Gal-Or and Ghose (2005), Gordon et al. (2005).

² H. Res. 107-414, 116 STAT. 745 Cong. Rec. 5395 (2002) (enacted). Retrieved on Apr. 9 2007, from <http://news.findlaw.com/hdocs/docs/gwbush/sarbanesoxley072302.pdf>

contents of risk factor disclosures using text mining techniques toward the end of the paper. Thus, our paper provides a comprehensive investigation involving both quantitative and qualitative analyses.

The rest of the paper is organized as the following. We first draw on the literature about information security and voluntary disclosure (section 2). The research framework and hypotheses are elaborated in section 3. In section 4, details of the cross-sectional analysis and the results are presented. Results from text mining are given in section 5. We conclude with discussion of contributions, limitations and venues for future research in section 6.

2. Literature Review

There are two major streams of literature that are directly related to our study. One is the research in business value and information security. The other is the literature on disclosures in accounting.

2.1 Information Security

A majority of information security literature focuses on technical issues but analytical and empirical studies in information security from an economic perspective are relatively limited. For instance, several studies have been done to address information security investments analytically (Gordon and Loeb 2002; Gordon et al. 2003). Studies have also pointed out that information security breaches can result in material impacts of business operation, including physical and intangible impacts such as negative company image and loss of reputation (Glover et al. 2001; Warren and Hutchinson 2000). Further, several empirical studies investigate the impact of information security events on business value. Based on different methodologies and different data sets, some of the results show that there exist significant negative impacts (Ettredge and Richardson 2003; Garg et al. 2003; Cavusoglu et al. 2004), while others do not

find such impact (Campbell et al. 2003; Hovav and D'Arcy 2003; Kannan et al. 2007). For example, Ettredge and Richardson (2003) investigate the impacts of the denial of service attacks which happened in February 2000 and attempt to determine which firm might suffer or benefit from similar incidents in the future. Their results demonstrate the existence of information transfer and show that the larger the firm, the larger the abnormal return. Kannan et al. (2007) also analyze short-term and long-term impacts of security announcements on market value and do not uncover a relationship between announcements and business value. Although our paper also considers security breach events, we focus on understanding the impact of information security disclosures.

2.2 Disclosures in Accounting

There is a rich body of literature examining voluntary disclosures in accounting. When there is no disclosure cost, full disclosure exists because investors believe that non-disclosing companies have the worst possible information (e.g. Grossman 1981; Milgrom 1981). However, if disclosure costs or uncertainty exist, companies will disclose only when the benefits exceed the costs (e.g. Verrecchia 1983; Dye 1985). The disclosure decision also depends on whether such disclosure will provide information to competitors and depends on mandatory disclosures (e.g. Verrecchia 1983; Darrough 1993; Eihorn 2005). Disclosure may also be used so as to reduce legal and reputation costs from bad news (Skinner 1994). Kasznik and Lev (1995) also find that firms are more probable to make a disclosure when the firm faces earnings disappointments. Specific to risk disclosures, one recent study by Jorgensen and Kirschenheiter (2003) has formally modeled managers' decisions on voluntarily disclosing a firm's risks. Furthermore, several empirical studies focus on the quality and credibility of the disclosures (e.g. Lang and Lundholm 1993; Penno 1997; Stocken 2000), the usefulness of disclosures (e.g.

Francis et al. 2002; Landsman and Maydew 2002), and other aspects of voluntary disclosures such as expectation adjustment, costs, analysts following, and signaling rationale (e.g. Ajinkya and Giff 1984; King et al. 1990; Elliott and Jacobson 1994; Lang and Lundholm 1996; Lev and Penman 1990).

In this paper, we link both the above streams of research. To the best of our knowledge, Sohail (2006) and Balakrishnan et al. (2008) are the only two studies that have also linked these two streams. Balakrishnan et al. (2008) focuses on the impact of SOX and whether the timeliness of information induced by SOX increase the quality of information disclosed to the market by investigating 8-K reports (important events not covered by previous annual or quarterly reports such as material disposition of assets or bankruptcy) and drawing relationship between the disclosure of 8-K reports and stock market reactions. However, our paper has a different focus. We focus on the relationship among risk factors disclosed in financial reports (10-K reports), information security incidents and stock price reactions to the incidents. In Sohail's paper, he demonstrates that security disclosures are positively related to stock price. However, his work solely focuses on disclosures but does not consider the relationship between the disclosures and subsequent information security incidents. By including the subsequent incidents, we are able to better understand how disclosures formulate investors' expectations and, in turn, affect the business value. In our paper, we not only analyze how the characteristics of information security incidents and disclosures in financial reports affect the valuation of a firm but also consider how investors react to disclosures and how firms can appropriately convey information security concerns or practices through disclosures.

3. Research Framework and Hypotheses Development

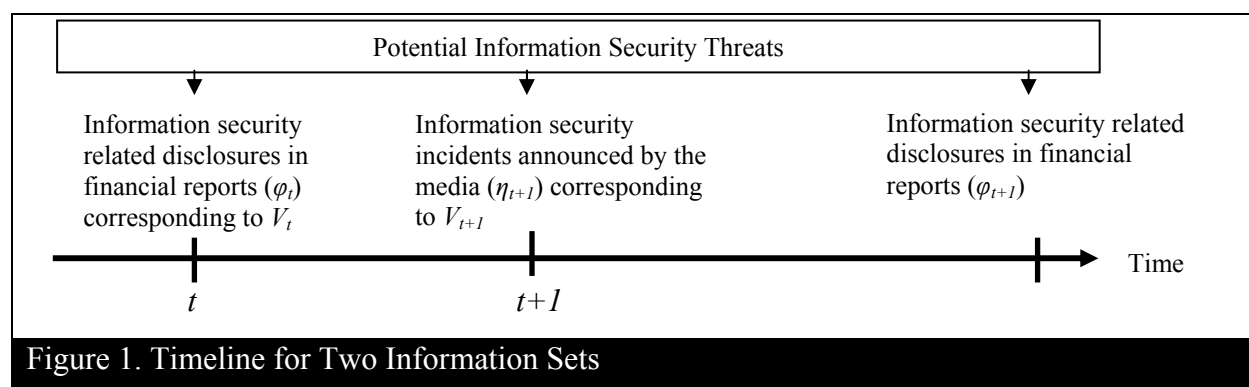
Based on the efficient market hypothesis (Fama 1970), a firm's business value at time t (V_t) can be expressed as the discounted value of expected future cash flows given all the available information until that time.

$$V_t = E \left\{ \sum_{i=t}^T \frac{x_i | \Phi_t}{\prod_{j=t}^i (1+r_j^t)} \right\} \quad (1)$$

In Equation (1), E is the expectation operator, T denotes the assumed terminal period which can be infinite, $x_i | \Phi_t$ is the net cash flow in period i given the information Φ_t available at time t , and r_j^t is the interest rate faced by the firm in period j at time t . Often, there is asymmetry in the information Φ_t available to the firm and its investors. Particularly, in this paper, the asymmetry is with respect to information security risks/threats the firm faces. The information security threats faced by the firm can be one of the following three types: (1) confidentiality, such as theft of source code or customer data, (2) integrity, such as a virus attack which deletes or alters files, or (3) availability, such as denial-of-service attacks (Bowen et al. 2006; Gordon et al. 2006). The threats can lead to both direct and indirect costs for the firm (Ettredge and Richardson 2003; Garg et al. 2003; Cavusoglu et al. 2004). The direct costs include the loss of productivity, the costs related to informing consumers, etc. The indirect costs include the loss of future transactions with consumers (and partners) that may be unwilling to trust the firm. Therefore, as with any other type of risk, the investors' uncertainty regarding the risks can negatively affect the expectation of the future cash flow and also the valuation of the firm. Given the uncertainty, each firm decides whether to disclose the threats to its future cash flows to the investors (Jorgensen and Kirschenheiter 2003).

In our information security context, we consider two different pieces of information (Φ_t in Equation (1)) regarding the threats a firm faces (the timeline is provided in Figure 1). The first

piece of information involves breach related information announced in the media and is represented by η_t . The second piece of information involves information security disclosures submitted by the firm in financial reports and represented by φ_t . Within the financial reports, information security related disclosures occur in two different places. The first is the disclosure of internal control and procedures mandated by Sarbanes-Oxley Act (SOX) section 404 (see Appendix A for an example).³ The second is the list of risk factors or possible uncertainties regarding forward-looking statements that may adversely affect a firm's future performance including information security related risk factors (see Appendix B for examples). Given the above characterization, we first investigate how φ_t and η_{t+1} affect the change in a firm's business value ($\Delta V=(V_{t+1} - V_t)|_{\eta_{t+1} \neq \phi}$). Then we analyze the nature of the textual content in φ_t and φ_{t+1} given that η_{t+1} is not null.



Now, let us focus on investigating which of the motivations discussed in the introduction is valid. The underlying rationale for the first motivation is that when the disclosure benefits outweigh the costs, a firm is willing to disclose so as to increase its business value (Verrecchia 1983; Dye 1985). One particular concern regarding information security disclosure is that it can

³ According to SOX Section 404 (see footnote 1), the internal control report should: “(1) state the responsibility of management for establishing and maintaining an adequate internal control structure and procedures for financial reporting, and (2) contain an assessment, as of the end of the most recent fiscal year of the issuer, of the effectiveness of the internal control structure and procedures of the issuer for financial reporting.”

expose the firm to the specific type of risks mentioned in the disclosure resulting in industrial espionage, loss of reputation and/or loss of competitive advantage (Gordon et al. 2005). Despite these concerns, firms disclose information security risk factors in their financial reports. This leads us to believe that, perhaps, the reason is that the accompanied litigation and reputation costs with information security incidents are even larger which is the second motivation presented in Section 1 (Skinner 1994; Kasznik and Lev 1995). Therefore, the impact on the business value will be larger if the disclosed threats are realized implying increased litigation and reputation costs lowering future cash flows and, in turn, the business value. Accordingly, we have our first hypothesis:

Hypothesis 1: For breached firms, as the number of internal control related items disclosed in the section of “Control and Procedures” and/or the number of disclosures of information security related risk factors (ϕ_i) increase, the impact of information security incidents on stock prices (ΔV) increases.

Hypothesis 1 plays an important role in the paper. It serves not only as the foundation for the cross-sectional analyses but also as the basis for the exploration of the contents within the disclosures in Section 5.

When hypothesizing about the impact of disclosure at the aggregate level, Hypothesis 1 fails to account for three other issues which are discussed below. The first issue relates to the realization of the expectations. Prior literature has investigated the investors' reaction to the realization of the expectations. For example, Begley and Fischer (1998), Bagnoli et al. (2002) study the investor reaction to whether a firm meets or misses the expected earnings report date. Similarly, Kasznik and McNichols (2002) study the reaction to realization, the so-called “meet or miss” earnings expectations. That is, whether the realization of an event meets investors'

expectations built from disclosures can result in different stock price reactions. In our context, meeting expectations refers to the realization of the actual warning, i.e. information security incidents. Therefore, we suspect that the “match” between disclosed risk factors and incidents is an important explanation to our argument in Hypothesis 1. Therefore, we present Hypothesis 2.

Hypothesis 2: For breached firms, if the information security risk factors disclosed by a firm in financial reports (φ_t) match the incident the firm suffers (η_t), as the number of disclosures of information security related risk factors (φ_t) increases, the impact of information security incidents on stock prices (ΔV) increases.

The second issue not covered in Hypothesis 1 is about new signals. The literature that our argument builds on only focuses on the effect of first time disclosures or new signals (e.g. Verrecchia 1983; Dye 1985; Verrecchia 2001). From Hypothesis 1, if the disclosures were not covered in previous years’ financial reports, the information set is new and will decrease the valuation of the firm. However, the disclosures in annual reports may not provide additional new information since the same content could have already been disclosed in previous years’ reports, i.e. φ_{t-1} and φ_t are the same. Thus, we further distinguish first-time and repeated disclosures in Hypothesis 3.

Hypothesis 3: If information security related risk factors stated in current year’s financial report (φ_t) have not already been disclosed in previous years’ financial report(s) (φ_{t-j}), as the number of disclosed information security related risk factors (φ_t) increases, the impact of information security incidents on stock prices (ΔV) increases.

The third issue that has not been considered in Hypothesis 1 relates to how expectations are formed. As shown by Katz (2001), how these risks are disclosed affects the formation of

expectations. In order to understand how these risk factors help investors adjust expectations, we perform cluster analysis on the disclosures using text mining techniques in Section 5. Text mining is the technique used to extract information from text through finding nontrivial patterns and trends (Tan 1999; Feldman and Sanger 2006). It has been widely used in different contexts, such as fraud detection, drug design, or customer support (Cecchini et al. 2007; Fan et al. 2006; Han et al. 2002). By applying text mining techniques on the contents of the risk factor disclosures, we are able to identify and categorize the risk factors that are specifically related to information security. The contents might also help us better explain the phenomenon under investigation. Moreover, how companies react to information security incidents via disclosures can further affect investors' expectations in subsequent periods. Therefore, we also compare the disclosures *before* (φ_t) and *after* (φ_{t+1}) the incident for breached companies to address this point. We discuss this issue also in Section 5.

In the following section, we test our hypotheses. Based on the results, we further investigate the disclosures in detail through text mining in order to provide additional explanations. We also build a tool that can aid managers in designing appropriate disclosures.

4. Cross-Sectional Analysis

In order to test our hypotheses empirically, we first identify information security incidents. For the firms experiencing the incidents, we extract information security related disclosures from financial reports, and the associated stock prices. Based on the data collected, we investigate the relationship between stock price reactions and the disclosures in financial reports.

4.1 Sample Selection

To identify incidents, we search for news articles with the following keywords from 1997 to 2006 in the *Wall Street Journal*, *USA Today*, the *Washington Post*, and the *New York Times* via

the Factiva database as well as in *CNet* and *ZDNet*. The keywords are: (1) security breach, (2) hacker, (3) cyber attack, (4) virus or worm, (5) computer break-in, (6) computer attack, (7) computer security, (8) network intrusion, (9) data theft, (10) identity theft, (11) phishing, and (12) cyber fraud. These keywords are similar to those used in prior studies (e.g. Campbell et al. 2003; Garg et al. 2003; Kannan et al. 2007). Only the samples with the following properties are retained in our dataset. First, the articles must enable us to identify a specific date of the security incident announcement. Second, only publicly traded firms are included in the analysis/sample. Last, annual reports (10-K or 20-F reports) or quarterly reports (10-Q reports, when annual reports are not available) of the sample firms must be available one period before and after the event from EDGAR Online⁴. While searching for the articles, we check whether the companies release such information on their websites before the media does, in which case we need to appropriately distinguish self-reported incidents and account for the timing of release of the information. In our sample, we do not encounter any self-disclosed information directly from the companies. The resulting sample consists of 106 firm-event observations. These breached firms are referred to the experimental group in the rest of the paper. Within our sample, there are 37, 31 and 44 incidents of confidentiality, integrity, and availability type incidents, respectively.⁵ Interestingly, our sample demonstrates that there seems to be a shift from integrity and availability type incidents to confidentiality type after 2003. The shift might result from the growing value of data for criminals, i.e. identity theft.

⁴ <http://www.sec.gov/edgar/searchedgar/webusers.htm>

⁵ Six observations are grouped into both the integrity type and the availability type. For example, the “I love you” worm not only destroys files (integrity type) but also slows down mail server systems (availability type). Furthermore, two raters performed the coding task. We do not consider these six observations for reliability since these six observations were dropped from the analysis for the types of incidents later. Given the high inter-rater reliability (Cohen’s $\kappa = 92.83\%$), we adopt the author’s coding results for later analysis.

Based on the selection process discussed above, we collect the following data: (1) Company information: company name, industry identification code (SIC code), and CUSIP number for each firm’s stock. (2) Security incident information: source of the news, date, and the article. (3) Financial reports information: which report (10-K, 20-F or 10-Q) and which year of the report we use, disclosures in the “Control and Procedures” section, and information security related and all risk disclosures. The descriptive statistics of our sample is provided in Table 1. It shows that, on average, there is a greater number of security related disclosure and total number of risk factors disclosed per firm-event observation after SOX (year 2002). Table 1 also demonstrates that the number of disclosed risk factors increases after information security incidents.

Table 1. Descriptive Statistics of Disclosures				
Risk Factor Disclosures <i>before</i> Incidents	Number of Security Related Disclosures		Total Number of Risk Factors Disclosed	
	before 2002 ^a	after 2002	before 2002	after 2002
Total	27	20	888	616
Average (stdev)	0.52 (1.057)	0.56 (1.027)	17.08 (9.083)	17.11 (10.353)
Max (min)	4 (0)	4 (0)	38 (5)	43 (0)
Risk Factor Disclosures <i>after</i> Incidents	Number of Security Related Disclosures		Total Number of Risk Factors Disclosed	
	before 2002	after 2002	before 2002	after 2002
Total	39	56	837	856
Average (stdev)	0.64 (1.304)	1.91 (1.469)	13.72 (9.830)	18.21 (11.205)
Max (min)	4 (0)	5 (0)	38 (0)	45 (0)

^a SOX was enacted in 2002

The number of information security related risk factors and the total number of risk factors are calculated as follows. We count how many factors (both the total number of risk factors and information security related risk factors) mentioned by the firm in annual reports or quarterly reports under the section of risk factors or the section of forward-looking statements.⁶ This measurement of disclosure frequency has long been used and discussed in the accounting

⁶ Again since two raters’ inter-rater reliability is high for both groups (Cohen’s $\kappa = 97.23\%$ and 100% respectively), the author’s coding results are used.

literature (e.g. Francis et al. 1994; Lang and Lundholm 2000; Jo and Kim 2007). What we have done can be illustrated as follows. For instance, as shown in Appendix B, one risk factor disclosed by Amazon in year 2000⁷ was “We face intense competition”. The other was “System interruption and the lack of integration and redundancy in our systems may affect our sales”. Thus, after looking into the content of the disclosures, we count one for information security related risk factors and two for total risk factors in this case. Since firms generally group several elements with similar consequences in one risk factor, we posit that investors also take these elements as a single factor and evaluate the impacts.

4.2 Regression Models

The impact of economic events on business value can be measured by the stock price reactions in a short time period according to the theory of market efficiency (Fama 1970; MacKinlay 1997). To capture the impact of security incidents on stock price, we apply the market model (which is described in detail in Appendix C) and obtained the cumulative abnormal return (CAR). To properly measure the impact of security incidents, samples with confounding events, such as earnings announcements, merger and acquisition, and stock splits, are first eliminated so as to avoid possible other causes to the stock price reaction. The resulting sample size is 88 firm-event observations for the experimental group. We examine the two-day period (window) around the event date (the date of announcement, denote as day 0) for the stock price reactions, i.e. -1~0, where -1 represents 1 day *before* the event date. The results for the window -1~+1 is similar to those for -1~0.

In order to test our hypotheses, we formulate a cross-sectional analysis by regressing CAR on the number of disclosures in financial reports by controlling the size effect. Since previous

⁷<http://www.sec.gov/Archives/edgar/data/1018724/000103221001500087/0001032210-01-500087.txt>

studies have shown that large firms are more able to endure shocks than small firms and invest more in security (Fama and French 1992; PriceWaterhouseCoopers 2002), we take the logarithm of a firm's net assets (denoted as *Size*) in the corresponding accounting period to control for the size effect.

Based on our discussion of Hypothesis 1, we first investigate the effect of the disclosures of internal control and procedures as well as information security risk factors (see Table 2 for a list of variables). We consider the disclosures of internal control and procedures by three different elements, i.e. how a firm evaluates its internal controls and procedures (*ConP₁*), how a firm manages its internal controls and procedures (*ConP₂*), and whether a firm changes its internal controls and procedures (*ConP₃*) in Equation (2), and by total number of elements (*ConP*) a firm discloses in the control report in Equation (3). For the disclosures of information security risk factors in Equations (2) and (3), *Sec* denotes the number of information security related risk factor and *Trisk* represents the total number of risk factors. From Hypothesis 1, we expect β_2 , β_3 , β_4 , and β_5 in Equation (2), and β_2 as well as β_3 in Equation (3) to be negative for the experimental group.

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 ConP_{1i} + \beta_3 ConP_{2i} + \beta_4 ConP_{3i} + \beta_5 Sec_i + \beta_6 Trisk_i + \varepsilon_i \quad (2)$$

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 ConP_i + \beta_3 Sec_i + \beta_4 Trisk_i + \varepsilon_i \quad (3)$$

Next, we investigate how the other characteristics of the disclosure and the media announcements affect the CAR. Specifically, for Hypothesis 2, we investigate whether the disclosure coincides with the media announcement affects the CAR in Equations (4) and (5). We introduce two measures. *MSec* represents the number of matched information security risk factors and *PSec* is the percentage of matched security risk factors, i.e. *MSec* divided by *Sec*. Based on Hypothesis 2, we expect β_2 in Equations (4) and (5) to be negative.

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 MSec_i + \beta_3 Trisk_i + \varepsilon_i \quad (4)$$

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 PSec_i + \beta_3 Trisk_i + \varepsilon_i \quad (5)$$

Table 2. List of Variables	
CAR	Cumulative abnormal return (define in Appendix A)
Size	Firm size which equals to the logarithm of net assets.
ConP	The number of elements a firm disclosed in the section of Internal Control report. There are three possible elements (ConP ₁ , ConP ₂ , and ConP ₃) which are explained below
ConP ₁	Dummy variable for whether a firm discloses how it evaluates its internal controls and procedures. 1 if disclose, 0 otherwise.
ConP ₂	Dummy variable for whether a firm discloses how it manages its internal controls and procedures. 1 if disclose, 0 otherwise.
ConP ₃	Dummy variable for whether a firm discloses if there is a change in its internal controls and procedures. 1 if disclose, 0 otherwise.
Sec	Number of information security risk factors disclosed in financial reports.
Trisk	Total number of risk factors disclosed in financial reports.
DConf	Dummy variable for confidentiality type incidents. 1 if the incident is a confidentiality type incident, 0 otherwise.
DInt	Dummy variable for integrity type incidents. 1 if the incident is an integrity type incident, 0 otherwise. (if DConf and DInt equal to 0, the incident is an availability type incident)
MSec	A subset of Sec. Number of matched disclosures.
PSec	A subset of Sec. Defined as MSec divided by Sec, i.e. the level of matched disclosures
NSec	A subset of Sec. Number of first time information security related risk factors disclosed in financial reports
ε	Residual term

The second characteristic we investigate is that whether the mitigation effect is limited to first time disclosures (*NSec*) through Equation (6) for Hypothesis 3. We trace all our information security related risk disclosures back to 1997⁸ (or as early as possible) and determine whether these pieces of information have been disclosed in previous years' financial reports. Based on Hypothesis 3, we expect β_2 to be negative in equation (6).

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 NSec_i + \beta_3 Trisk_i + \varepsilon_i \quad (6)$$

4.3 Results

The relationships between the disclosures of internal control and procedures as well as the disclosures of information security risk factors and CAR are given in Table 3. From the results

⁸ We pick the year 1997 because 1997 is roughly the year that the Internet starts to get popular. Before 1997, we can hardly get information security related disclosures. Furthermore, 1997 is the earliest year that we can retrieve financial reports from EDGAR Online for most of our observations.

for Model 2 and Model 3 in Table 3, we show that the impact of the disclosures of internal control and procedures is not significant (the coefficients for $ConP_1$ to $ConP_3$ and $ConP$ are not significant). But the number of security related risk disclosures (Sec) significantly and negatively affects CAR. Therefore, our first hypothesis is partially supported. This finding supports our hypothesis that the greater the number of information security related risk disclosures, the larger the impact of information security incidents on stock prices. That is, the investors actually treat these disclosures as possible bad news that might happen in the future. The investors believe that managers have some private information about the weakness which might be used against the firm. Once the firm faces actual information security incidents, stock price reaction will be larger.

Table 3. Results for The Disclosure of Internal Control and Procedures Using Full Sample					
Variables	Model 2	Model 3	Model 4	Model 5	Model 6
Intercept	-0.05	-0.04	-0.03	-0.04	-0.04
Size	0.00	0.00	0.00	0.00	0.00
ConP ₁	0.03				
ConP ₂	-0.01				
ConP ₃	-0.02				
ConP		0.00			
Sec	-0.02***	-0.02***			
MSec			-0.05***		
PSec				-0.08***	
NSec					-0.02***
Trisk	0.00	0.00	0.00	0.00	0.00

* significant at 10% ** significant at 5% ***significant at 1%

So far, we have not yet considered the effect when the disclosed information security risk factors coincide with the incidents. For Hypothesis 2, the results are given in Table 3 (Model 4 and Model 5). Our finding supports Hypothesis 2 that whether the incident matches the disclosed risk factor and the level of match ($MSec$ and $PSec$) significantly affect stock price reactions. As shown by the magnitude of the coefficients, investors react more because of the match. This finding suggests that once the disclosed risk factor really happens, investors have concerns not only about that specific factor but also about whether other risk factors will also

realize. This result is further investigated in Section 5. In the previous analyses, we treated all the disclosures as first time disclosures. That is, we did not consider whether the disclosures have been disclosed in previous year's reports. Next, for Hypothesis 3, we only investigate the effect if the disclosures are first time disclosures as shown in Table 3 (Model 6). The result demonstrates a significantly negative relationship if the disclosures have not been announced before (*NSec*) which supports our third hypothesis.

In summary, the results above demonstrate that investors adjust their expectation about a firm's future profitability using the firm's disclosure about the uncertainty of information security events. By taking the disclosed information into account, the firm's stock price reactions are larger because the warnings of bad news are realized. However, since the disclosures of internal control and procedures are relatively standardized and have less "real" information, these disclosures do not provide useful information for investors.

4.4 Robustness Tests

In order to rule out possible explanations to our results, we perform the following robustness tests. First, in order to argue that our results are not prevailing to all the companies, for every breached firm, we find one of its publically trading competitors that does not have any breach announcements. We gather this information from Yahoo! Finance and the Hoover's Database. If several competitors can be selected, we choose the one with similar market capitalization and with financial reports available. We perform the same analysis using the controlled firms and do not find any significant results.

Second, we consider different windows for CAR and different subsamples. For different windows, we checked the following four different periods: (1)-30~-1, (2)-30~+1, (3)-7~+1, and (4)+1~+30, but we do not find any significant results. Also, since firms started to disclose

controls and procedures after SOX, we also investigate the subsample *after* 2002 (including 2002) (subsample 1) to estimate the impact in order to rule out possible systematic differences across time. However, we do not observe any significant results for subsample 1 either.

Third, Lev and Pennman (1990) have shown that firms in the same industry might have similar disclosure policies. Therefore, we also accounted for this industry effect as a robustness check. However, our results are not sensitive to the industry effect. Furthermore, the relationship between information security disclosures and incidents can also be affected by the types of the incidents. Thus, for Hypothesis 1, the effects of different types of incidents are also investigated in Equation (7), where *DConf* stands for confidentiality type incidents and *DInt* represents integrity type incidents. When *DConf* and *DInt* are both 0, the incident is an availability type incident. However, our results remain the same.

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Sec_i + \beta_3 Trisk_i + \beta_4 DConf_i + \beta_5 DInt_i + \varepsilon_i \quad (7)$$

Last, in Equations (2) and (3), we do not account for any interaction term. According to the accounting literature (e.g. Verrecchia 2001; Dye 2001; Bagnoli and Watts 2007), mandatory and voluntary disclosures are substitutes if a manager's private information is correlated to the firm's liquidation value. Voluntary risk disclosures provide supplemental information to mandatory disclosures if "the probability of voluntary risk disclosure is decreasing in the mandated amount of risk disclosures" (Bagnoli and Watts 2007, p.904). That is, if the mandatory disclosure can provide the information a firm wants to disclose, given the same disclosure between mandatory and voluntary disclosures, the firm will not voluntarily disclose additional information since the costs outweigh the benefits. However, given our observations of the number of voluntary risk disclosures and the purpose of risk factor disclosures on conveying uncertainties of future cash flows, these two information security related disclosures are neither substitutes nor

complementary to each other. In order to further confirm our argument about the interaction relationship between the disclosures of internal control and procedures and information security risk factors, we also examine their relationships with CAR separately via Equations (8) to (10). Based on our argument, the results from Equation (8) to (10) should be similar to the ones in Equations (2) and (3). Thus, we expect β_2 , β_3 and β_4 in equation (8), and β_2 in Equations (9) and (10) to be negative.

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 ConP_{1i} + \beta_3 ConP_{2i} + \beta_4 ConP_{3i} + \varepsilon_i \quad (8)$$

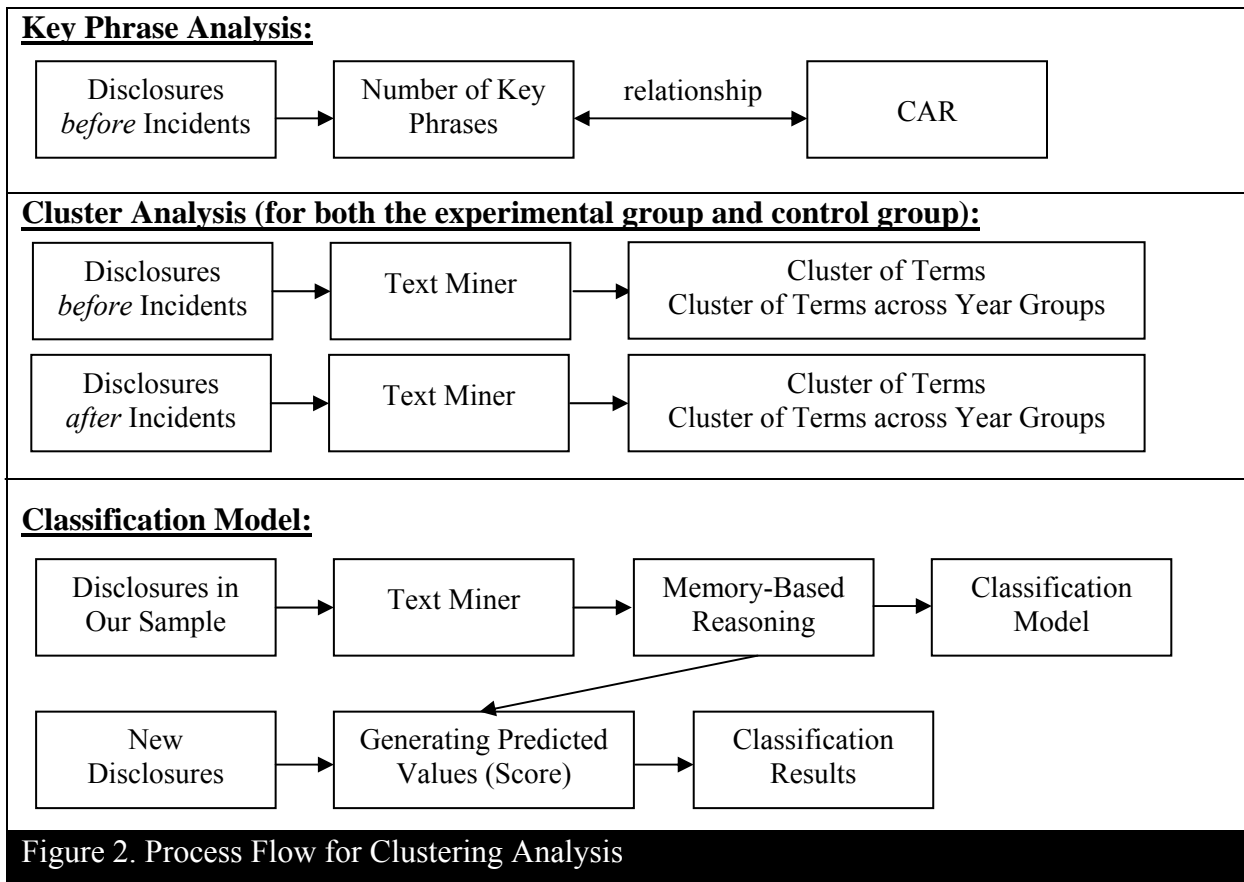
$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 ConP_i + \varepsilon_i \quad (9)$$

$$CAR_i = \beta_0 + \beta_1 Size_i + \beta_2 Sec_i + \beta_3 Trisk_i + \varepsilon_i \quad (10)$$

As expected, the results for Equations (8) to (10) are similar to those from Equations (2) and (3).

5. Text Mining

The analysis thus far has focused only on the quantitative attributes of the disclosure. Having developed an understanding of the reaction of the investors and the motivation for the reaction, we focus on textual mining of the data to further understand how information security risk factors form expectations. Specifically, we first perform a content analysis to link the textual data to stock price reactions. Then we perform cluster analysis on the disclosures using standard text mining techniques (see Appendix D). Then, a classification model is built to classify disclosures based on the disclosures in financial reports. The process flows for each of these three analyses are shown in Figure 2.



5.1 Key Phrase Analysis

This key phrases analysis section is a means to transition from cross-section analysis to the text mining analysis. The results presented in this section (Section 5.1) are still preliminary. Note that, until now, the analysis focused only on quantitative measures of the disclosures (such as the number of disclosures). While such an analysis is useful, we are interested in further delving the issue and analyzing the characteristics of textual content disclosed. For that, we first perform a content analysis and extract key phrases from information security risk factors, and regress CAR on the number of these key phrases. This analysis is also another basis for text mining analysis in Section 5. The key phrases are those explicitly mentioned in the disclosures as threats or the major factors that result in the disclosed impacts. The terms are “break-in”, “denial-of-service attacks”, “hacker”, “sabotage”, “secure transmission”, “system failure,

interruption or disruption”, “unauthorized access”, “virus/worm”, and other general security or breach statements. We also count how many times these key phrases show up in disclosures. Based on the coding scheme we build to perform this analysis, we asked another independent rater to code the same disclosures in order to check the reliability of our coding results. Since two raters’ inter-rater reliability is high (Cohen’s $\kappa = 89.14\%$), the author’s coding results are used.

The only significant result we observe is that CAR is positively associated with the number of times “break-in” is disclosed in the risk factors (coefficient = 0.04, $p = 0.09$) while the association is negative with the number of times “system failure, interruption or disruption” (coefficient = -0.03, $p = 0.05$) appears. This may perhaps so because statements about “break-in” do not provide any specific information for investors, it helps investors understand future uncertainties instead to cover possible future incidents. But when firms discuss “system failure, interruption or disruption” in the disclosures, the disclosure conveys some information that let investors believe the managers have private information about future incidents. It seems that different disclosures will result in different investors’ perceptions. We further explore this issue using text mining tools in Section 5.

In addition to the above analysis, we also examine whether the breached firms change their disclosures *after* experiencing information security incidents in terms of the number of these key phrases. The results show that the number of general security statements, the number of disclosures about virus and worms, and the number of disclosures about DoS attacks decreases. Instead, the firms disclose additional specific other factors, such as unauthorized access, cyber attack, and spam attack. It seems that firms adjust their disclosures not only to reflect emerging

risk factors but also to cover the incidents they have experienced which will also be further explored in Section 5.

5.2 Cluster Analysis

In this section, we further explore the contents of the disclosures and how the contents explain our findings in previous sections using text mining techniques. We first perform a cluster analysis on the disclosures before incidents. Then we investigate whether these clusters change across time and whether these clusters change after the firm experiences incidents. Last, we use concept links to explain the observations from the clusters.

Table 4 demonstrates the clustering results on the contents of the disclosures under three different levels of reduced dimensions, i.e. high, medium, and low resolutions in the SAS output. When performing the cluster analysis, a high resolution refers to use all the information while a low resolution uses only two-thirds of the information (SAS Institute Inc 2004). Further, a high resolution generally summarizes data better but might include more noise at the same time (SAS Institute Inc 2004). By changing different levels of information we use to generate clusters, we are able to see whether the computed cluster is consistent.

In Table 4, each row represents one cluster. Within each cluster, there are five terms. A term with the plus (+) sign represents a group of equivalent terms. For example, both “ability” and “abilities” are considered equivalent. The percentage is the frequency of a set of terms divided by the total frequency. The root mean squared standard deviation (RMS Std.) for cluster k equals to $\sqrt{W_k/[d(N_k - 1)]}$, where W_k is the sum of the squared distances from the cluster mean to each of the N_k documents in cluster k , and d is the number of dimensions. Further, in Table 4, only those clusters with frequency percentage over 15% are presented. As expected, since these are the disclosures of information security risk factors, from Table 4, we can see

terms with negative meanings such as “harm” and “failure” (in a gray area), and the subjects that may be affected, such as “system” and “customer” (in a gray area). From the frequency percentage, it seems that the breached firms disclose more specific factors than the firms in the control group do. Furthermore, the terms for the control firms are more about operation and actions such as “implement”, “protect”, and “require” (in a gray area). We will further explore the clusters using concept links, which will be defined later, to provide meanings to the terms.

Table 4. Text Mining Results of Information Security Related Risk Factors before Incidents ^a			
Terms	Freq.	Percentage	RMS Std.
Experimental Group			
Resolution: High			
+ attack, + harm, + have, other, + type	9	17%	0.152
+ event, + failure, operating results, + result, + site	9	17%	0.131
+ continue, + experience, increased, infrastructure, traffic	8	15%	0.133
Resolution: Medium			
+ event, + failure, + site, + system, web	9	17%	0.178
+ customer, + product, + protect, + revenue, software	9	17%	0.163
+ breach, confidential information, public networks, secure transmission, transmission	9	17%	0.176
Resolution: Low			
+ attack, + experience, + harm, + number, other	20	37%	0.163
Control Group			
Resolution: High			
+ depend, + failure, + interrupt, + interruption, + system	11	42%	0.220
+ implement, other, + protect, + require, + transaction	10	38%	0.207
+ impact, information, not, + process, 's	5	19%	0.237
Resolution: Medium			
+ implement, + protect, + require, + transaction, + transmission	7	27%	0.04
+ affect, computer systems, + failure, + result, + system	6	23%	0.213
+ control, + employee, + failure, + process, potential	5	19%	0.212
+ customer, + damage, + harm, + interrupt, + product	4	15%	0.113
+ depend, + interrupt, + interruption, power loss, + result	4	15%	0.000
Resolution: Low			
+ affect, + failure, financial, + security, + system	14	54%	0.254
+ depend, + harm, + interrupt, + result, + system,	6	23%	0.069
+ implement, + protect, + require, + transaction, + transmission	6	23%	0.015

^a The descriptive terms over the percentage of 15% are reported.

By performing the cluster analysis on different year groups, we further investigate whether the text mining results vary from year to year. Since we do not have enough observations for all the year groups for the control group, the results are only for the experimental group. Table 5 shows that firms start to consider “customers” and “users” in the disclosures after 2003. It seems that different types of incidents might be considered across time.

Table 5. Text Mining Results of Information Security Risk Factors by Year Groups before Incidents^a

Year Group	Terms	Freq.	Percentage	RMS Std.
Resolution: High				
1997-1998	N/A ^b			
1999-2000	critical, + event, + site, + system, web	8	32%	0.257
2001-2002	N/A			
2003-2004	+affect, +business, +customer, security, service	9	64%	0.289
2005-2006	+computer, +network, other, +protect, significant	6	75%	0.067
Resolution: Medium				
1997-1998	N/A			
1999-2000	data, + delay, + experience, + harm, + virus	7	28%	0.113
	+ breach, confidential, public networks, secure, transmission	7	28%	0.071
2001-2002	N/A			
2003-2004	+disruption, +event, +experience, +result, +user	8	57%	0.354
2005-2006	+computer, +network, other, +protect, significant	6	75%	0.067
Resolution: Low				
1997-1998	N/A			
1999-2000	data, + delay, + experience, similar, + virus	11	44%	0.098
2001-2002	N/A			
2003-2004	+affect, +business, +customer, +harm, information	7	50%	0.357
	+cause, +event, +experience, +user, +virus	7	50%	0.208
2005-2006	+computer, +network, other, +protect, significant	6	75%	0.067

^a The descriptive terms with the largest frequency are reported for the experimental group.
^b There are only one document in 1997-1998 and six documents in 2001-2002. The Text Miner cannot converge.

As discussed above, we also attempt to understand whether firms change their disclosure policies about information security risk factors after information security incidents. When we compare the clusters *before* the incidents (Table 4) and the clusters *after* the incidents (Table 6), we see similar terms such as “harm” and failure” as well as new terms such as “reputation”.

Table 7 shows, for experimental group, the cluster with the largest frequency percentage for different year ranges, by considering the disclosures after information security incidents. It seems that the terms “data” and “customer” starts to be considered in disclosures after 2003. This evidence confirms that there is a shift from availability type incidents to integrity and confidentiality type ones. This evidence also shows that firms might learn about new risk factors from the environment and disclose them in the financial reports.

Table 6. Text Mining Results of Information Security Risk Factors after Incidents ^a			
Terms	Freq.	Percentage	RMS Std.
Experimental Group			
Resolution: High			
+demand, infrastructure, internet, not, +product	10	15%	0.137
+continue, +depend, +experience, internet, +provide	10	15%	0.131
Resolution: Medium			
+business, information, not, security, +service	29	45%	0.177
+computer, +experience, +failure, +interruption, +result	16	25%	0.171
+disruption, +interruption, +loss, +telecommunication, +system	15	23%	0.164
Resolution: Low			
+depend, +experience, internet, not, +user	18	28%	0.173
+disruption, +event, +interruption, +loss, +system	13	20%	0.171
+breach, information, prevent, reputation, security	12	18%	0.156
Control Group			
Resolution: High			
+customer, +depend, +harm, +interrupt, +system	10	24%	0.183
+demand, internet, +maintain, +number, +provide	10	24%	0.173
+access, +implement, +require, unauthorized, unauthorized access	8	19%	0.161
+affect, +breach, +computer, +customer, +have	7	17%	0.193
adversely, +affect, +business, +failure, +loss	7	17%	0.191
Resolution: Medium			
ability, +affect, +damage, +failure, financial	10	24%	0.204
+computer, +customer, +event, +failure, +have	9	21%	0.152
Resolution: Low			
+affect, +business, +computer, +failure, +have	14	33%	0.124
+access, +breach, +implement, +require, unauthorized	10	24%	0.130

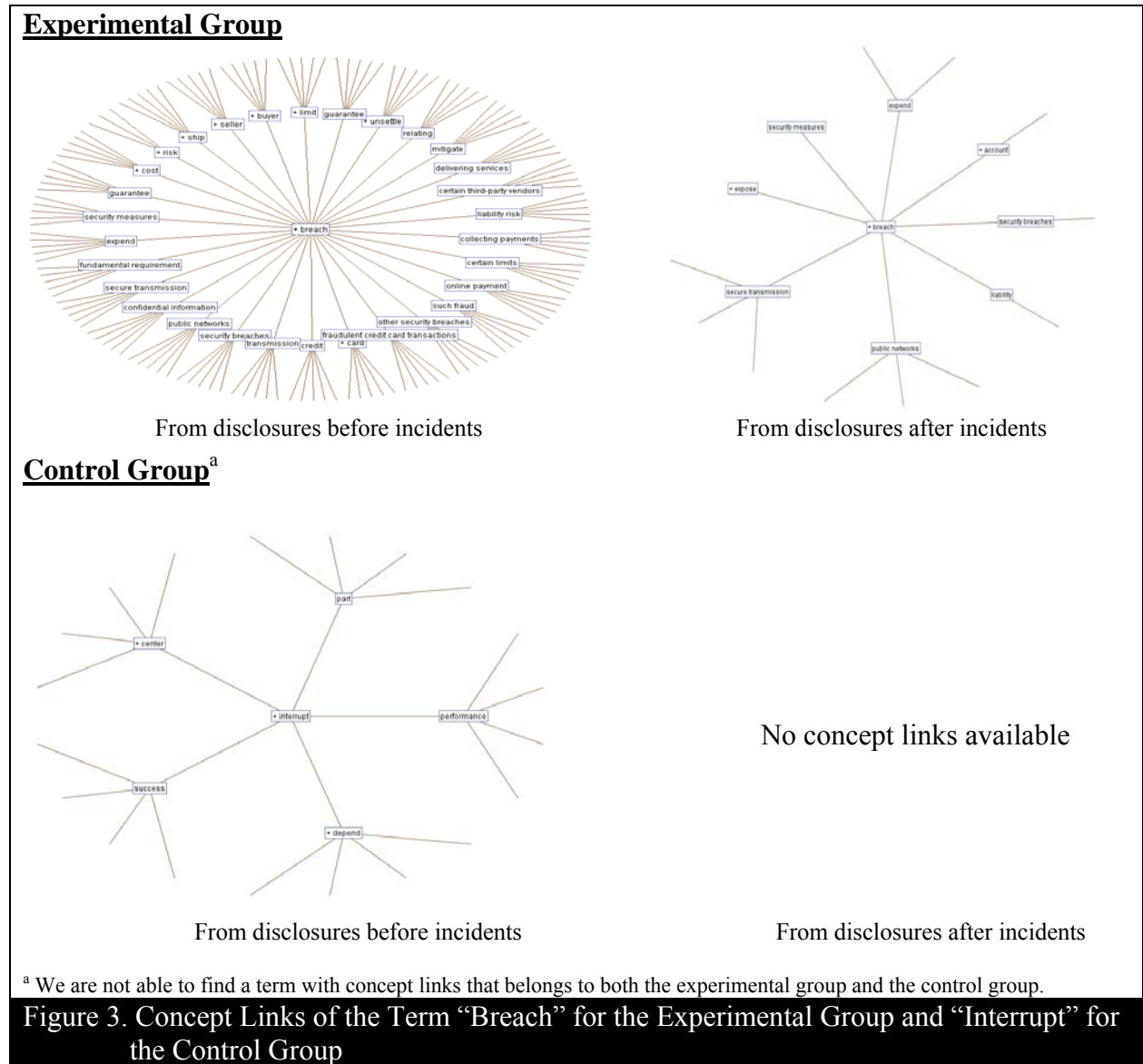
^a The descriptive terms over the percentage of 15% are reported for the experimental group.

Table 7. Text Mining Results of Information Security Risk Factors by Year Groups after Incidents ^a				
Year Group	Terms	Freq.	Percentage	RMS Std.
Resolution: High				
1997-1998	No disclosures			
1999-2000	+experience, +have, +interruption, +provide, +service	10	53%	0.263
2001-2002	+affect, +event, +failure, +operation, +system	6	75%	0.151
2003-2004	+failure, +interruption, +network, +product, +service	13	68%	0.240
2005-2006	data, +disruption, +experience, +failure, +have	11	58%	0.249
Resolution: Medium				
1997-1998	No disclosures			
1999-2000	+experience, +interruption, +loss, +provide, +user	13	68%	0.291
2001-2002	+affect, +event, +failure, +operation, +system	6	75%	0.151
2003-2004	+disruption, +failure, +interruption, +provider, +service	11	58%	0.267
2005-2006	+disruption, +event, +experience, +failure, +service	10	53%	0.255
Resolution: Low				
1997-1998	No disclosures			
1999-2000	+experience, +have, +interruption, +loss, +provide	11	58%	0.266
2001-2002	+affect, +event, +failure, +operation, +system	6	75%	0.151
2003-2004	adversely, +affect, +customer, other, +system	11	58%	0.299
2005-2006	data, +disruption, +experience, +failure, +service	9	47%	0.220

^a The descriptive terms with the largest frequency are reported for the experimental group.

In order to better understand the terms in the clusters, we further connect the terms in the cluster with other phrases in the disclosures. For example, the term “traffic” in the cluster

usually co-occurs with “speed”, “infrastructure”, and “response”. Also, as expected, the terms “attack” and “denial” are often disclosed together. This co-occurrence relationship can be captured by concept links (see Appendix D).



The relationship between the term “breach” for the experimental group in the cluster with other phrases in the disclosures, i.e. concept links, provides interesting insights as shown in Figure 3. The term “breach” mainly links to confidentiality related terms in the statements, such as “confidential information”, “secure transmission”, and “fraudulent credit card transactions”.

That is, when firms discuss the term “breach” in risk factors, confidentiality related terms usually occurred together which reflect the importance of these terms on a firm’s future performance. It seems that the general term “breach” is used to describe the security policy or practices of many specific risk factors, especially confidentiality risk factors. On the other hand, for the control group, for example, the term “interrupt” co-occurs with “performance” and “success”. It seems that the control group discloses relatively more general statements than the breached firms do. This result also confirms that the breached firms focus more on threats or attacks but the firms in the control group pay more attention on operations or processes.

We further compare whether there is any change in concept links in order to better explain the change in disclosures. Thus, we also investigate the relationship between the term “breach” with other terms in the disclosures *after* the incidents for the experimental group in Figure 3. Figure 3 shows that the term “breach” in the disclosures after the incidents becomes more like a general term. That is, when the firm needs to discuss the impacts of specific risk factors, such as confidentiality risk factors, the firm tends to use the phrases specific to that risk factor instead of using the term “breach” frequently. Combining what we found before and after incidents, it seems that firms are more specific or more focused in disclosing risk factors after experiencing information security incidents. Interestingly, we are not able to find any concept links for the control group after the incidents. It seems that the firms in the control group still disclose relatively general statements and more diversified risk factors than the breached firms do.

From the above analysis, our findings suggest that firms learn from the incidents and respond to the incidents by disclosing more focused risk factors to financial report users. Based on the accounting literature, managers in breached firms attempt to disclose more specific risk factors after experiencing incidents in order to reduce the possibility of future lawsuits or the loss

of reputation due to information security incidents which appears to be consistent with Hypothesis 1. Furthermore, these firms are still willing to disclose more risk factors even they know that investors will punish them in terms of stock price reactions if they face information security incidents. According to the literature in disclosures, the punishment in stock price must be smaller than the legal and reputation costs these firms might occur without disclosures which drives the firms to disclose more factors. However, there remains one issue that we have not solved. That is, what kinds of disclosures will result in different investors' perceptions? We further explore this issue through a classification model.

5.3 Classification Model

As shown in Figure 2 at the beginning of this section, in addition to the qualitative comparison above, we further build a classification model based on the disclosures of information security risk factors and determine whether we can distinguish different stock price reactions to information security incidents. The classification results can allow us to understand the difference of the disclosure practices between the groups of firms and to uncover the underlying terms within disclosures that result in different perceptions. However, the results attributing CAR to disclosure is preliminary.

To build the model, the disclosures in our sample (70 documents) are categorized into three groups. The first group (named as the positive group) consists of all the disclosures associated with a positive stock price reaction more than +3% after information security incidents. The second group (named as the stable group) has all the disclosures associated with the stock price reaction within -3% to +3% after information security incidents. The third group (named as the negative group) has all the disclosures associated with the stock price reaction more than -3% after information security incidents. We use 3% because it results in the most diverse

distribution of groups. The purpose of the model is to distinguish between these three groups. All these 80 documents are partitioned into three parts: training, validation and test. The percentage values for these three data sets are 80%, 10%, and 10% respectively. Based on the text mining results, a classification model is trained, validated, and tested using Memory Based Reasoning (MBR). Although different methods have been used to classify text documents, such as Young et al. (1985), Hayes et al. (1990), and Rau et al. (1991), MBR has the advantage of not requiring manual definition of topics but achieving at least moderate precision and high recall, i.e. the classifier can find relevant documents and assign them to the correct category (Masand et al. 1992; SAS Institute Inc 2004). After the model is trained, the validation results show that the model correctly classifies the disclosures of the negative group about 65% of the time. As a robustness check, we resample the training, validation, and test data sets ten times based on the 80%, 10%, and 10% values and obtain similar results. This model is then applied to predict the categories of 24 new disclosures from 50 firms across different industries with different sizes that are not in the current sample and without breach announcements. The firms are, for example, Kroger, Motorola, JDS Uniphase, Delta Airlines, Blockbuster, and etc. Based on the patterns within the disclosures in our sample, the “score” node determines which disclosure belongs to either one of the three groups based on the model developed. The classification results demonstrate that about 83% of the time our model categorize these new disclosures into the stable group. This model is also applied to all the documents after the incidents. The model groups the document into the stable group about 61% of the time.

Interestingly, when we link the clusters with groups, we found that more than 60% of the time the cluster with the terms “reputation”, “+customer”, “security”, “information”, and “+breach” will be considered a stable group. The terms “+depend”, “implement”, “+develop”,

“+require” are more related to positive group. It seems that these action terms are more related to a positive stock price reaction to information security incidents. This model explains why some firms are not punished even if they disclose information security incidents. First, as the firms in the control group disclose, action terms or terms about processes do not create a perception that the managers are attempting to preempt the announcement of bad news in order to reduce future legal costs. Second, disclosures focus more on general statements about company image and security information or policy also seems to convey the information to the investors without formulating the perception of hiding bad news. This result also confirms the disclosure pattern in the control group and the result in the key phrases analysis.

6. Conclusions and Discussion

This paper investigates the relationship between information security related disclosures in financial reports and information security incidents. We use two different measures for information security related disclosures. One is the disclosures of internal control and procedures mandated by SOX. The other one is the disclosures of risk factors. Based on the observations of information security incidents we obtain from 1997 to 2006, our results do not provide enough evidence for the relationship between the disclosures of internal control and procedures and CAR. But the findings do demonstrate that the disclosures of information security risk factors statistically significantly increase the impact of information security incidents. After further investigating the argument of expectation formulation in detail, we find that the impact of information security incidents on stock price reactions depend on whether the incidents match the content of the disclosures. We also demonstrate that there indeed exists a difference between first time disclosures and acknowledgements. The text mining results point out that breached firms react to information security incidents by disclosing more specific and

additional risk factors in subsequent financial reports. Combining the text mining results with our classification model, we show that firms need to disclose generic and more actionable information when they provide information security risk factors. By doing so, the firm can not only reduce possible future lawsuit and reputation losses but also lower the impact of information security incidents. By disclosing in other patterns, the disclosures can be seen as a warning of future incidents which can harm its business value.

This research has implications for both researchers and practitioners. For researchers, our findings provide explanations that can be considered when conceptualizing risk management strategies especially after the implementation of SOX. Furthermore, for daily operations, these risk factors convey important information for investors and can be used for estimating the impacts on those firms. This paper also demonstrates possible issues in the literature of voluntary disclosure. As shown by the differences between the disclosure practices across our experimental and control group, there might be a different disclosure practice than the ones discussed in the literature. More importantly, there can be some additional firm characteristics or market impacts that lead to different disclosure practices.

For practitioners, the results shed light on how they can convey security practices to their customers and investors more effectively. We observe that standardized disclosures of information security related issues provide relatively little information. It is not difficult to tell that the information those disclosures provide is meant to satisfy the requirements of SOX as opposed to voluntary disclosures. By properly reflecting possible security concerns, a firm should be able to convey their security practices and concerns to investors without being considered as a warning of subsequent incidents. Moreover, when firms try to generate the disclosing information, executives should also consider the effectiveness of information security

governance. That is, the firm can at the same time identify the existing and emerging risk factors and assess the impact of such factors. Based on the identification and assessment, executives are able to evaluate whether the current information security policies and practices are adequate and convey some of the information to investors for business value evaluation and investment decisions. Last, although our results do not focus on the impact of spillover effect and do not find it significant on CAR, our preliminary results demonstrate the control group might also be affected by security breach announcements for firms in the experimental group. Therefore, when considering the possible impacts of information security incidents, managers also need to take other firms effect into account. Moreover, firms in the same industry can also cooperate to mitigate the possible impacts of information security incidents.

The analyses in the paper suffers from a few limitations. One of the major limitations of our study is sample size. Although we attempt to capture as large of a sample as possible, it is still problematic to collect a larger data set base on our filtering processes and our research questions. A larger data set allows us to get different perspective of the text mining results from different industries. A larger data set also makes the classification model more reliable. Furthermore, many firms might suffer from information security incidents that are not disclosed to the public. Obviously, we are unable to incorporate this information into our sample. Second, we implicitly assume that the stock price truly reflects a firm's business value. Although the stock price for high-tech firms might be biased, we only look at the price change in a short time period. Thus, we believe that our results still hold even with this possibility. Third, we adopt a simple coding scheme for the disclosures. Although we believe that a more complicated coding scheme does not alter our main results, a finer coding scheme for all the disclosures that can be applied to different industries may provide more details than the present scheme.

Possible future extensions are as follows. First, in our paper, we implicitly assume that the disclosures are credible and truly reflect a firm's practices. However, some firms might disclose lots of information but invest little. On the other hand, some other firms might invest in information security but refuse to disclose such investments to the public. Therefore, this anomaly is worth further investigation. Second, a larger data set can be used to provide more meaningful text mining results for both information security risk factors and business risk factors. The text mining analysis of business risk factors can also provide a first glance on how these risks affect different businesses. Third, as different media becomes popular information sources for investors, we can further consider other media sources, such as blogs, to investigate the relationship among different information sources, information security incidents, and stock price reactions. Last, the spillover effect can be investigated in detail by considering how the information is transferred and the major factors that result in the spillover effect.

Reference

- Ajinkya, B. B., M. J. Gift. 1984. Corporate managers' earnings forecasts and symmetrical adjustments of market expectations. *J. of Accounting Res.* **22**(2) 425-444.
- Balakrishnan, K., A. Ghose, P. Ipeiritis. 2008. The impact of information disclosure on stock market returns: the Sarbanes-Oxley Act and the role of media as an information. Working Paper, New York University.
- Bagnoli, M., S. G. Watts. 2007. Financial reporting and supplemental voluntary disclosures. *J. of Accounting Res.* **45**(5) 885-913.
- Bagnoli, M., W. Kross, S. G. Watts. 2002. The information in management's expected earnings report date: a day late, a penny short. *J. of Accounting Res.* **40**(5) 1275-1296.
- Begley, J., P. Fischer. 1998. Is there information in an earnings announcement delay? *Rev. of Accounting Studies* **3** 347-363.
- Bowen, P., J. Hash, M. Wilson. 2006. *Information security handbook: a guide for managers*, NIST Special Publication 800-100.
- Campbell, K., L. A. Gordon, M. P. Loeb, L. Zhou. 2003. The economic cost of publicly announced information security breaches: empirical evidences from the stock market. *J. of Computer Security* **11** 431-448.
- Cavusoglu, H., B. Mishra, S. Raghunathan. 2004. The effect of internet security breach announcements on market value of breached firms and internet security developers. *Internat. J. of Electronic Commerce* **9**(1) 69-105.
- Cecchini, M., H. Aytug, G. J. Koehler, P. Pathak. 2007. Detecting Management Fraud in Public Companies. Working Paper, University of South Carolina.
- CERT. 2007. *CERT/CC Statistics 1988-2006*, Retrieved Apr. 9 2007, from http://www.cert.org/stats/cert_stats.html.
- CSI/FBI. 2007. *The CSI/FBI computer crime and security report in 2006*, Retrieved Apr. 9 2007, from <http://abovesecurity.com/doc/CommuniquesPDF/FBISurvey2006>.
- Darrrough, M. N. 1993. Disclosure policy and competition Cournot vs. Bertrand. *The Accounting Rev.* **68**(3) 534-561.
- Dye, R. A. 1985. Disclosure of Nonproprietary Information. *J. of Accounting Res.* **12**(1) 123-145.
- Dye, R. A. 2001. An evaluation of 'essays on disclosure' and the disclosure literature in accounting. *J. of Accounting and Econom.* **32** 181-235.
- Eihorn, E. 2005. The nature of the interaction between mandatory and voluntary disclosures. *J. of Accounting Res.* **43**(4) 593-621.
- Elliott, R., P. Jacobson. 1994. Costs and benefits of business information disclosure. *The Accounting Horizons* **8**(4) 80-96.
- Ettredge, M. L., V. J. Richardson. 2003. Information transfer among internet firms: the case of hacker attacks. *J. of Inform. Systems* **17**(2) 71-82.
- Fama, E. 1970. The behavior of stock market prices. *J. of Finance* **25** 383-417.
- Fama, E., K. French. 1992. The cross-section of expected stock returns. *J. of Finance* **47**(2) 427-465.
- Fan, W., L. Wallace, S. Rich, Z. Zhang. 2006. Tapping the power of text mining. *Comm. of the ACM* **49**(9) 77-82.
- Feldman, R., J. Sanger. 2006. *The text mining handbook: advanced approaches in analyzing unstructured data*, UK: Cambridge University Press.
- Francis, R., D. Philbrick, K. Schipper. 1994. Shareholder litigation and corporate disclosure. *J. of Accounting Res.* **32**(2) 137-164.

- Francis, J., K. Schipper, L. Vincent. 2002. Expanded disclosures and the increased usefulness of earnings announcements. *The Accounting Rev.* **77**(3) 515-546.
- Foster, G. 1981. Intra-industry information transfers associated with earnings releases. *J. of Accounting and Econom.* **3**(3) 201-232.
- Gal-Or, E., A. Ghose. 2005. The economic incentives for sharing security information. *Inform. Systems Res.* **16**(2) 186-208.
- Garg, A., J. Curtis, H. Halper. 2003. Quantifying the financial impact of IT security breaches. *Inform. Management & Computer Security* **11**(2) 74-83.
- Glover, S., S. Liddle, D. Prawitt. 2001. *Electronic commerce: security, risk management, and control*, NL: Prentice Hall.
- Gordon, L. A., M. P. Loeb. 2002. The economics of information security investment. *ACM Transac. on Inform. and System Security* **5**(4) 438-457.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn. 2003. Sharing information on computer systems security: an economic analysis. *J. of Accounting and Public Policy* **22**(6) 461-485.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn, R. Richardson. 2005. *Tenth annual CSI/ FBI computer crime and security survey*. Computer Security Institute, 1-26.
- Gordon, L. A., M. P. Loeb, W. Lucyshyn, T. Sohail. 2006. The impact of the Sarbanes-Oxley Act on the corporate disclosures of information security activities. *J. of Accounting and Public Policy* **25** 503-530.
- Grossman, S. J. 1981. The information role of warranties and private disclosure about product quality. *J. of Law and Econom.* **24**(3) 461-483.
- Hayes, P. J., P. M. Anderson, I. B. Nirenburg, L. M. Schmandt. 1990. TCS: a shell for content-based text categorization. *Proc. of the 2nd IEEE Conf. on AI Applications*, Santa Barbara, 320-326.
- Healy, P. M., K. G. Palepu. 2001. Information asymmetry, corporate disclosure, and the capital markets: a review of the empirical disclosure literature. *J. of Accounting and Econom.* **31**(1-3) 405-440.
- Hovav, A., J. D'Arcy. 2003. The impact of denial-of-service attack announcements on the market value of firms. *Risk Management and Insurance Rev.* **6**(2) 97-121.
- Han, J., R. Altman, V. Kumar, H. Mannila, D. Pregibon. 2002. Emerging scientific applications in data mining. *Comm. of the ACM* **45**(8) 54-58.
- Jo, H., Y. Kim. 2007. Disclosure frequency and earnings management. *J. of Financial Econom.* **84**(2) 561-590.
- Jorgensen, B. N., M. T. Kirschenheiter. 2003. Discretionary risk disclosures. *The Accounting Rev.* **78**(2) 449-469.
- Kannan, K., J. Rees, S. Sridhar. 2007. Market reactions to information security breach announcements: an empirical study. *Internat. J. of Electronic Commerce* **12**(1) 69-91.
- Kaszniak, R., B. Lev. 1995. To warn or not to warn: management disclosures in the face of an earnings surprise. *The Accounting Rev.* **70**(1) 113-134.
- Kaszniak, R., M. F. McNichols. 2002. Does meeting earnings expectations matter? Evidence from analyst forecast revisions and share prices. *J. of Accounting Res.* **40**(3) 727-759.
- Katz, S. B. 2001. Language and persuasion in biotechnology communication with the public: How not to say what you're not going to say and not say it, *AgBioForum* **4**(2) 93-97.
- King, R., G. Pownall, G. Waymire. 1990. Expectations adjustment via timely management forecasts: review, synthesis, and suggestions for future research. *J. of Accounting Lit.* **9** 113-144.

- Lang, M. H., R. J. Lundholm. 1993. Cross-sectional determinants of analyst ratings of corporate disclosures. *J. of Accounting Res.* **31** 216-271.
- Lang, M. H., R. J. Lundholm. 1996. Corporate disclosure policy and analyst behavior. *The Accounting Rev.* **71**(4) 467-492.
- Lang, M. H., R. J. Lundholm. 2000. Voluntary disclosure and equity offerings: reducing information asymmetry or hyping the stock? *Contemporary Accounting Res.* **17**(4) 623-662.
- Landsman, W., E. Maydew. 2002. Has the information content of quarterly earnings announcements declined in the past three decades? *J. of Accounting Res.* **40**(3) 797-807.
- Lev, B., S. H. Pennman. 1990. Voluntary forecast disclosure, nondisclosure, and stock prices. *J. of Accounting Res.* **28**(1) 49-76.
- Li, F. 2006. Annual report readability, current earnings, and earnings persistence. Working Paper, University of Michigan.
- MacKinlay, A. C. 1997. Event studies in economics and finance. *J. of Econom. Lit.* **35**(1) 13-39.
- Masand, B., G. Linoff, D. Waltz. 1992. Classifying news stories using memory based reasoning. *Proc. of the 15th Annual Internat. ACM SIGIR Conf. on Res. and Development in Inform. Retrieval*, Copenhagen, Denmark, 59-65.
- Milgrom, P. R. 1981. Good news and bad news: representation theorems and applications. *Bell J. of Econom.* **12**(2) 380-391.
- Penno, M. 1997. Information quality and voluntary disclosure. *The Accounting Rev.* **72**(2) 275-284.
- PriceWaterhouseCoopers. 2002. *Inform. Security Breaches Survey 2002 – A Technical Report*. Prepared by PriceWaterhouseCoopers for the Department of Trade and Industry.
- Rau, L. F., P. S. Jacobs. 1991. Creating segmented databases from free text for text retrieval. *Proc. of the 15th Annual Internat. ACM SIGIR Conf. on Res. and Development in Inform. Retrieval*, Chicago, IL, 337-346.
- Sandoval, G., T. Wolverson. 2000. Leading web sites under attack. Retrieved April 17, 2007, from http://news.com.com/Leading+Web+sites+under+attack /2100-1017_3-236683.html.
- SAS Institute Inc. 2004. *Getting started with SAS® 9.1 text miner*. Cary, NC: SAS Institute Inc.
- Shadish, W. R., T. D. Cook, D. T. Campbell. 2002. *Experimental and quasi-experimental designs for generalized causal inference*, MA: Houghton Mifflin Company.
- Skinner, D. J. 1994. Why firms voluntarily disclose bad news. *J. of Accounting Res.* **32**(1) 38-60.
- Sohail, T. 2006. *To tell or not to tell: market value of voluntary disclosures of information security activities*. Unpublished doctoral dissertation, University of Maryland, Maryland.
- Stocken, P. 2000. Credibility of voluntary disclosure. *RAND J. of Econom.* **31**(2) 359-374.
- Tan, A. H. 1999. Text mining: the state of the art and the challenges. *Proc. of the PAKDD'99 Workshop on Knowledge discovery from Advanced Databases*, Beijing.
- Tran, K., R. L. Rundle. 2000. Hackers attack major internet sites, cutting off Amazon, Buy.com, eBay. *The Wall Street Journal*. Retrieved March 2, 2007, from FACTIVA database.
- Verrecchia, R. E. 1983. Discretionary disclosure. *J. of Accounting and Econom.* **5**(3) 179-194.
- Verrecchia, R. E. 2001. Essays on disclosures. *J. of Accounting and Econom.* **32**(1-3) 97-180.
- Warren, M. J., W. E. Hutchinson. 2000. Cyber attacks against supply chain management systems. *Internat. J. of Physical Distribution and Logistics Management* **30** 710-716.
- Young, S. R., P. J. Hayes. 1985. Automatic classification and summarization of banking telexes. *Proc. of the 2nd IEEE Conf. on AI Applications*, Miami Beach, FL, 402-409.

Appendix A. An Example of the Disclosures of Internal Control and Procedures

Excerpt from Yahoo's annual report for year 2005, retrieved on Apr.23, 2007

source: http://www.sec.gov/Archives/edgar/data/1011006/000110465906014033/a06-3183_110k.htm

“Evaluation of Disclosure Controls and Procedures

The Company's management, with the participation of the Company's principal executive officer and principal financial officer, has evaluated the effectiveness of the Company's disclosure controls and procedures (as such term is defined in Rules 13a-15(e) and 15d-15(e) under the Securities Exchange Act of 1934, as amended (the “Exchange Act”) as of the end of the period covered by this report. Based on such evaluation, the Company's principal executive officer and principal financial officer have concluded that, as of the end of such period, the Company's disclosure controls and procedures are effective in recording, processing, summarizing and reporting, on a timely basis, information required to be disclosed by the Company in the reports that it files or submits under the Exchange Act.

Management's Report on Internal Control Over Financial Reporting

The Company's management is responsible for establishing and maintaining adequate internal control over financial reporting as defined in Rules 13a-15(f) and 15d-15(f) under the Exchange Act. Under the supervision and with the participation of the Company's management, including its principal executive officer and principal financial officer, the Company conducted an evaluation of the effectiveness of its internal control over financial reporting based on criteria established in the framework in Internal Control—Integrated Framework issued by the Committee of Sponsoring Organizations of the Treadway Commission. Based on this evaluation, the Company's management concluded that its internal control over financial reporting was effective as of December 31, 2005.

Because of its inherent limitations, internal control over financial reporting may not prevent or detect misstatements. Also, projections of any evaluation of effectiveness to future periods are subject to the risks that controls may become inadequate because of changes in conditions, or that the degree of compliance with the policies or procedures may deteriorate.

The Company's independent registered public accounting firm has audited management's assessment of the effectiveness of the Company's internal control over financial reporting as of December 31, 2005 as stated in their report which appears on page 58.

Changes in Internal Control Over Financial Reporting

There have not been any changes in the Company's internal control over financial reporting (as such term is defined in Rules 13a-15(f) and 15d-15(f) under the Exchange Act) during the most recent fiscal quarter that have materially affected, or are reasonably likely to materially affect, the Company's internal control over financial reporting.”

Appendix B. Examples of Risk Factors

Excerpt from Amazon's annual report for year 2000, retrieved on Apr.23, 2007

source: <http://www.sec.gov/Archives/edgar/data/1018724/000103221001500087/0001032210-01-500087.txt>

“We Face Intense Competition

The e-commerce market segments in which we compete are relatively new, rapidly evolving and intensely competitive. In addition, the market segments in which we participate are intensely competitive and we have many competitors in different industries, including the Internet and retail industries.

Many of our current and potential competitors have longer operating histories, larger customer bases, greater brand recognition and significantly greater financial, marketing and other resources than we have. They may be able to secure merchandise from vendors on more favorable terms and may be able to adopt more aggressive pricing or inventory policies. They also may be able to devote more resources to technology development and marketing than us.

As these e-commerce market segments continue to grow, other companies may enter into business combinations or alliances that strengthen their competitive positions. We also expect that competition in the e-commerce market segments will intensify. As various Internet market segments obtain large, loyal customer bases, participants in those segments may use their market power to expand into the markets in which we operate. In addition, new and expanded Web technologies may increase the competitive pressures on online retailers. The nature of the Internet as an electronic marketplace facilitates competitive entry and comparison shopping and renders it inherently more competitive than conventional retailing formats. This increased competition may reduce our operating profits, or diminish our market segment share.”

“System Interruption and the Lack of Integration and Redundancy in Our Systems May Affect Our Sales

Customer access to our Web sites directly affects the volume of goods we sell and thus affects our net sales. We experience occasional system interruptions that make our Web sites unavailable or prevent us from efficiently fulfilling orders, which may reduce our net sales and the attractiveness of our products and services. To prevent system interruptions, we continually need to: add additional software and hardware; upgrade our systems and network infrastructure to accommodate both increased traffic on our Web sites and increased sales volume; and integrate our systems.

Our computer and communications systems and operations could be damaged or interrupted by fire, flood, power loss, telecommunications failure, break-ins, earthquake and similar events. We do not have backup systems or a formal disaster recovery plan, and we may have inadequate insurance coverage or insurance limits to compensate us for losses from a major interruption. Computer viruses, physical or electronic break-ins and similar disruptions could cause system interruptions, delays and loss of critical data and could prevent us from providing services and accepting and fulfilling customer orders. If this were to occur, it could damage our reputation.”

Appendix C. Stock Price Reactions from Information Security Incidents

In our study, the market model is used to capture the impact of security incidents.

$$R_{it} = \beta_0 + \beta_1 R_{mt} + \varepsilon_{it} \quad (\text{A-1})$$

where R_{it} denotes company i 's return at period t which equals to $(p_t - p_{t-1}) / p_{t-1}$. Dividends and stock splits are excluded here because (1) they are rare events and (2) we have already considered confounding events. Thus, stock return of a certain company equals to the change in stock price or the capital gain. R_{mt} stands for the corresponding market return at period t and is estimated by the CRSP equally weighted index. The CRSP equally weighted index is the average of the returns of all trading stocks in NYSE, AMEX and NASDAQ. β_0 and β_1 are the parameters and estimated in a 255-day periods ending at 45 days before the estimation window we choose by ordinary least square (OLS) method. We calculate the abnormal return (AR) from the market model:

$$AR_{it} = R_{it} - \hat{\beta}_0 - \hat{\beta}_1 R_{mt} \quad (\text{A-2})$$

As shown by equation (A-2), abnormal return is the return that cannot be captured by the market as a whole or the ex post return over the event window minus the normal return. The total effect of an economic event on stock price is reflected in mean cumulative abnormal return, which is the summation of abnormal returns for company-event observations in the window we choose, i.e. $(\sum_{t=1}^N \sum_{t_0}^{t_1} AR_{it})/N$, where t_0 and t_1 are the beginning and the ending trading day for the window we choose. Cumulative abnormal return (CAR, $\sum_{t_0}^{t_1} AR_{it}$) for each observation is used for the cross-sectional analysis.

Appendix D. Cluster Analysis and Concept Links

The cluster analysis is performed as follows using SAS[®] 9.1 Text Miner. First, text parsing decomposes the sentences into terms and creates a frequency matrix as a quantitative representation of the input documents. When decomposing the documents, we choose to rule out definite as well as indefinite articles, conjunctions, auxiliaries, prepositions, pronouns and interjections since these terms do not help provide meaningful results in our context. This matrix also shows the weight for the terms. The weight for term i in document j (w_{ij}) is the multiplication of the frequency weight (L_{ij}) and the term weight (G_i). In our study, the frequency weight is the logarithm of the frequency (f_{ij}) of term i in document j plus one, i.e. $L_{ij} = \log_2(f_{ij} + 1)$. The term weight of term i (G_i) is calculated as $1 + \sum_j (p_{ij} \log_2(p_{ij}) / \log_2(n))$, where $p_{ij} = f_{ij} / g_{ij}$, g_i is the number of times term i appears in the data set, and n is the number of documents in the data set. These two methods put more weights on words that show in few documents and generally give the best results (SAS Institute Inc 2004). For dimension reduction, we use the single value decomposition (SVD) method. SVD generates the dimensions that best represent the original frequency matrix. The singular value decomposition of a frequency matrix (A) is to factorize the matrix into matrices of orthonormal columns and a diagonal matrix of singular values, i.e. $A = U\Sigma V^T$. Then the original documents are projected to matrix U (SAS Institute Inc 2004). Through matrix factorization and projection, SVD forms the dimension-reduced matrix. In our analysis, we set the maximum reduced dimensions to be one hundred (as default) and test three different levels of reduced dimensions (high, medium and low resolutions) as a robustness check. The resulting SVD dimensions are further used for cluster analysis. We then divide our data into disjoint groups using expectation maximization clustering by setting the maximum clusters to be forty (as default). The expectation maximization method is an iterative process that estimates the parameters in the mixture model probability density function which approximates that data distribution by fitting k cluster density function to a data set. The mixture model probability density function evaluated at point x equals $\sum_{h=1}^k \omega_h f_h(x|\mu_h, \Sigma_h)$, where μ_h , Σ_h are the mean vector and covariance matrix for cluster h under Gaussian probability distribution. For each observation x at iteration j , whether x belongs to a cluster h equals to $(\omega_h^j f_h(x|\mu_h^j, \Sigma_h^j)) / (\sum_i \omega_i^j f_i(x|\mu_i^j, \Sigma_i^j))$ (SAS Institute Inc 2004). The iteration terminates if the likelihood value of two iterations is less than $\varepsilon > 0$ or a maximum of five iterations are reached (SAS Institute Inc 2004). The text mining results are discussed in section 4.3.2.

The concept links are determined based on the following criteria when all three of them are met: (1) Both

terms occur in at least n documents, where n equals $\text{Max}(4, A, B)$. A is the largest value of the number of documents that a term appears in divided by 100 and B is the 1000th largest value of the number of documents that a term appears in for concept links (SAS Institute Inc 2004), (2) Term 2 occurs when term 1 occurs at least 5% of the time (SAS Institute Inc 2004), and (3) The relationship between terms is highly significant (the chi-square statistic is greater than 12) (SAS Institute Inc 2004).