

## A Study of Keywords based on The Word Frequency Effect Theory in Video Lectures of Software Engineering Education for Detecting Mind

Jaechoon Jo

Dept. Computer Science and Engineering  
Korea University  
Seoul, Korea  
e-mail: jaechoon@korea.ac.kr

Heuseok Lim

Dept. Computer Science and Engineering  
Korea University  
Seoul, Korea  
e-mail: limhseok@korea.ac.kr

**Abstract**— The increased popularity of Massive Open Online Courses (MOOC) and e-learning has constantly increased video-based online education platforms. There are also many video lectures for software engineering education in online education platforms. Although online lectures have many advantages, there are also limitations. We performed a verification research to see if high frequency words can detect mind wandering to resolve existing limitations. In this verification study, experiments to identify whether high frequency words can represent the software engineering video lecture, the minimum number of words needed to detect mind wandering, and whether mind wandering detection standards should be changed according to the length of the video lecture. The results of this study confirmed that mind wandering can be detected through high frequency words and they can be used as an important feature in various learning analysis investigations to resolve existing limitations of online education.

**Keywords**- Attention; Online Education; e-Assessment; Word embedding; Forgetting curve; Working memory

### I. INTRODUCTION

The market size of e-learning is 3,485 billion won (over three billion USD) in Korea and is growing every year [1]. With the addition of the mobile device environment, the e-learning market seems to be growing constantly. E-learning has advantages such as ‘users can learn anytime and anywhere’ and ‘users can control playback speed and repetitions according to their learning pace.’ On the other hand, it needs users’ concentration to get a good result. Therefore, many different assessment methods have been researched and developed. Furthermore, teachers want to verify that learners have really watched the online video lectures. In the conventional e-learning environment, learning was checked through a quiz, which imposes a cognitive burden on both teachers and learners. The biggest problem of e-learning is decreased concentration [2]. Unlike offline lectures, where instructors and learners interact, online lectures provide information unilaterally. To supplement this shortcoming, various studies to enhance concentration have been performed [3]. Additionally, with video-based online lectures where many learners need to be managed, it is very important to automatically assess if

learners have watched the lecture without mind wandering [4].

Therefore, in a previous study, the minimum learning judgement system which detects concentration level through a simple word game was developed rather than conventional methods such as a quiz, pop up events or the length of video playback [5]. The word game asks if the presented word was used in the video lecture. This system allows the instructor to check if there was mind wandering by means of a word game. In this study, ‘Mind Wandering’ signifies a learner who is not concentrating on the video lecture while it is being played. The word game abstracts words used in the video lecture automatically and detects mind wandering using high frequency words.

Thus, a verification study is needed to determine the importance of high frequency words in the detection of mind wandering using software engineering video lectures. This study was conducted to verify whether high frequency words can be used as an important feature to detect mind wandering before actually using them to detect mind wandering in software engineering video lectures. Additionally, we sought to verify and propose that the words used in video lectures can be used as an important feature in various research areas, such as automatic quiz generation, to detect whether users watched and concentrated on the video lectures. Research questions for this study were the following (RQs):

RQ 1. Can high frequency words represent the software engineering video lecture?

RQ 2. How many words should be used for the word game to detect mind wandering?

RQ 3. Should the detection standards be different according to the length of the video lecture?

### II. BACKGROUND

We developed a Minimum Learning Judgement System (MLJS) that can automatically detect mind wandering in a video-based online lecture [5]. In this study, ‘Minimum Learning’ represents the status when learners have achieved minimal learning behavior without mind wandering through an online video lecture. Minimum learning behavior is achieved when learners have actually watched the video lecture regardless of their understanding level. Based on the word frequency effect theory, minimum learning is

automatically detected using data created during cognitive processing [6]. Words for the word game are created automatically as shown in Figure 1.

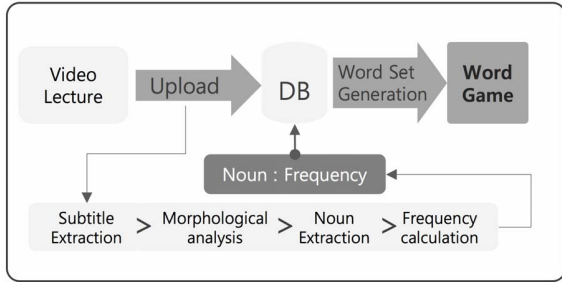


Figure 1. Word generation model

When an instructor uploads a video lecture to the minimum learning judgement system, the system automatically extracts the subtitle. Then it extracts words that are nouns through a morpheme analyzer in the extracted subtitle. The frequencies of the extracted nouns are calculated and the nouns and the calculated frequencies are stored in a database. When the word game is initiated, the minimum learning judgement system compares the words in the video lectures that the user has previously watched and the words in the current video lecture, calculates the weight of the words and finally selects high frequency words which shall be used in the word game. Finally, the selected words are combined with words (randomly picked from the database) which were never used in the video lecture and one of the words is presented to the learner. The learner selects O or X according to their memory of watching the video lecture. Figure 2 shows a screen of the word game.

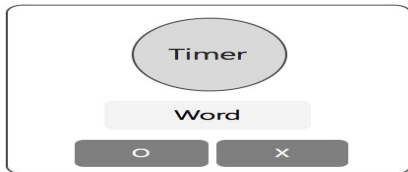


Figure 2. A screen of the word game

### III. EXPERIMENTS AND RESULTS

Various kinds of experiments have been performed to identify if high frequency words in a software engineering video lecture can be used as a feature to detect mind wandering. The aforementioned three research questions were verified in this experiment. For the first research question, we investigated whether high frequency words can be used as keywords that represent the software engineering online video lecture and whether they can be used to detect mind wandering. For this, a word embedding technology was applied. For the second research question, we examined the number of words that should be used to appropriately detect mind wandering. Finally, for the third research question, we verified whether mind wandering detection standards should vary according to the length of the video lecture.

#### A. High Frequency Words

A study of the importance of high frequency words in a software engineering video lecture was executed. Mind wandering can be detected through a word game using high frequency words in an online video lecture. This study examined whether high frequency words used in a word game are an appropriate feature to detect mind wandering.

For the verification experiment, the word2vec method was used. High frequency words abstracted from the online video lecture, keywords entered by experts and random words not related to the lecture were embedded in a vector space. Similarity of the words were calculated and compared. To embed words in the vector space, an existing open source Gensim was used [7]. The Korean Wikipedia was used as a model study. Figure 3 illustrates the experimental process.

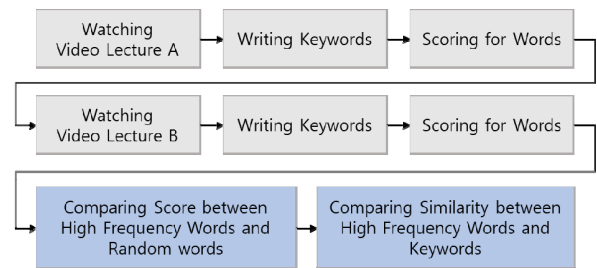


Figure 3. The experimental process of high frequency words

To assess the importance of high frequency words, an expert group was composed of members who satisfied the following conditions:

- Those who have a master's or doctorate degree in the discipline related to the software engineering video lectures, and,
- Those who have more than one teaching experience in software engineering.

The expert group wrote a list of keywords after watching online video lectures and evaluated the importance of high frequency words and random words. Two video lectures were used for this process. Lecture A was related to the software engineering education discipline and Lecture B was a general lecture that was not related to any specific discipline. After the experiments, the scores of high frequency and random words were compared. All the words (high frequency words, random words and keywords) were embedded in a vector space and their similarity was compared. Table 1 provides the experts' scores of importance for high frequency and random words that were randomly picked from the database.

TABLE I. RESULT OF AVERAGE SCORE

List	Word Group	Score	SD
Video Lecture A	High Frequency Words	86.857	8.713
	Random Words	29.714	10.991

Video Lecture B	High Frequency Words	89.142	10.766
	Random Words	38.857	10.804

High frequency words were evaluated as important in both video lectures (A and B) with resulting scores that were at least 80 points higher. On the other hand, random words, with a score of 30 points, were evaluated as not important. The relatively lower score is due to random words having a very low relationship with the video lecture. Table 2 provides the level of similarity between the word groups in a vector space.

TABLE II. RESULT OF SIMILARITY SCORE

Word Group	Video Lecture A	Video Lecture B
Keywords – High Frequency Words	1.96525	1.60309
Keywords – Random Words	2.72769	2.50064

The average similarity vector value between the keywords that experts selected and high frequency words was between 1 and 2, which indicates that the values are very similar to each other but not completely identical. On the other hand, the average similarity vector value between the keyword and random word groups was 2 or greater, which confirms that random words are not as similar to the keywords as were the high frequency words. Since the similarity between the vector values of each word in the vector space is a relative value, it is difficult to make an absolute comparison.

While investigating the scores of importance for high frequency words, it was discovered that high frequency words can be used as a feature to detect mind wandering; although, they might be insufficient as absolute keywords that represent the online software engineering video lecture. Figure 4 displays the similarity between the word groups in the vector space.

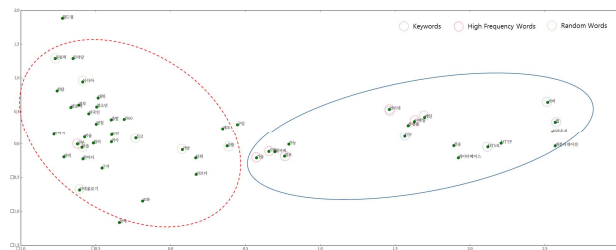


Figure 4. The vector space of word groups  
Blue line encircles keywords and high frequency words, while red dotted line encircles random words

### B. The Number of Words in the Word Game

To detect mind wandering, high frequency and random words were used in a word game. An experiment was conducted to determine the optimal number of words to include in the word game. Mind wandering does not detect how much learners understand the online video lecture;

however, it detects whether learners allow their minds to wander or if they attended to the video lecture. When learners paid attention, they could naturally remember high frequency words in the video lecture. According to a previous study, people store information they want to remember temporarily in a working memory. Typically, the working memory can store 7 words [8]. Based on conventional theory, seven words should be used in the word game. To see if that number of words is appropriate, a verification experiment was conducted. A comparative analysis was made to see if the detection of mind wandering varies according to the number of words. Figure 5 shows the experimental process.

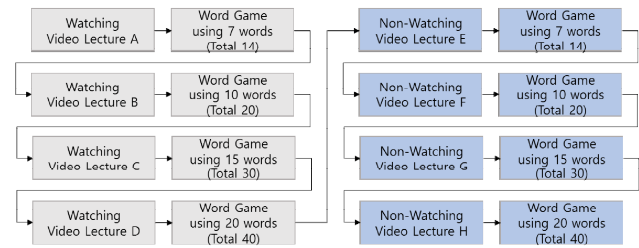


Figure 5. The experimental process for the number of words

We used a total 8 software engineering video lectures. Video lectures were classified as watched and unwatched and then 7, 10, 15 and 20 high frequency words were abstracted from the video lectures for the word game. Fourteen college students participated in the experiment. Participants watched the online video lectures and played the word game 8 times. Four word games were played after they watched video lectures. Afterwards, participants played another four word games without watching the video lectures. They checked the titles of the video lectures before playing the word game.

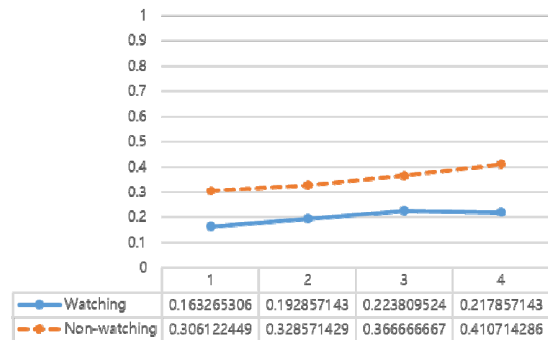


Figure 6. The result of average rate of wrong answers  
At x axis, "1" shows the result when 14 words were used, "2" for 20 words, "3" for 30 words and "4" for 40 words. The y-axis represents percentages.

The length of all video lectures was about 10 minutes. We selected 10-minute video lectures based on the results of a previous study that reported that 5 to 10 minutes of online video lectures are the most efficient in terms of learners' concentration levels [9]. It is also expected that 10-minute video lectures are primarily used in most online education

programs. Figure 6 shows the average rate of wrong answers by the number of words.

Although the rate of incorrect answers increased as the number of words in the word game increased, it was almost constant compared to the increase in the number of words. In other words, although the number of words increased, it only affected the correct versus incorrect answer rate. Accordingly, the use of an increased number of words does not have a significant advantage and only increases the cognitive burden and time consumption of learners to play the word game.

Figure 7 demonstrates that the wrong answer rate stays constant and the correct answer rate only increased according to the increase in the number of words in the word game that was completed after watching the video lectures. Therefore, the difference between the correct and incorrect answers increased as the number of words increased. On the other hand, the difference between correct and incorrect answers did not change in the word game for the video lectures that were not watched even when the number of words increased. This outcome is due to the similar number of correct and incorrect answers in the word game because they guessed the correct answers in the word game since they did not watch the video lectures. The probability that their response was the correct answer was 50% even without watching the video lectures. Since the wrong answer rate did not significantly change when the number of words increased, we concluded that the number of words does not have an effect on the detection of mind wandering. Conclusively, it was confirmed that the optimal number of words in a word game is 7 considering the cognitive burden and time consumption of learners; although, an increased number of words may slightly improve the detection accuracy.

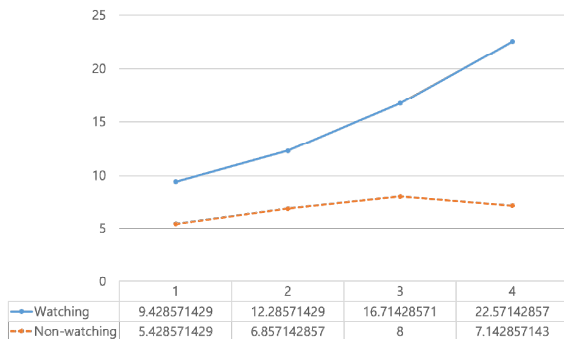


Figure 7. The average gap between the correct and wrong answers. On the x-axis, "1" signifies the result when 14 words were used, "2" for 20 words, "3" for 30 words and "4" for 40 words. The y-axis represents percentages.

### C. Video Lecture Length

All of the previous experiments used video clips that were approximately 10 minutes in length. Although most online video lectures are about 10 minutes long, it may be worthwhile to verify if mind wandering detection standards change according to the length of the video lectures. According to forgetting curve by Ebbinghaus that hypothesized the decline of memory retention over time [10], forgetting can happen with word memories since learning

takes longer when the length of the video increases. In other words, if we apply the forgetting curve theory to this experiment, an 80% correct answer rate is expected with a 10-minute video lecture, 70% with a 20-minute lecture, and 50% with a 60-minute lecture. However, the forgetting curve theory by Ebbinghaus represents the forgetting rate of meaningless words. Therefore, there may be a difference in the forgetting rate of high frequency words. Thus, we performed an experiment to examine if the criteria for detecting mind wandering changed according to the length of video lectures.

For this experiment, we selected 4 software engineering video lectures that were 10, 20, 30 and 60 minutes in length. After watching those video lectures, participants were asked to play word games. Since the average length of video-based online lectures is 10 to 20 minutes, videos longer than 60 minutes were inappropriate for the purpose of this study. Therefore, videos up to 60 minutes in length were included in this study. A total of 10 college students participated. Table 3 provides the results of this experiment.

TABLE III. RESULT OF WORD GAME

Video Lecture Length	Incorrect Avg.	SD
10 minutes	1.5	1.118
20 minutes	1.4	0.917
30 minutes	1.7	0.900
60 minutes	1.4	1.020

The results demonstrate that the wrong answer rate (forgetting rate) is constant regardless of the video length. High frequency words abstracted from the video lectures were meaningful so they did not affect the forgetting rate. According to Jenkins and Dallenbach, forgetting does not solely depend upon the simple lapse of time; it is also affected by the amount of experiences during the retention period [11]. Accordingly, it seems that the wrong answer rate is constant in word games in this system because the word games were performed right after watching the video lectures and there were relatively few experiences during the retention period. Additionally, according to cue-dependent forgetting, forgetting happens not in memory but by failure to recall information during the retrieval stage; therefore, if cues are provided, learners experience better recall than without cues [12]. The word game detects if the presented word is used in the video lecture or not. In other words, words are presented to the learners' memories and retrieval is dependent on the cue regardless of the length of the video lecture. This finding explains the constant wrong answer rate (forgetting rate). Conclusively, the length of the video lecture does not affect the detection of mind wandering through the use of high frequency words.

### IV. CONCLUSION

With the rapid growth of the e-learning industry and popularity of Massive Open Online Courses (MOOC), video-based online lectures are constantly increasing every year. Since the interaction between a teacher and a learner is

difficult and learners should drive their own learning in a video-based online lecture, it is very important for a teacher to know if learners have actually learned. This paper performed a verification experiment to easily detect mind wandering called 'minimum learning' that is unlike conventional methods. If a learner has watched a video lecture without mind wandering, it means that minimum learning occurred.

The minimum learning judgement system detects mind wandering through the high frequency words used in video lectures; so, diversified studies on high frequency words are required. According to the results, the first research question addressed the importance of high frequency words. The results demonstrated that while high frequency words do not exactly match the keywords that represent the software engineering video lectures, they can be used as features to identify mind wandering and determine minimum learning. The second research question, which was related to the incorrect answer rate and the number of words used in a word game, was addressed. Increasing the number of words makes the task more time consuming and represents a burden to cognitive processing; therefore, seven was chosen as the appropriate number of words for the task. The third research question addressed how the incorrect answer rate was affected by the duration of the software engineering video lecture. The results demonstrated that the length of the video lecture did not affect the assessment of minimum learning.

It is confirmed that high frequency words can be used as a feature to detect mind wandering. Additionally, the result of this study shows that high frequency words can be used for studies on automatic quiz generation, automatic learning detection and concentration improvement.

#### ACKNOWLEDGMENT

The National Research Foundation of Korea (NRF) grant funded by the Korea government (No. 2016R1A2B 2015912).

#### REFERENCES

- [1] Ministry of trade industry & energy, and National IT industry promotion agency, "2015 survey of korea e-learning industry," April 2016.
- [2] Joohye Kang, Minjea Park, Yujung Yun, Hoyang Choi, Seongwon Park, and Kwangsu Cho, "Study on e-Learner's Attention improvements," Korean HCI Association Conference 2014, pp. 353-356, 2014.
- [3] Evan F. Risko, Dawn Buchanan, Srdan Medimorec, and Alan Kingstone, "Everyday attention: Mind wandering and computer use during lectures. *Computers & Education*, vol. 68, pp. 275-283, October 2013.
- [4] Faber, Myrthe, Robert Bixler, and Sidney K. D'Mello, "An automated behavioral measure of mind wandering during computerized reading," *Behavior Research Methods*, pp. 1-17, February 2017.
- [5] Jaechoon Jo, and Heuseok Lim, "How to Judge Learning on Online: Minimum Learning Judgment System", *The 9th International Conference on Educational Data Mining*, pp. 597-598, 2016.
- [6] Savin, H. B., "Word-frequency effect and errors in the perception of speech," *Journal of the Acoustical Society of America*, vol. 35, pp. 200-206, 1963.
- [7] Rehurek, R., and Sojka, P., "Software framework for topic modelling with large corpora," In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pp. 46-50, 2010.
- [8] Miller GA., "The magical number seven plus or minus two: Some limits on our capacity for processing information," *Psychological Review*, vol. 63, no. 2, pp. 81-97, March 1956.
- [9] Guo, P. J., Kim, J., and Rubin, R., "How video production affects student engagement: An empirical study of mooc videos," In *Proceedings of the first ACM conference on Learning@ scale conference*, pp. 41-50, March 2014.
- [10] Ebbinghaus, H., "Memory: A contribution to experimental psychology," *Annals of neurosciences*, vol. 20, no. 4, pp. 155-156, 2013.
- [11] Jenkins, J. G., and Dallenbach, K. M., "Obliviscence during sleep and waking," *Journal Psychology*, vol. 35, pp. 605-612, 1924.
- [12] Tulving, E., "Cue-dependent forgetting," *American Scientist*, vol. 62, pp. 74-82, 1974.