

Violent Agreement Amongst Systems and Software Engineers, CSER 2008

DeWitt T. Latimer IV¹

University of Southern California
941 W. 37th Place, SAL 300, Los Angeles, CA 90089-0781
dlatimer@ieee.org

¹The views expressed in this report are those of the author and do not reflect the official policy or position of the United States Air Force, Department of Defense, or the U.S. Government.

Abstract

Recent discussions of how to integrate systems and software engineers begs the question to define the difference between these two disciplines. The disciplines of electrical, mechanical, and civil engineering do not seem to be engaged in this discussion. By performing some basic lexical analysis between the published bodies of knowledge, we observe that software engineers and systems engineers use similar, but not identical concepts to describe their discipline of expertise. Through lexical analysis, systems and software engineering seem more similar to each other than to civil engineering. In the end, the differences in systems and software engineering are concluded to be one where the two disciplines are in violent agreement on the principles, but differ in exact utilization of concepts.

Introduction

The current discussion about needing to reconcile systems and software engineering (as stated in (Baldwin 2007)) brings to mind the question about what needs to be reconciled. At a high level, the CMMI (CMMI Product Team 2006) integrated the practices of systems and software engineering into a common process model. However, (Baldwin 2007) states "software continues to be separate from the systems engineering practices." How can both states seem to exist? The difference may be in how systems and software engineers perceive similar activities.

A difference in perception by the engineers in how to deal with issues like "requirements" (which (Baldwin 2007) notes poor traceability to system requirements) may lead to an inability of the disciplines to form a joint plan to address the root issue. If a systems engineer's concept of requirements is sufficiently different from a software engineer's, then the two may have problems creating solutions. Indeed, such difference may become the focus of a "violent agreement" as to the nature of the problem or potential solutions.

A violent agreement is a popular culture phrase which indicates that there is some agreement about the goal or objective, but the participants cannot realize they agree because they differ strongly in terminology or approach (Oxymorons 2008 and Urban Dictionary 2008). This may be the case when systems and software engineers agree there is a problem in the traceability of requirements, but do not realize their perception of a requirement is too different.

The goal of this paper is to examine common and divergent vocabulary between systems and software engineering, as expressed by their consensus bodies of knowledge. This examination

will help clarify what needs to be integrated between the disciplines. Further, by examining words in common to the two disciplines, this paper identifies how underlying usage differences may contribute to the lack of integration between the disciplines. Ultimately this paper cannot conclude which discipline is "more correct", but the specific findings may be useful in addressing future engineering integration discussions.

Related Work

Engineering Knowledge

A body of knowledge is an important part of being a professional discipline (Abran 2004). Some engineering professional societies have created handbooks, bodies of knowledge, or guides to their bodies of knowledge in an effort to advocate to members and non-members the role of their discipline.

Three engineering professional societies publish bodies of knowledge for their engineering practice. The International Council on Systems Engineering (INCOSE) publishes a Systems Engineering Handbook (SEH) (Haskins 2006). The IEEE Computer Society publishes a guide to the Software Engineering Body of Knowledge (SWEBOK)(Abran 2004). To offer a contrasting engineering view, we will consider the current draft of the body of knowledge of civil engineering (ASCE 2007) of the American Society of Civil Engineers (ASCE). Mechanical, electrical, and other engineering disciplines do not have publicly available bodies of knowledge.

Lexical Analysis

One of the simplest forms of lexical analysis is word frequency. Numerous plagiarism detection tools use word frequency analysis to queue deeper analysis or inform plagiarism rating functions (Lukashenko 2007). The idea is that if words are occurring with similar frequencies, then the texts may be similar. When performing lexical analysis, potential words in the document are called "tokens" before human judgement is applied (to ensure that parser errors do not create non-word tokens). In a similar fashion, "type" is used to connote a class to which a token belongs. In doing frequency analysis, the type of a token is a token, but not in the body. For example, in the sentence "the dog enters the room" there are 5 tokens, but only 4 types ("the" appears twice as a token).

Zipf's law (Li 1992) indicates that the frequency of a word is inversely related to the word's rank. The usage of words falls off on an inverse exponential scale. For example, all words considered, in SWEBOK, "the" appears 4,417 times, or slightly less than 1.5 times as often as the next occurring words, "software" (2,966), "and" (2,874), and "of" (2,856). By the tenth word in frequency order, "IEEE" only occurs one-fifth as often with 861 occurrences. In a text of 84,221 words, there are 6,603 different words in the document; however, fifty percent of the text consists of only 80 of those words.

Initially, contextual indices have been used to build concordances for topics as diverse as interpreting the dead sea scrolls (Abegg 2003) or researching the differences between two words. (Sweterlitsch) used near neighbour word information and relationships to differentiate between the usage of "someone" and "somebody" when the words are used in a spoken English corpus. This notion of looking at the neighbouring words of a key word is called "key word in context". Using a concordance, built from examining the key words in context, researchers have been able to determine the meanings and purposes of words in the dead sea scrolls.

The results in this paper were achieved through the use of two principal tools, both freely available. The Simple Concordance Program (Reed 2007) and OpenOffice Calc (an open source spreadsheet application, similar to Excel).

Building the Token Data Sets

First Impression of the BOKs

INCOSE's Guide to the System Engineering Body of Knowledge consists of 5,219 different tokens types with 57,384 tokens. This gives us 0.091 types per token. The work is arranged in 185 pages for an average of 310 tokens per page.

The IEEE Computer Society's Guide to the Software Engineering Body of Knowledge consists of 6,603 different token types with 84,221 tokens. This gives us 0.078 types per token. This work is arranged in 202 pages for an average of 417 tokens per page.

The ASCE Civil Engineering Body of Knowledge consists of 4,451 different token types with 49,170 tokens. This gives us 0.091 types per token. The work is arranged in 191 pages for an average of 257 tokens per page.

Preparing BOKs

The BOKs, however, may be difficult to evaluate in an automatic fashion lexically without some basic preparation. To enable full lexical counts, abbreviations and acronyms need to be expanded and page/document formatting text needs to be examined or cut.

Abbreviations are frequently used in engineering (such as "CM" for "Configuration Management"). Since this study is examining both the frequency and the concordance of tokens, the expansion of acronyms ensure that words like "configuration" are appropriately weighted in the total lexical counts and also that its concordance with words like "management" is noted appropriately. Some abbreviations, however, are kept for clarity. Abbreviations that are names of organizations, product names (such as CMMI), Y2K, or used in document navigation (such as citation abbreviations) are not changed, as such abbreviations are themselves tokens in the body of knowledge.

However, not all words in the documents are equal, for example, SEH uses extensive header and footer information to indicate copyright, document title, and document version information. Such information is replicated on every page, but does not contribute to the understanding of the body of knowledge. Such page headers and footers are cut from examination.

All other information, forwards, thanks, annexes, et. al. are kept despite potentially limited importance to the body of knowledge. Such removals would be much more subjective than the automatic removals above and may accidentally cut some portion of the body of knowledge.

Word Exclusions and Mergers

To determine the important concepts of the body of knowledge, some words are systematically excluded or merged with other words for frequency analysis. Excluded words include: words with less than 10 occurrences, articles, pronouns, prepositions, conjunctions, forms of the verb "to be", forms of the verb "to have", forms of the verb "to be able", single letters, document navigation words ("chapter", "section", "title", "figure", and "table" being the most common), "who", "what", "where", "when", "why", "how", "http", "com", "www", "so", "also", "that", "must", "most", "many", "there", "those", "this", and "may". Plural words are

merged with the singular form. Adjective or adverb forms are left as they are (although future work may seek to do a full stem merging, the act presently requires too many situational judgements to effectively perform automatically, e.g. "use" and "used" cannot be merged when discussing "use cases").

Although these words are excluded for frequency analysis, they will all come back into play when examining concordance of selected tokens.

The "Massaged" Bodies of Knowledge

With the above modifications, the documents appear slightly different. The new token, token type, and type/token ratios are below.

INCOSE's Guide to the System Engineering Body of Knowledge consists of 612 different tokens types with 23,185 tokens. This gives us 0.026 types per token.

The IEEE Computer Society's Guide to the Software Engineering Body of Knowledge consists of 817 different token types with 38,614 tokens. This gives us 0.021 types per token.

The ASCE Civil Engineering Body of Knowledge consists of 537 different token types with 19,873 tokens. This gives us 0.027 types per token.

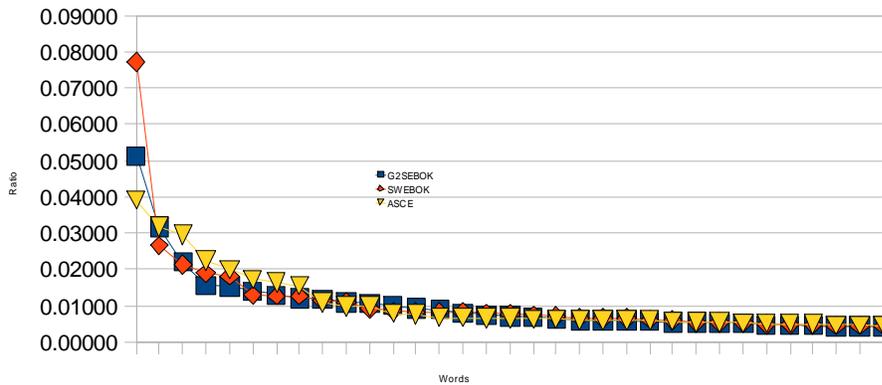


FIGURE 1: Graph of word frequency in SEH, SWEBOK, and ASCE

frequency of word usage that is fairly similar between the three documents.

Figure 1 illustrates the token ratio for the top 40 tokens. Beyond the 40th token, the tail continues the same gentle slope downward. As expected from (Li 1992), these texts exhibit regular Zipf-like properties regarding the decreasing

Lexical Analysis

Comparison of Token Frequencies

With the new token frequencies, by examining the curvature of the resulting graphs in figure 1, the top 10 tokens in all three bodies of knowledge do appear to be the most frequently occurring words. Further, Table 1 lists all the words and which bodies of knowledge they share in each other's top 10. The only word common to all three BOKs is "engineering", while SWEBOK and ASCE also share "knowledge". SEH and SWEBOK share 3 additional top words: "process", "requirement", and "management".

TABLE 1: Top ten common word tokens in SEH, SWEBOK, and ASCE

<i>Top 10 Words for Each Body of Knowledge</i>			
Rank	SEH	SWEBOK	ASCE
1	system	software	<i>engineering</i>
2	process	process	<i>knowledge</i>
3	project	<i>engineering</i>	civil
4	requirement	requirement	body
5	management	management	outcome
6	risk	quality	professional
7	<i>engineering</i>	<i>knowledge</i>	level
8	enterprise	standard	engineer
9	life	testing	committee
10	cycle	configuration	problem

Bold Italics are common to all
Italics are common to SWEBOK-ASCE
Bold are common to SEH-SWEBOK
 No tokens common to SEH-ASCE
 (in the top 10)

Looking further, Table 2 lists all the words which are common in the top 0.2% of each body of knowledge. Observe there are 20 words that are common between all three bodies of knowledge. In addition, there are 21 words in common between SEH and SWEBOK, 12 words in common between SEH and ASCE, and 12 words in common between SWEBOK and ASCE.

TABLE 2: Common word tokens between SEH, SWEBOK, and ASCE

<i>Common Tokens Appearing More than 0.2% in the document</i>			
Common To All (20)	SEH-SWEBOK (22)	SEH-ASCE (12)	SWEBOK-ASCE (12)
system	requirement	technical	knowledge
process	risk	policy	area
project	cycle	new	related
management	activity	criteria	practice
engineering	cost	work	program
life	analysis	required	technology
design	product	appropriate	society
development	support	communication	edition
information	quality	team	defined
other	control	material	issues
organization	change	impact	resource
model	verification	study	body
level	maintenance		
need	standard		
engineer	configuration		
example	data		
plan	approach		
problem	document		
based	planning		
provide	review		
	user		
	item		

Of note is that SWEBOK and ASCE both use the same knowledge taxonomy to describe their desired level of understanding for items within their body of knowledge that SEH doesn't use. This explains the presence of the words "knowledge" and "area" on the SWEBOK-ASCE common list and their absence from the SEH.

Comparison of Lexical Concordances

The words "requirement" and "management" are selected for examination for concordance differences in SEH and SWEBOK. "Requirement" is selected as an engineering term of significance in both disciplines. The hope is that "requirement" is used in a similar fashion between the disciplines to indicate that the concept is used in a similar fashion. "Management" is selected because both bodies of knowledge extensively reference the PMBOK, so again, we would expect to see the term being used in a consistent fashion. Although "process" ranks high in both, "process" is used considerably more frequently than the other two words (by 30% and 50% respectively) and is a word common to many other disciplines.

The goal isn't to determine the "correct" concordance, but to determine if there is a difference in the usage of words in the two disciplines.

To facilitate this analysis, references to the words are excluded if they appear to be part of the formatting of the document or are the names of conferences (this is done specifically for SWEBOK, which extensively cites conference papers), as these references do not add to the understanding of how the word is put into active use. Titles of papers are maintained for

concordances, as the paper titles convey specific nuanced understandings of the perception and understanding of the word. Further, we ignore preposition and transitive verbs to find the nearest left and right neighbour. Finally, we only consider concordances within the same sentence or bullet, thus if the word starts a sentence or bullet, its left association will be "START"; similarly, if the word ends a sentence or bullet, its right association will be "STOP" (neither of these words appear as left or right associations in the documents for "requirement" or "management").

For each word, we compare the frequency of the neighbouring words that occur more than 1% of the time in SEH and the corresponding frequency of that same word in SWEBOK to see if the words are drawn from the same underlying population. A paired, single-tail, t-test is then performed. A score of below 0.05 will indicate that we can reject the hypothesis that left and/or right words are drawn from the same distribution. A score below 0.20 indicates evidence towards rejecting the similar distribution hypothesis, but the result is partially inconclusive. Rejecting the similar distribution hypothesis in this second case may cause what is called a "type II" statistical error (meaning that we falsely call two distributions different when they are the same). Scores over 0.80 indicate that there is some evidence that the distributions are the same (but may result in a type I error, marking distributions as similar when they are not) and scores over 0.95 indicate that the similarity hypothesis can be accepted.

"Requirement":

appears in SEH 362 times (the plural form appearing 613 times, or 93.1%) and in SWEBOK 658 times (the plural form appearing 347 times or, 95.8%). Table 3 lists the top left and right concordances (those occurring more than 1%) for the two works. Taking the top 14 left concordances from SEH (those occurring more than 1%, but excluding "the" and "a"), the t-test score is 0.0041. This score

TABLE 3: Frequency of words appearing to the left and right of "requirement" in SEH and SWEBOK

WORD	LEFT		WORD	RIGHT	
	SEH	SWEBOK		SEH	SWEBOK
the	0.0994	0.1687	STOP	0.2431	0.1778
stakeholder	0.0829	0.0046	analysis	0.0387	0.0228
START	0.0580	0.0243	system	0.0387	0.0046
performance	0.0442	0.0030	definition	0.0331	0.0030
stakeholders	0.0249	0.0000	management	0.0249	0.0000
design	0.0221	0.0061	verification	0.0249	0.0000
Non-functional	0.0193	0.0258	elicitation	0.0138	0.0137
compliance	0.0138	0.0000	architectural	0.0110	0.0015
maintenance	0.0138	0.0015	specifications	0.0110	0.0699
derived	0.0138	0.0015	validation	0.0110	0.0152
interface	0.0138	0.0015			
resources	0.0110	0.0000			
technical	0.0110	0.0000			
functional	0.0110	0.0076			
user	0.0110	0.0198			
a	0.0110	0.0502			

indicates that we can reject the similar distribution hypothesis, i.e. the words to the left of "requirement" come from the different distributions. For the right concordance, taking the top 10 words, the t-test score is 0.099, which mildly supports rejecting the similarity hypothesis. If the two distributions are combined (left and right distributions in to one list), the resulting t-test score is 0.0025. We can say that there is a good probability that these two bodies of knowledge use the word "requirement" in different contexts.

"Management": appears in SEH 353 times and in SWEBOK 700 times (note that the exclusion considerations for this section are different from the token frequency section). Neither document uses the plural of management. Table 4 lists the top left and right concordances for the two works, again those words occurring more than 1%. performing the same analysis as in requirements ends up with t-test scores of 0.285 for the left concordance (15 words), 0.147 for the right concordance (16 words), and 0.168 for the combined list. While these results are not

conclusive, the fact that they do not support the similarity hypothesis is unexpected, as both SEH and SWEBOK extensively cite PMBOK as their source for management practices. This may mean that the two disciplines are not applying the PMBOK in similar fashions, or that the purpose or function of management may be slightly different between the disciplines.

TABLE 4: Frequency of words appearing to the left and right of "management" in SEH and SWEBOK

	LEFT			RIGHT	
WORD	SEH	SWEBOK	WORD	SEH	SWEBOK
risk	0.1331	0.0429	STOP	0.2068	0.2171
quality	0.1303	0.1071	process	0.1813	0.0614
configuration	0.1246	0.3614	plan	0.0567	0.0371
information	0.0992	0.0043	processes	0.0227	0.0314
project	0.0992	0.0900	information	0.0198	0.0071
resource	0.0652	0.0043	systems	0.0170	0.0214
investment	0.0538	0.0000	project	0.0142	0.0143
processes	0.0397	0.0000	guidelines	0.0142	0.0000
requirements	0.0283	0.0000	institute	0.0142	0.0114
engineering	0.0255	0.0771	quality	0.0142	0.0029
environment	0.0255	0.0000	customer	0.0113	0.0014
enterprise	0.0227	0.0014	system	0.0113	0.0043
START	0.0170	0.0286	activities	0.0113	0.0329
decision	0.0170	0.0000	plans	0.0113	0.0100
business	0.0113	0.0000	tools	0.0113	0.0357
			resource	0.0113	0.0000

Observations and Future Work

At a high level, we can see that top list of top words are more similar between SEH and the SWEBOK than either of them are to ASCE. This is likely due to the acceptance of similar philosophies in the management and organisation of their practices, or possibly due to the influence of the CMMI or other engineering process research efforts. Likely, system and software engineers use very similar words in the generation of their documents.

However, this similarity ends at the token frequency level. When examining the context in which two words were used, there is a surprising level of difference. For the term "requirements", we note that SWEBOK is about six times more likely to have the word associated with "specification" than in SEH. SWEBOK talks about "requirements engineering", a term not used in SEH, which uses "requirements management" (conversely, not used in SWEBOK). SWEBOK references "quality requirements" four times more often than SEH. These differences are among those that lead to the statistical difference in the concordance for "requirements" in the two documents. This evidence may indicate that when a system and software engineer are talking requirements, the related associations may be very different. This can lead to the engineers talking at cross-purposes.

With "management", the fact that the term is not used in a similar fashion in the documents is slightly disturbing, as the practice of management is cited to a common, authoritative source by both documents. SWEBOK is three times more likely to refer to "management activities", whereas SEH is six times more likely to refer to "management process". SWEBOK is eight times more likely to use the phrase "management techniques", whereas SEH uses "management guidelines" (not appearing in SWEBOK). This difference appears to have some grounding in whether the respective disciplines look at management more as art ("activities" and "techniques") or a science/engineering domain ("process" and "guidelines").

Overall, there appears to be evidence that systems and software engineering share a common vocabulary, but use the words in different contexts. This supports the anecdotal observation that system and software engineers tend to be arguing about semantics, or as stated previously, "violent agreement". The engineers can agree what to call something, but not that they are meaning the same thing.

Future work should focus on the professional engineering licensure examination guides for different disciplines. This additional work can further clarify how engineering disciplines are expected to vary between each other and where similarities are expected. This work only used one engineering discipline, civil, as a control. Another comparison could be to compare certification/licensure specification documents (typically 2-4 pages) in a more qualitative fashion, to examine the difference and similarities in expected skills of engineers.

Another point of comparison that may yield guidance addressing issues in lack of similarity may be to involve the management bodies of knowledge. As mentioned, the PMBOK is cited by both SEH and SWEBOK. Additionally, ASME has a body of knowledge for engineering management (ASME 2007). The comparison of these works in this context may bring into focus the portion of the engineering practice in the bodies of knowledge that is engineering management.

Finally, given the heritage, comparison of the SEH and SWEBOK with the CMMI (which previously reconciled different systems and software engineering process improvement guides) or another engineering process improvement guidance (from the ISO 9000-family) may yield insight into how the disciplines organize their practices.

References

- Abegg, M., et. al, *The Dead Sea Scrolls Concordance, Volume 1*, Brill, ISBN 9004125213, 2003.
- Abran, A., et. al. (eds), *Guide to the Software Engineering Body of Knowledge, 2004 Version*, IEEE Computer Society, Washington, DC, 2004.
- ASME, *A Guide to the Engineering Management Certification -- Body of Knowledge*, ASME Press, New York, NY, USA, ISBN: 0791802477, 2006.
- Baldwin, Kristin, "Software and Systems Integration: Sponsor Overview". The 22nd International Annual Forum on COCOMO and Systems/Software Cost Modeling, October, 2007.
- Body of Knowledge Committee (ASCE), *Civil Engineering Body of Knowledge for the 21st Century: Preparing the Civil Engineer for the Future, Second Edition*", American Society for Civil Engineers, draft 10, dated 15 July 2007.
- CMMI Product Team, "CMMI for Development, Version 1.2", Carnegie Mellon University, Software Engineering Institute, Technical Report, CMU/SEI-2006-TR-008, Pittsburgh, PA, USA, August 2008.
- Haskins, C. (ed.), *Systems Engineering Handbook: A guide for System Life Cycle Processes and Activities, Version 3*, International Council on Systems Engineering, INCOSE-TP-2003-002-03, June 2006.
- Li, W., "Random Texts Exhibit Zipf's-Law-Like Word Frequency Distribution". *IEEE Transactions on Information Theory*, volume 38, issue 6, pp. 1842–1845, 1992.
- Lukashenko, R., et. al., "Computer-Based Plagiarism Detection Methods and Tools: An Overview", *International Conference on Computer Systems and Technologies*, Bulgaria, 14-15 June 2007.

"Oxymoron Forum: Violent Agreement", <http://www.oxymorons.info/reference/oxymorons/forum-m.asp?oxyid=143&msgid=2298>, visited on 1 Feb 2008.

Reed, A. *Simple Concordance Program*, <http://www.textworld.com/scp>, visited on 1 Feb 2008, 2007

Sweterlitsch, T., "Someone vs. Somebody: Word/Phrase Pairs Exposed", <http://www.andrew.cmu.edu/course/76-451/tom.html>, visited on 1 Feb 2008, no date given.

"Urban Dictionary: violent agreement", <http://www.urbandictionary.com/define.php?term=violent+agreement>, visited on 1 Feb 2008.

Biography

Captain DeWitt Latimer IV, USAF earned his PhD in Computer Science at the University of Southern California. He is advised by Prof Barry Boehm and Prof Gaurav Sukhatme. His research focuses on investigating the nature of acquiring autonomous robotic systems. He earned his MS degrees in Robotics (2001) and Civil Engineering (2002) at Carnegie Mellon University. He is a senior member of the IEEE and ACM, a member of INCOSE, ASCE, and AFCEA, and was awarded the CSDP credentials from the Computer Society.