

# Equality of cumulative votes

K. Rinkevičs<sup>a</sup>, R. Torkar<sup>a,b</sup>

<sup>a</sup>*Blekinge Institute of Technology, Sweden*

<sup>b</sup>*Simula Research Laboratory, Norway*

---

## Abstract

**Context.** Prioritization is an essential part of requirements engineering, software release planning and many other software engineering disciplines. Cumulative Voting (CV) is known as a relatively simple method for prioritizing requirements on a ratio scale. Historically, CV has been applied in decision-making in government elections, corporate governance, and forestry. However, CV prioritization results are of a special type of data—compositional data.

**Objectives.** The purpose of this study is to aid decision-making by collecting knowledge on the empirical use of CV and develop a method for detecting prioritization items with equal priority.

**Methods.** We present a systematic literature review of CV and CV analysis methods. The review is based on searching electronic databases and snowball sampling of the found primary studies. Relevant studies are selected based on titles, abstracts, and full text inspection. Additionally, we propose Equality of Cumulative Votes (ECV)—a CV result analysis method that identifies prioritization items with equal priority.

**Results.** CV has been used in not only requirements prioritization and release planning but also in e.g. software process improvement, change impact analysis and model driven software development. The review presents a collection of state of the practice studies and CV result analysis methods. In the end, ECV was applied to 27 prioritization cases from 14 studies and identified nine groups of equal items in three studies.

**Conclusions.** We believe that the analysis of the collected studies and the CV result analysis methods can help in the adoption of CV prioritization method. The evaluation of ECV indicates that it is able to detect prioritization items with equal priority and thus provide the practitioner with a more fine-grained analysis.

## 1. Introduction

Software products are becoming larger and more complex. Each product is usually affected by a large number of factors such as functional requirements, quality attributes, or software process improvement issues. Since time, funding, and resources are limited, it is seldom possible or even desirable to fully address all the factors. Therefore, the level of attention to a particular factor must be decided according to its importance (e.g. business value), cost, risk, volatility, dependencies between the factors and other such criteria. These type of decisions are made by product stakeholders: users, clients, managers, sponsors, developers, and other persons associated with the product. In order to make decisions regarding a large number of factors it is highly advisable to prioritize the factors in a systematic way [1].

One of the prioritization methods used in software engineering is Cumulative Voting (CV) [2]. The main advantage of CV is that it is relatively simple and fast, yet produces priorities in ratio scale [1, 3]. This allows us to not only determine what prioritization items are more important but also how much more important they are. (Ratio scale prioritization is particularly important in software release planning and cost-value analysis [4, 5].)

Prioritization is usually performed by multiple stakeholders where individual priorities are combined into a single priority list. Each stakeholder's preferences may have different weight in the final priority. Such prioritization provides more information than just the priorities of factors. In the end, it may be useful to analyze the results of the prioritization to assess disagreement between stakeholders, measure stakeholder satisfaction with the results or find distinct groups of stakeholders.

The purpose of this study is to help industry practitioners and academia researchers in adopting, using and developing CV, while the importance of prioritization in software engineering and the prospectiveness of CV constitutes a need to do further research in this area.

This study presents a systematic literature review on the empirical use of CV and CV result analysis methods. A new method for CV result analysis, called Equality of Cumulative Votes (ECV), is proposed. The method identifies prioritization items with *equal* priority. ECV is evaluated using a considerable amount of data, which was obtained from the primary studies

35 identified by the systematic review (through the kindness of the authors of  
36 said studies).

## 37 **2. Background**

38 This section presents definitions and places this study in a context. In the  
39 coming sections we will cover: a description of software requirements priori-  
40 tization methods; examples of CV result analysis methods; and a description  
41 of compositional data analysis and CV.

### 42 *2.1. Prioritization Methods*

43 Some of the most popular prioritization methods are the analytical hierar-  
44 chy process (AHP), cumulative voting (CV), ranking, numerical assignment,  
45 top-ten, the planning game, minimal spanning tree, bubble sort and binary  
46 search tree [1, 6]. Ranking and numerical assignment methods perform prior-  
47 itization on an ordinal scale. AHP and CV are, on the one hand, considered  
48 to be harder to use and also more time consuming compared to other methods  
49 but, on the other hand, produce priorities in ratio scale.

50 Ratio scale priorities have several advantages over ordinal scale priorities.  
51 Ratio scale shows not just the order of items but also relative distance be-  
52 tween them. This enables the priority of a group of items to be calculated  
53 by summing up the priorities of individual items [4]. It is possible to say  
54 that one item or set of items has higher priority than another set of items.  
55 Supposing stakeholders have to choose between several low priority items  
56 and one item with higher priority; with ordinal scale, the item with highest  
57 priority will always be selected first. However, if priorities are given on a  
58 ratio scale, it is possible that lower priority items will be selected if their  
59 cumulative priority is higher.

60 Finally, the ratio scale allows the combining of multiple priority factors  
61 by calculating ratios between them. One example of this is the cost-value  
62 ratio that shows which requirements give more value for less money [5].

### 63 *2.2. Prioritization Result Analysis*

64 Disagreement between stakeholders happens when two or more stakehold-  
65 ers have assigned a different priority to one prioritization item. If the level of  
66 disagreement is high it may indicate potential conflicts between stakeholders.  
67 Such conflicts may be of technical character, as well as social or cultural.

68 The satisfaction a stakeholder has with the final prioritization results is  
69 determined by the difference between the results and the individual priorities  
70 of the stakeholder. A smaller level of difference leads to higher satisfaction.  
71 In the end, stakeholder satisfaction is important because it is necessary to  
72 achieve stakeholder commitment.

73 In some cases a part of stakeholders may form a group of some kind  
74 and, therefore, prioritize requirements similarly. It may be useful to detect  
75 whether a group of stakeholders has different preferences compared to other  
76 stakeholders. As an example, in [7], domain experts, technical experts, man-  
77 agers, project managers, testers, and developers use CV to prioritize software  
78 process improvement issues and the CV results are analyzed using disagree-  
79 ment charts and satisfaction charts. Finally, principal component analysis  
80 (PCA) is used to identify distinct groups of stakeholders.

81 The same items can be prioritized by the same stakeholders multiple times  
82 from different perspectives. In this case it is useful to determine correlation  
83 between the priorities in different perspectives to assess the differences be-  
84 tween the perspectives. As an example, in [8], CV is used by developers,  
85 testers and managers to prioritize quality attributes. The same quality at-  
86 tributes are prioritized from two perspectives: the perceived situation today  
87 and the perceived ideal situation. Correlation between the two perspectives  
88 is evaluated using the Spearman rank correlation matrix. This allows an  
89 analysis of how well the company balances the priorities of software quality  
90 attributes.

91 In [9] change impact issues are prioritized by developers, testers, man-  
92 agers, and system architects. The prioritization is done with respect to three  
93 perspectives: strategic, tactical, and operative. In order to determine corre-  
94 lation between the perspectives, CV results are analyzed using the Kruskal-  
95 Wallis test. In [10] the results of [9] are further analyzed using PCA, bi-plot,  
96 and ternary plot. In this case, PCA is used to find correlated issues, bi-  
97 plot shows variance, correlation, difference between the priorities of issues,  
98 and the viewpoints of stakeholders, while ternary plots are used to show the  
99 relative number of issues that received high, medium, and low priority.

100 As can be seen above, from the examples above, prioritization has been  
101 performed with various stakeholders, using different perspectives and, in the  
102 end, also analyzed using various techniques. We will next describe in more  
103 detail one of the more common methods to manage prioritization issues—  
104 cumulative voting—which has been used in software engineering for some  
105 time. (CV has its roots in corporate governance and biology.)

### 106 2.3. Cumulative Voting

107 CV is a prioritization method for prioritizing a list of items [2]. CV  
108 has many synonyms in literature: hundred (100) dollar (\$) method/test and  
109 hundred (100) point method. Before being applied in software engineering  
110 CV was used for political elections [11] and corporate governance [12]. CV  
111 has also been applied in e.g. decision making in forestry [13], voting in social  
112 networks [14] and in computer algorithms for consensus clustering [15] (as a  
113 method for combining the results of different clustering algorithms).

114 In CV a stakeholder is given 100 points, imaginary dollars or units of  
115 percentages that can be spent on the prioritization items. In the simplest  
116 case, the stakeholder can spend any amount of points on any number of items  
117 as long as the total amount adds up to 100. The more points assigned to an  
118 item, the higher the priority of the item (and implicitly, the lower priority  
119 to the other items). The stakeholder may spend all points on just one item  
120 or distribute them among all or some of the items. Once again, this is the  
121 simplest case; other variants exist, which we will see next.

122 Often prioritization is done by more than one stakeholder. The final prior-  
123 ity of an item can be calculated by adding up the points each stakeholder has  
124 spent on it. Sometimes the vote of some stakeholders may be more important  
125 than the votes of others. For example, a manager may be more influential or  
126 shareholders may have different amount of shares. In such a case the prior-  
127 ities of each stakeholder may be multiplied by an individual coefficient or a  
128 stakeholder may be given a more points to perform the prioritization.

129 Worth mentioning in this context is that it is advisable to randomize the  
130 order of items in a prioritization list. This is necessary in order to minimize  
131 the effect of order on the prioritization results, which has shown to have an  
132 effect [16].

#### 133 2.3.1. Benefits and Drawbacks of Cumulative Voting

134 Compared to analytical hierarchy process (AHP), CV is faster and easier  
135 to learn and use [1, 3]. AHP benefits from consistency check, but CV does  
136 not require this because all prioritization items are evaluated simultaneously  
137 [3].

138 There are, however, a few problems with CV. First of all, it cannot be  
139 repeated for the same stakeholders and prioritization items due to stakeholder  
140 bias [2] (c.f. Section 2.3.4). Secondly, CV becomes more difficult to use when  
141 the number of prioritization items increases [17].

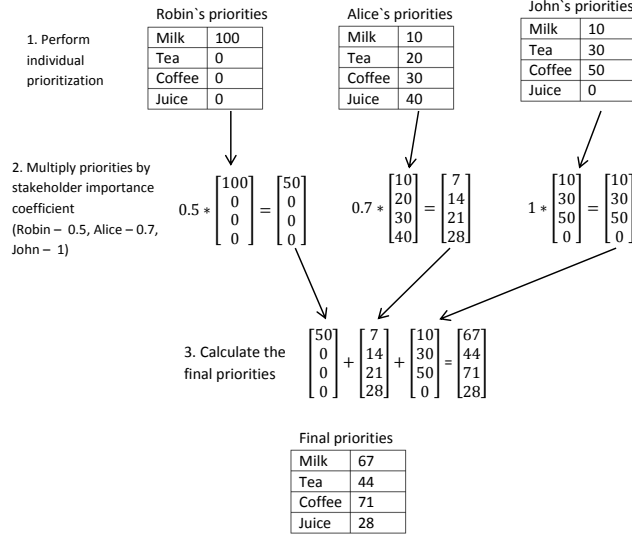


Figure 1: Example of CV with several stakeholders.

### 2.3.2. Example of Cumulative Voting with Several Stakeholders

Let us next give an example of CV with several stakeholders. Suppose Robin, Alice, and John are three friends who want to buy some beverages in a store. They have different preferences but do not want to buy too many drinks. Therefore, they decide to use CV to decide what to buy. Each of the friends distributes 100 points between four items: milk, tea, coffee, and juice (Step 1 in Figure 1). In this case each of them will spend a different amount of money on the purchase, hence, their priorities are multiplied by different coefficients (Step 2 and the stakeholder importance coefficient in Figure 1). The final beverage priorities are calculated by summing up the weighted priorities of stakeholders (Step 3 in Figure 1).

### 2.3.3. Stakeholder Bias

Prioritization using CV may be biased if a stakeholder knows the preferences of other stakeholders. She may manipulate the results by spending more points on items that are important to her but not to the other stakeholders. On the one hand, stakeholder bias makes it unreasonable to repeat CV with the same prioritization items and stakeholders. On the other hand, this property of CV may be useful in giving more power to important minority stakeholders, such as security experts or software testers. Suppose the

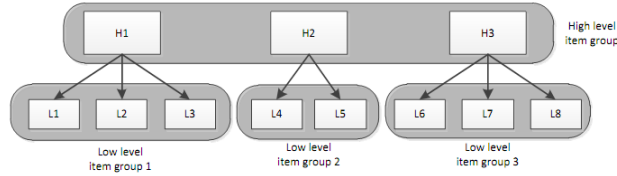


Figure 2: Example of prioritization item hierarchy.

161 same software requirements are prioritized for a second time using CV. A  
 162 developer might know that all vital functionality is selected by other stake-  
 163 holders, but his toy feature is left out. In effect, the developer could spend  
 164 all his points on this feature to put it in the next release.

165 Stakeholder bias may be mitigated by setting a maximum priority that  
 166 can be assigned to an item. This way each stakeholder is forced to distribute  
 167 the money between several prioritization items [4].

168 Another bias is that people in general tend to assign round priority values.  
 169 This is likely caused by lack of objective judgement criteria. Either way it  
 170 seems to be a problem not acknowledged by many since all prioritization is  
 171 largely based on expert opinion.

#### 172 2.3.4. Scalability of Cumulative Voting—Hierarchical Cumulative Voting

173 The standard CV approach has a low scalability. If the number of prior-  
 174 itization items is high, stakeholders may lose sight of the bigger picture and  
 175 assign priorities to a limited number of items. One, unsophisticated, solution  
 176 to the problem is to provide more points for prioritization (1,000 or 10,000  
 177 instead of 100); however, one could take another approach.

178 When the number of prioritization items is high they can usually be  
 179 grouped hierarchically by forming a tree structure (Figure 2) and, thus,  
 180 parent-child dependencies will exist between many items.

181 In [4] the authors propose a method for prioritizing hierarchically struc-  
 182 tured items called Hierarchical Cumulative Voting (HCV). It may be seen  
 183 as combination of the hierarchical part of the Analytical Hierarchy Process  
 184 (AHP) [1, 18] and the CV prioritization method. Since items are prioritized  
 185 in smaller sets, stakeholders do not lose sight of the bigger picture during  
 186 prioritization, and the prioritization of a large number of requirements is  
 187 considered easier.

### 2.3.5. Compensation Factors

HCV deals with the problem of prioritization scalability but it comes at a cost. Low level item groups may consist of different numbers of items, but the number of points spent on each group is the same, i.e. in a small-sized group, the same amount of points is distributed among fewer items. Hence, items in smaller groups are statistically more likely to have a higher priority, on average, compared to items in larger groups. To balance this difference each low level prioritization item can be multiplied by a compensation factor [4].

As an example, suppose an item ( $A$ ) in a group of 10 items is assigned 60 points. Hence,  $A$  will receive 600 compensated points. In this case it is impossible for any item in a group smaller than 6 items to compete with  $A$ . Even if item ( $B$ ) in a group of 5 is assigned the maximum number of points (100), the maximum compensated priority value  $B$  can receive is 500.

In [17] the authors suggest that compensated prioritization is more favorable compared to uncompensated. But neither compensated nor uncompensated prioritization is perfect and, as a general rule, it is better to keep the size of prioritization item groups similar.

### 2.3.6. HCV Execution

According to [4], HCV is conducted with the following steps (Steps 4–5 are optional):

1. Construct hierarchy. Prioritization items need to be divided into one high and several low level item groups. Each low level item group is child to exactly one high level item. And each high level item has one low level item group. One low level item may belong to several item groups. Even if parts of the items are not logically connected they can be grouped separately and assigned a fake parent item, e.g. ‘misc. items’. HCV does not, as far as we know, provide any instructions for creating a requirements hierarchy.
2. Each high and low level item group is prioritized separately using CV. The stakeholder may prioritize all item groups at once or one by one. But it should be possible to prioritize groups in any order and repeatedly, because the stakeholder might learn more about the items while performing the prioritization.



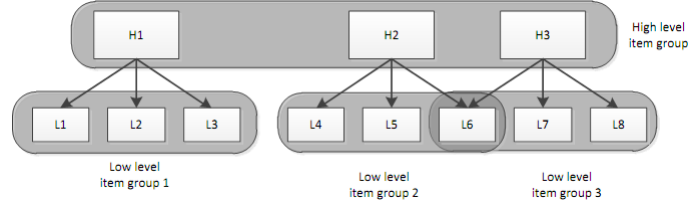


Figure 3: Overlapping prioritization item hierarchy example.

In particular the stakeholder is likely to learn more about a high level item when prioritizing its low level item group [19]. Some stakeholders may prioritize only part of the groups and each group may be prioritized by different stakeholders.

3. The priority of each low level item is normalized by dividing it with the sum of all low level priorities of each item in all groups.
4. The final priority of each low level item is calculated by multiplying it with the priority of its parent high level item.
5. Then one applies the compensation factor to all low level requirements as described in Section 2.3.5.
6. Finally, when multiple stakeholders have performed the prioritization, priorities of low level items are combined as in standard CV.

It is possible that one low level item is child of more than one high level requirement and, thus, belongs to two or more low level requirement groups (see Figure 3). Such requirements participate in the standard HCV prioritization process and are prioritized two or more times with each group they belong to. At the end of the prioritization they receive several priority values. These values must be summed together to form the final priority of the item. (This is done because the item adds value to both parts of the hierarchy.)

#### 2.3.7. Example of Hierarchical Cumulative Voting

Suppose six requirements for a mobile phone operating system need to be prioritized: ‘reminder alarm’, ‘specify repeated event’, ‘hide contact’, ‘add picture to phonebook’, ‘search contact’, ‘make video call’. Three high level requirements can be identified: ‘Calendar’, ‘Phonebook’, ‘Call’. The low level requirements are then grouped as sub-requirements of high level requirements

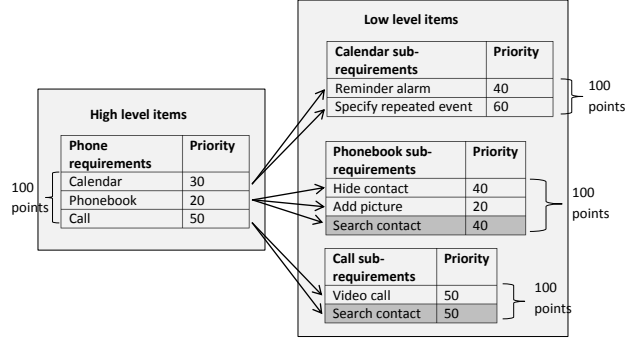


Figure 4: Example of hierarchical cumulative voting with requirement hierarchy.

Table 1: Example of hierarchical cumulative voting.

Phone requirements	Compensation factor	Sub-requirements	Priority calculation	Final priority
Calendar	2	Reminder alarm	$40 \times 30 \times 2$	2400
Calendar	2	Specify repeated event	$60 \times 30 \times 2$	3600
Phonebook	3	Hide contact	$40 \times 20 \times 3$	1600
Phonebook	3	Add picture	$20 \times 20 \times 3$	800
Phonebook & Call	3 & 2	Search contact	$40 \times 20 \times 3 + 50 \times 50 \times 2$	7400
Call	2	Video call	$50 \times 50 \times 2$	2500

as shown in Figure 4. The ‘Search contact’ requirement is a sub-requirement and has two parent requirements: ‘Phonebook’ and ‘Call’. The computation of the final priorities of requirements is shown in Table 1.

After requirements are grouped, and a hierarchy is defined, each group of requirements are then prioritized using CV. The final priority of a low level requirement is computed by multiplying the priority of the requirement with the priority of its parent high level requirement and the compensation factor. The compensation factor in this particular case is the number of elements in a group, two for the ‘calendar’ and ‘call’ sub-requirements and three for the ‘phonebook’ sub-requirement.

#### 2.4. Compositional Data Analysis

CV results can be seen as a special type of data, i.e. compositional data. Compositional data does not contain absolute values. It shows only the relative weight of a component compared to the whole. In [10] the authors

261 propose the use of compositional data analysis for the statistical analysis of  
 262 CV.

263 A compositional data item is a vector  $(x)$  of positive components with a  
 264 constant sum  $k$ :

$$x = (X_1; X_2; \dots; X_n) \text{ where } x_i \geq 0 \text{ and } \sum_{j=1}^n x_j = k \quad (1)$$

265 The property of the sum of the items being restricted is called the constant  
 266 sum constraint. In CV, priorities assigned by a stakeholder to the items of  
 267 a prioritization set is a compositional data vector with a constant sum of  
 268 100. The value of  $k$  (i.e. 100 in this case) is arbitrary and does not affect  
 269 the analysis of the data because the information is contained in the ratios  
 270 between the components of the vector. The vector can sum up to any number  
 271 but still hold the same data, i.e. vectors (1, 2, 7) and (10, 20, 70) are in this  
 272 case considered equivalent.

273 The priority of an item is relative to the priority of the other items in  
 274 the set. Hence, the priority of an individual item is meaningless without  
 275 context, i.e. the complete set of items. The same item may receive different  
 276 priority when put in two different prioritization sets. If the item is put in a  
 277 set of items with high priority it will receive a lower relative priority. This  
 278 also holds true the other way around i.e. if the item is put in a set with low  
 279 priority items its priority will be higher.

280 Compositional data analysis has, however, serious limitations. Ordinary  
 281 unconstrained variables are free to take any positive or negative values,  
 282 whereas, compositional data values can only be positive and have a con-  
 283 strained maximum value. Moreover, components of compositional data vec-  
 284 tors are not independent from each other. The fact that an item is assigned  
 285 70 priority points means that the next item can take only values between 0  
 286 and 30. Hence, there is a negative correlation between the items.

287 Standard parametric statistical tests require that data vectors have mul-  
 288 tivariate normal distribution. Vector  $X = (X_1, X_2, \dots, X_n)$  is considered to  
 289 have multivariate normal distribution if any linear combination of its parts  
 290 is normally distributed, and linear combination is defined by:

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n \quad (2)$$

291 where  $Y$  is the product of lineal combination and  $a_i$  is any real number.  
 292 Now, since the sum of priorities assigned in CV must add up to 100, or any

293 other constant number, at least one linear combination of  $X$  is not normally  
 294 distributed because it must always add up to 100:

$$Y = 1 \cdot X_1 + 1 \cdot X_2 + \dots + 1 \cdot X_n = 100 \quad (3)$$

295 In our opinion, the above indicates, quite strongly, that CV results do  
 296 not follow a multivariate normal distribution and, hence, it follows that they  
 297 should be analyzed using non-parametric statistical tests [20].

#### 298 2.4.1. Problem of Zeroes

299 Compositional data analysis requires that ratios between any components  
 300 in a vector can be computed. But computing a ratio with a zero value is,  
 301 in this case, meaningless. This is a problem since CV allows stakeholders to  
 302 assign zero priorities to some prioritization items (we would even strongly  
 303 argue that this is very common).

304 In compositional data there are two types of zeroes: essential and rounded.  
 305 Essential zeroes mean that a data component is not present. Rounded zeroes  
 306 mean that the component is present but its value is very low. We, as others  
 307 have before us, conjecture that zeroes in CV results are rounded because the  
 308 priority of an item is a completely abstract notion and the instrument for  
 309 measuring priority is human judgement [10].

310 Before compositional data analysis can be applied to CV results, we must  
 311 first remove zeroes in the data. One approach can be to forbid stakeholders to  
 312 assign zero priorities. This approach is used in e.g. [7]. But this can add some  
 313 unnecessary complexity to the prioritization process and, explicitly, delimits  
 314 an expert's freedom. In [10] the authors propose the use of a multiplicative  
 315 replacement strategy (as defined in [21]) for CV result analysis.

This method replaces rounded zeroes with small values using the expres-  
 sion

$$r_j = \begin{cases} \delta_j, & \text{if } x_j = 0, \\ (1 - \frac{\sum_{k|x_k=0} \delta_k}{c})x_j, & \text{if } x_j > 0, \end{cases} \quad (4)$$

316 where  $\delta_j$  is the imputed value and  $c$  is the constant sum constraint. In  
 317 order for the total sum of components to stay constant, the equation sub-  
 318 tracts some value from the items with a priority higher than zero. More is  
 319 subtracted from components with higher values than from components with  
 320 lower values (and the value of the imputed  $\delta_j$  is arbitrary).

### 321 2.4.2. Isometric log-ratio transformation

322 In order to apply standard statistical methods to compositional data it  
 323 must be transformed to remove the inherent correlation of the values. Com-  
 324 positional data analysis proposes special transformations that change the  
 325 compositional data values to unconstrained real values. One such transfor-  
 326 mation is isometric log-ratio (*ilr*) transformation (as proposed by [20, 22]):

$$\begin{aligned} z &= (z_1, \dots, z_{D-1}), \\ z_i &= \sqrt{\frac{i}{i+1}} \log \frac{\sqrt[i]{\prod_{j=1}^i x_j}}{x_{i+1}} \text{ for } i = 1, \dots, D-1 \end{aligned} \quad (5)$$

327 where  $x$  is the vector that is being transformed and  $z$  is the vector that  
 328 is created. It should be noted that  $z$  is shorter than  $x$  by one element.

329 After compositional data vectors are transformed using zero replacement  
 330 and *ilr*, any standard statistical tests can be applied.

## 331 3. Related Work

332 A systematic review of requirements prioritization methods is presented  
 333 in [23]. The study focuses on prioritization method comparison and selects  
 334 eight relevant studies. Two of the studies use CV. These studies are also  
 335 revealed by the systematic literature review conducted as part of this study.  
 336 In [23] the author concludes that there is little research on requirements  
 337 prioritization and studies usually deal with a small number of requirements.

338 The systematic literature review presented in this paper does not reveal  
 339 any CV result analysis methods that allows to identify prioritization items  
 340 with equal priority. Thus, this problem is not addressed in any way.

## 341 4. Methodology

342 This section covers the research questions of this study and the methods  
 343 used to answer them.

### 344 4.1. Selection of Research Methods

345 The main purpose of this study is to collect knowledge on the use of CV  
 346 in order to help software engineers and researchers in adopting it.

One way of collecting this knowledge is to conduct an empirical study. A survey in a large number of software companies can be used to quantify the level of adoption of CV in industry (similarly to the study by [24]), while a case study can be used to receive qualitative feedback on the use of CV [25].

Knowledge on the empirical use of CV can also be obtained from existing studies. This may be done by means of a systematic literature review. Several studies have used CV in industry as well as in academic settings. Nevertheless, there are no studies that provide an overview of the current state of the practice in this field (as reported by research studies). Therefore, before continuing with the refinement of CV and conducting new empirical studies (i.e. case study or experiment), a systematic literature review would be required.

This paper proposes a new method for CV result analysis, called Equality of Cumulative Votes (ECV). (ECV groups prioritization items into groups of items with similar priority.) As will be presented later, the systematic review did not reveal any methods that solve this problem; however, ECV needs to be evaluated and, hence, applied to CV results.

There are two options to obtain CV results in order to test ECV. One is to conduct a new empirical study. The second option is to collect CV results from existing studies. The latter approach also has the added benefit of trying to replicate the results from previous studies and, if data from several other studies are used, a larger amount of data can be obtained. Moreover, the generalizability of the evaluation increases when prioritization results from different sources and domains are used. On the other hand, the main benefit of conducting a separate empirical study is the possibility to control the conditions of CV.

In our study we evaluated ECV by obtaining data from previously conducted studies as found by the systematic literature review. In order to obtain the data, authors of relevant primary studies were contacted.

In short, this study consists of two parts: a systematic literature review (SLR) of CV and an evaluation of ECV based on the data from the primary studies found in the SLR.

#### *4.2. Research Questions*

The systematic review should focus on catching studies that empirically use CV. Information about place, time, scale, and domain of the studies should be collected and the results of the review will hopefully aid academic researchers by identifying paths for further investigation of CV. Hence, the first research question is:

384 **RQ 1.** What is the state of practice in empirical studies that use CV?

385 The level of trust in research results considering CV is determined by the  
386 quality of the studies that use CV, hence this study includes an evaluation  
387 of the quality of primary studies identified by the systematic review.

388 Next, a valuable aspect of decision-making is the analysis of prioritization  
389 results. Thus, the second research question is:

390 **RQ 2.** What CV result analysis methods have been presented in papers as  
391 identified by RQ 1?

392 Finally, the evaluation of ECV answers the third research question:

393 **RQ 3.** Is ECV capable of identifying prioritization items with equal priority?

## 394 **5. Systematic Literature Review**

395 This section presents the design of the systematic literature review. For  
396 the results of the execution please see Section 7.1 and 7.2.

397 Table 2 presents an overview of activities performed during the systematic  
398 literature review. The review protocol was developed by one researcher and  
399 evaluated by another researcher. Studies were searched for in two iterations.  
400 The first search was performed using databases. The second search was  
401 performed using snowball sampling [26] (snowball sampling examines the  
402 references of primary studies revealed by the first search). References that  
403 are relevant to the review, i.e. they pass the selection criteria, are then added  
404 to the set of primary studies.

405 The search for papers was performed by a single researcher. Study se-  
406 lection, on the other hand, was performed by two researchers. First, one  
407 researcher examined all found studies. Next, another researcher re-examined  
408 all studies classified as primary studies in addition to 20 randomly selected  
409 excluded studies to ensure the quality of the selection.

410 To ensure the quality of the review, the quality evaluation and data ex-  
411 traction was performed independently by two researchers. Inter-rater analy-  
412 sis was performed using Krippendorff’s Alpha statistics [27, 28].

Table 2: Review activities.

Review phase		Researchers involved
Trial search in databases		A
Develop review protocol		A
Evaluate review protocol		B
Paper search and selection from databases	Search in databases	A
	Search string validation	A
	Selection based on metadata	A and B
	Selection based on full text	A and B
Pilot data extraction (3 papers)		A
Paper selection from the reference lists	Selection based on metadata	A and B
	Selection based on full text	A and B
Data extraction		A and B
Data synthesis		A



413 *5.1. Data Sources and Search Strategy*

414 The SLR was designed based on the guidelines by Kitchenham [29]. First  
415 a trial search in electronic databases was conducted. In order to scale the  
416 review to a manageable, yet sufficient size, databases were searched with dif-  
417 ferent search strings. Relevant papers that were found during the trial search  
418 were used to extract additional search strings. The trial search revealed that  
419 the number of studies that use CV is not very large. Therefore, we decided  
420 to include not only software engineering studies but also studies in other re-  
421 search areas, such as forestry or corporate governance, since one key aspect  
422 we intended to investigate was analysis methods for CV.

423 Since CV is frequently used in studies without mentioning this in the  
424 abstract, full text search in databases is preferable. Unfortunately not all  
425 databases support full text search. Full text search was performed in the  
426 IEEE Xplore and Springer Link databases. In ACM Digital Library, In-  
427 spec/Compendex, ISI Web of Knowledge, and SCOPUS only metadata was  
428 searched. The search strings used, consisting of a Boolean expression (A or  
429 B or C or D or E or F or G), where:

- |                           |                               |
|---------------------------|-------------------------------|
| 430 (A) Cumulative voting | 434 (E) hundred dollar method |
| 431 (B) 100 dollar method | 435 (F) hundred dollar test   |
| 432 (C) 100 dollar test   |                               |
| 433 (D) 100 point method  | 436 (G) hundred point method  |

437 Search strings contained only synonyms of CV and they did not limit the  
438 research area to software engineering. The search was performed indepen-  
439 dently using each of the search strings in each database. All search results  
440 were combined and documented using reference management software. The  
441 quality of the search strings and the selection of electronic databases were  
442 validated against a previously known core set of papers—[3, 10, 30, 31]—  
443 checking that all papers from the core set were found by the search.

444 *5.2. Study Selection*

445 To select relevant papers a set of criteria were designed. The criteria for  
446 paper selection are presented in Tables 3 and 4.

447 Papers were selected in two phases: based on metadata and based on full  
448 text.

Table 3: Paper search and selection in the databases.

Selection phase	Inclusion criteria	Number of papers selected
Search in databases	published 2001–2011 (databases last accessed Feb. 20, 2011)	256
	contains search strings	
Selection based on metadata	exclude duplicates and tables of contents	177
	written in English	
Selection based on full text	full text is available	127
	study involves empirical use of CV or presents analysis of empirical use of CV	58
	CV is done by humans and not software	25

Table 4: Paper selection from the reference lists of the selected papers.

Selection phase	Inclusion criteria	Number of papers selected
Selection from references	papers included in the reference lists of relevant papers found in databases	467
Selection based on metadata	written in English	462
	reference is already revealed by search in databases	450
Selection based on full text	full text is available	329
	study involves empirical use of CV or presents analysis of empirical use of CV	15
	CV is done by humans and not software	

Obviously, the main criterion for inclusion of a paper is that it must present empirical use of CV or present an analysis of the results of using CV. However, there are papers that pass this criterion but are not relevant for this review. CV is frequently used in computer algorithms. There is a significant difference between the way humans and computers make decisions. Since this review is concerned with human decisions we excluded papers that present CV that is not performed by humans. In addition, only papers that were written in English were selected and duplicate studies were automatically excluded by the citation management software used in this review. We searched for papers between 2001–2011. By then performing a snowball sampling of these papers we are convinced that we have a representative sample and, furthermore, that the bulk of the studies are relevant from a software engineering perspective.

### 462 5.3. *Quality Evaluation*

463 The goal of quality evaluation is to determine the best primary studies  
464 according to some measure of quality. Since the number of studies that use  
465 CV is not large, quality evaluation was not used as an exclusion criterion.

466 The quality of a study obviously depends on the correctness of the study  
467 process including planning, operation, analysis and interpretation of the re-  
468 sults (is the study right?) The correctness of the process can be measured  
469 by evaluating the description of the study or replicating the study. Thus,  
470 to gain the trust of industry practitioners and other researchers, the process  
471 of the study must be rigorously described. In short, the description must  
472 facilitate replication of the study as well as the presentation of limitations  
473 and validity threats.

474 Even the most correct and rigorously described study is useless if it does  
475 not contribute to the industry or research community (is it the right study?)  
476 The topic of the research ought to address important goals and issues. The  
477 findings of the study should also be significant, i.e. there must be a high  
478 probability of the results of the study being true. The significance of the  
479 findings depends on how realistic the study is, the correctness of the process  
480 and the results of the study, as well as the statistical significance of the  
481 findings.

482 **Realism** of a study depends on the context, scale, and subjects of the  
483 study. The study should be conducted in a **setting** that is similar or equal  
484 to the setting in which the findings of the study are intended to be used.  
485 Hence, studies that are conducted in an industrial setting are in many cases  
486 valuable. The **subjects** of a study should be similar to the people who are  
487 supposed to use the findings of the study. The subjects ought to have appro-  
488 priate work experience, role in the organization, skills, cultural background,  
489 motivation, and so forth. The **scale** of a study refers to the size of the study  
490 objects. In the case of this systematic review the scale of a study is mea-  
491 sured as the number of prioritization items. Study in academia may have a  
492 large number of prioritization items. At the same time, an industrial study,  
493 with professionals as subjects, may involve a smaller number of prioritization  
494 items.

495 Each study may have a different level of realism. Some studies involve  
496 industry practitioners in an academic setting to simulate real word practice in  
497 a laboratory environment. Other studies may involve academic researchers  
498 that execute a project. For example, researchers may be developing open

499 source software. On the reality scale these studies are somewhere in between  
500 the purely academic and industrial studies.

501 The **type** of the research study can be considered as a criterion for the  
502 evaluation of study realism. [32] suggest that study designs that are more  
503 rigorous (e.g. experiments) are more realistic than observational studies (e.g.  
504 case study) due to a higher level of control. On the other hand [33] rate study  
505 designs based on other criteria, i.e. how frequently each type of study design  
506 is used in an industrial or academic setting. If a study design is used more  
507 in an industrial setting, then it is considered more realistic. For instance, in  
508 software engineering, case studies are frequently used in industrial settings,  
509 whereas, experiments are usually performed in academia using students as  
510 subjects. Therefore, the authors argue, case studies are more realistic than  
511 formal experiments [33]. Obviously the effect of study design on the study  
512 realism may be interpreted in different ways. Therefore, we will not use this  
513 parameter in our quality evaluation.

514 The statistical significance of the results of a study can be used to evaluate  
515 the significance of the study findings. This measure will not be used, because  
516 the studies that are evaluated belong to very different research areas, i.e. the  
517 significance levels of the findings of the studies are not directly comparable  
518 for meta-analysis. Additionally, sometimes no result is more interesting than  
519 a significant result. If a study's results do not conform to the expectations  
520 of researchers, this may reveal important gaps in existing knowledge.

521 The ultimate goal of research, at least in software engineering, is in many  
522 cases industry impact. However, most of the time ideas need to be devel-  
523 oped and validated in academia before industry professionals will risk to  
524 adopt them. Therefore, academic impact is important as well. Academic  
525 impact is usually measured by the number of citations. Academic impact is  
526 also measured for particular researchers, using the number of papers she has  
527 published and the number of times her papers have been cited. This measure  
528 will not be used in our quality evaluation because it is somewhat biased. The  
529 number of citations is likely to be lower for newer papers and the number  
530 of papers that a researcher has published gives little information about the  
531 actual quality or impact of her research.

### 532 5.3.1. *Rating of the Studies*

533 The quality evaluation in our review is based on the evaluation of: (i)  
534 Study realism. (ii) Study scale. (iii) Availability of raw results of CV. (iv)  
535 Quality of the research methodology.

536 Realism of the studies is rated in three aspects: subjects, setting, and  
537 scale. The subjects and setting is rated according to Table 5. The total  
538 rating of study realism is determined by summing up the ratings of the two  
539 aspects. For instance, if a study is conducted with industry professionals  
540 as subjects in an academic context the study will receive rating 1 (out of 2  
541 maximal points).

542 In order to rate the scale of a study the number of prioritization items was  
543 counted. If a paper presents several prioritization cases only the prioritization  
544 with the largest number of the prioritization items is considered. If HCV is  
545 used all of the prioritization items on different levels are counted together.  
546 However, if an item is present in several groups in the hierarchy it is counted  
547 only once.

548 The availability of raw results from the application of CV is rated sepa-  
549 rately because it is especially important for our purposes (and for most other  
550 researchers in order to replicate a study). The data availability rating criteria  
551 is given in Table 6. If the data of a study is not available it is not possible  
552 to validate the results of the study and, hence, the credibility of the findings  
553 is lower. Ideally the data collected in the study should be presented directly  
554 in the paper. An alternative may be to make the data freely available online  
555 and reference the online source.

556 The quality of the research methodology of a paper is rated according to  
557 a checklist presented in Appendix C. The checklist is based on guidelines  
558 for presenting research studies (as presented in [34, 35]) and the guidelines  
559 for quality evaluation of research studies as presented in [29, 33]. Evaluation  
560 is done with regard to the rigor of the description and correctness of the  
561 research process and reasoning. Checklist items represent issues that research  
562 studies should implement and present in a research paper. The checklist also  
563 contains item descriptions or questions that are used to evaluate the quality.  
564 Each item in the checklist is rated according to criteria presented in Table 7.  
565 The final rating of correctness of the research process of a study is computed  
566 by summing up the ratings assigned to all items in the checklist.

567 Study rating criteria was validated during a trial data extraction. Two  
568 researchers each rated three randomly selected papers. Afterwards, differ-  
569 ences in ratings were discussed and study rating criteria were updated to  
570 avoid differences in interpretation.

571 As a result of the rating each study was assigned four rating values on  
572 an ordinal scale. In order for us to perform a more advanced analysis of the  
573 quality evaluation results these ratings were then converted into ratio scale

Table 5: Rating of study reality level.

Aspect	Contribute to relevance (rating 1)	Do not contribute to relevance (rating 0)
Subjects	Industry professionals	Academia students or teachers, or other
Context	Industrial	Academia

Table 6: Research data availability rating.

Rating	Study rating criteria
0	CV results was not provided in the paper and we was unable to obtain the results from the authors.
1	CV results are not provided in the paper but the data was obtained from the authors. Part of the data is lost or corrupted.
2	CV results are not provided in the paper but all the data was obtained from the authors.
3	All CV results are included in the paper or reference is given to online source where all the data can be accessed.

Table 7: Rating of correctness of research process.

Rating	Study rating criteria
0	No description provided.
1	Only basic information is provided about the checklist item. Or significant validity threats exist with regard to this item.
2	Description is sufficient. Some minor questions are left unanswered. Validity threats may exist but they are not likely to affect the results of the study.
3	Description is rigorous and clear. Questions presented in quality evaluation checklist in Appendix C are answered. Decisions of the study are well justified, alternatives are discussed. No unhandled validity threats can be identified.

Table 8: Example of rating values.

Study	Realism	Research data availability	Correctness of research process	Number of prioritization items
ST1	2	0	15	6
ST2	1	3	20	69
ST3	0	3	10	6

Table 9: Example of ranking values.

Study	Reality level	Research data availability	Correctness of research process	Number of prioritization items
ST1	2	0	1	0
ST2	1	1	2	2
ST3	0	1	0	0

574 ranks. For each study, the number of studies that had received lower ratings  
575 were counted. The resulting number is the rank of the study; thereby, the  
576 quality of a study is expressed as four rank values.

577 An example of rating values is shown in Table 8. Table 9 shows ranking  
578 values computed for the studies in Table 8. We can observe that study  
579 realism level rating for ST3 is 0. There are no studies that have a lower  
580 study realism. Therefore, realism ranking for ST3 is 0. ST1 on the other  
581 hand has the highest realism rating. Since ST1 has higher reality level than  
582 both ST2 and ST3 it is assigned reality level rank 2.

#### 583 5.4. Data Extraction

584 The goal of data extraction is to understand how and why CV is used  
585 and how CV results are analyzed in research studies. Ultimately, this will  
586 allow us to answer the first and second research questions in our study.

587 Data extraction was documented with the help of spreadsheet software.  
588 Extracted data items are available from [36].

## 589 6. Equality of Cumulative Votes

590 In the previous section we described the execution of the systematic lit-  
591 erature review. In order to perform a more thorough analysis later we here  
592 present the design of ECV before presenting the results of the systematic  
593 literature review. For the results of the evaluation of ECV please see Sec-  
594 tion 7.3 (ECV is implemented in the *R* programming language [37] and the  
595 code can be found at [38].)

596 In CV stakeholders may assign similar or equal values to several prior-  
597 itization items. As a result the difference between the items is small. The  
598 variation in priorities is caused not only by the difference between prioritization  
599 items but also by human error and lack of information. For instance,  
600 people tend to simplify the task of prioritization by assigning rounded values  
601 to items or giving equal values to several items [39].

602 During prioritization it may be beneficial to know which items are equal.  
603 A common example is software release planning where requirements are distributed  
604 among several product releases. If two or more requirements are  
605 considered equal they can be freely interchanged between the releases, and  
606 other criteria, such as cost or effort, may be used as sole indicators for planning  
607 that particular release.

#### 608 *6.1. Testing Equality of Two Items*

609 There are two ways to determine which prioritization items have similar  
610 priority. One approach is to find items that are different and consider other  
611 items as equal. Another approach is to find items that are equal.

612 The first approach uses statistical tests to evaluate differences between  
613 e.g. two sample means, in order to determine that two items are different.  
614 Samples in this case consist of priorities assigned by all stakeholders to a  
615 particular prioritization item. The number of stakeholders that perform the  
616 prioritization is frequently small. Hence, the size of the sample is very often  
617 too small for statistical tests to detect a significant difference and the tests,  
618 thus, identify too many equal items to make any useful conclusions.

619 ECV, in contrast, uses the second approach. It finds items that are  
620 similar and the rest of the items are considered different. This method tests  
621 the probability of the difference between the means of two items being smaller  
622 than the given value. In short, ECV tests the probability of the means of two  
623 prioritization items differing by less than 25%. If the probability is higher  
624 than 70% the items are considered equal.

625 The input to ECV is an  $n \times p$  matrix  $A$  that contains the raw results of  
626 the prioritization. The columns of the matrix represent prioritization items  
627 while rows represent stakeholders. ECV performs the following operations  
628 for the priorities of each of the two prioritization items:

- 629 1. Replace zeroes in CV results.
- 630 2. Transform the data using *ilr* transformation.



- 631 3. Determine distribution function using kernel density estimation.
- 632 4. Use the distribution function to find the probability that the difference  
633 between two prioritization items is smaller than 25%.
- 634 5. Form groups of equal prioritization items.

635 Since CV results are compositional data, zeroes in  $A$  must be replaced  
636 with other values. This is done using the multiplicative replacement strategy  
637 which is described in Section 2.4.1. Next, two columns are extracted from  
638 matrix  $A$  to create the new matrix  $B$ :

$$B = [a_{*,k} a_{*,l}] \quad (6)$$

639 where  $a$  is an element of matrix  $A$ , and  $k$  and  $l$  are the columns that  
640 represent items that are tested for equality.

641 The  $ilr$  transformation is then applied to each row of the matrix  $B$  and  
642 the new vector  $C$  is obtained. The equation for calculating elements of  $C$   
643 using  $ilr$  transformation is:

$$c_i = ilr(b_{i1}, b_{i2}) = \sqrt{0.5} \log(b_{i1}/b_{i2}) \quad (7)$$

644 where  $c_i$  is the  $i^{th}$  element of  $C$  and  $b_{i1}$  and  $b_{i2}$  are the first and second  
645 elements in the  $i^{th}$  row of  $B$ . Each value  $c_i$  represents a ratio between  $k$  and  $l$ .  
646 The mean of the values of  $C$  can be interpreted as an average ratio between  
647 the items that expresses the difference between the items.

648 After the data is transformed into log-ratios statistical test can be applied.  
649 The purpose of the test is to determine what the probability is of the relative  
650 difference between two prioritization items  $k$  and  $l$  being less than 25%. This  
651 means determining the probability of the ratio  $k/l$  between the items  $k$  and  $l$   
652 as being in the range of  $\frac{3}{4}$  to  $\frac{4}{3}$ . Or in terms of log-ratios it means determining  
653 the probability of  $ilr(k, l)$  being between  $ilr(3, 4)$  and  $ilr(4, 3)$ . Hence, the  
654 objective of the test is to determine the probability of the sample mean (i.e.  
655 mean value of  $C$ ) laying between the two values.

656 The probability that the mean takes a particular value can be expressed  
657 in the form of a cumulative distribution function. The probability of the  
658 mean being between two values  $a$  and  $b$  (where  $a$  is smaller than  $b$ ) can be  
659 determined by subtracting the probability of the mean being smaller than  $a$   
660 from probability of the mean being smaller than  $b$ .

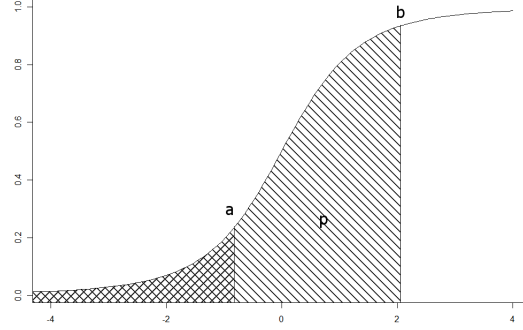


Figure 5: Cumulative distribution function of the ratio  $k/l$  between the items  $k$  and  $l$  (area  $p$  denotes probability that  $k/l$  is between  $\frac{3}{4}$  and  $\frac{4}{3}$ .)

661 However, CV result data may or may not be normally distributed. If  
 662 the data is normally distributed a Student's  $t$ -test can be used; otherwise, a  
 663 non-parametric estimation of the distribution function is needed.

664 In our case, the CV result data obtained from the primary studies identi-  
 665 fied by the systematic review, were tested for normality using the Anderson-  
 666 Darling test. The tests we performed indicated, quite strongly, that in most  
 667 of the prioritization cases the data is not normally distributed. Hence, our  
 668 recommendation is that, in general, a non-parametric approach should be  
 669 used to determine the probability density function, and one such, common,  
 670 approach would be to use the kernel density estimation. (In our implemen-  
 671 tation of ECV in the  $R$  programming language, kernel density estimation is  
 672 performed using the package *ks*.)

673 To determine the probability of  $\bar{x}$  being between  $a$  and  $b$  the following  
 674 equation is used:

$$p = P(b) - P(a) \quad (8)$$

675 where  $P$  is the cumulative distribution function obtained by applying  
 676 kernel density estimation on  $ilr$ -transformed priority values denoted by vector  
 677  $C$ . Variable  $a$  is equal to  $ilr(3, 4)$  and  $b$  is equal to  $ilr(4, 3)$ . (A graphical  
 678 interpretation of Equation (8) is presented in Figure 5.) The area that is  
 679 denoted by letter  $p$  represents the probability computed by the equation.

680 After both prioritization items are tested for equality it may be convenient  
 681 to display the equality of different items in the form of a table. Please see  
 682 Table 10 for an example.

Table 10: Example of an equality table.

prioritization items	i1	i2	i3	i4
i1	equal	equal	-	equal
i2	equal	equal	-	-
i3	-	-	equal	-
i4	equal	-	-	equal

## 6.2. Grouping Prioritization Items

When equal items are determined they must be divided into groups of equal items. Division must be performed in such a way that each two items in a group are equal. The test for equality of the items described in Section 6.1 is non-transitive. Hence, if prioritization item  $A$  is equal to  $B$  and  $B$  is equal to  $C$  then it does not automatically imply that  $A$  is equal to  $C$ . Therefore, there may be several ways to group the equal items. The two possible division criteria that we have considered in this study are:

1. Maximize the number of items that have a group.
2. Maximize the number of items in each group.

## 7. Results

This section presents the results of this study including the systematic literature review and the application of ECV on industry and academic data collected from the primary studies. Data extracted from primary studies and the results of the quality evaluation are available in [36].

### 7.1. State of Practice in Empirical Studies that use CV or Analyze the Results of CV (RQ 1)

The study search resulted in 634 unique studies. The search in databases revealed 180 papers, while an additional 454 papers were discovered using snowball sampling. The study selection resulted in 40 primary studies. Hence, 94% of the studies were excluded by the selection criteria. Snowball sampling revealed 15 (36%) out of all primary studies. The study selection criteria and the number of papers excluded by each criterion are shown in Tables 3 and 4. In total 163 of 634 studies were excluded because full text was not available.

All results of the study selection are available online and can be obtained by contacting the authors of this paper. For each study we specify keywords

710 and databases that were used to find the study. If a study has been excluded,  
711 the exclusion criteria are provided.

712 The number of papers revealed by each search string and database is  
713 presented in Table 11. It should be noted that several papers were found by  
714 more than one search string or in more than one database. Table 11 shows  
715 that the search string ‘cumulative voting’ was the most frequently used in  
716 the research community to denote CV. Therefore, researchers should use or  
717 reference this term when discussing CV.

718 To perform snowball sampling we examined the references of primary  
719 studies that were found during the database search. References were used  
720 to search for the papers in the Google and Google Scholar search engines.  
721 Studies that were found in the search and passed the study selection criteria  
722 were added to the set of primary studies.

723 After the primary studies were selected, data extraction and quality evalu-  
724 ation was performed by two researchers. One researcher examined all studies  
725 while the second researcher did quality evaluation and data extraction for  
726 10% of the studies. The studies were randomly selected. Inter-rater agree-  
727 ment were calculated by means of Krippendorff’s alpha coefficient. Agree-  
728 ment for data extraction results was 0.86 and agreement for the quality evalu-  
729 ation was 0.73. According to [28] it is common to require agreement above 0.8  
730 and the lowest acceptable agreement is 0.667. Therefore, we conclude that  
731 the agreement calculated for this study is sufficient. Ratings of the study  
732 setting, correctness, research data availability, and number of prioritization  
733 items are presented in Figure 6.

734 Table 12 shows the studies with the highest quality according to our cri-  
735 teria. These studies show a high level of rigor in a realistic setting. Moreover,  
736 authors of the studies manifest confidence by providing raw data for further  
737 use and evaluation.

738 Figure 7 shows a bubble chart of the distribution of studies over research  
739 areas and time. The figure shows that CV was, as far as we know, first ap-  
740 plied some time ago in research of government elections. Nowadays, though,  
741 CV has been adopted in a wide range of software engineering areas, most  
742 frequently in requirements engineering and software release planning. Eight  
743 studies use CV in academia while the remaining 32 studies report on using  
744 CV in industry.

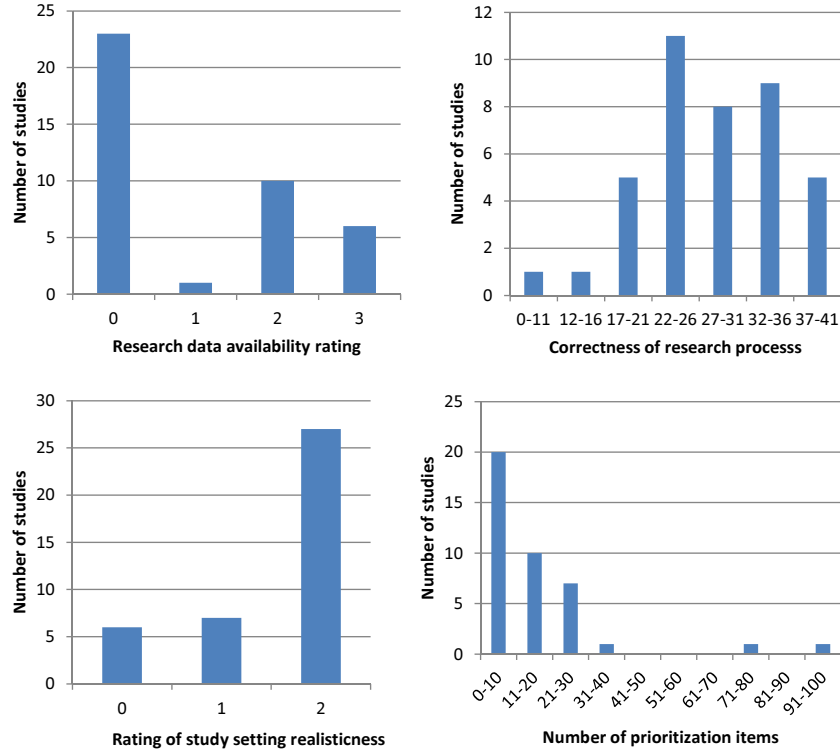
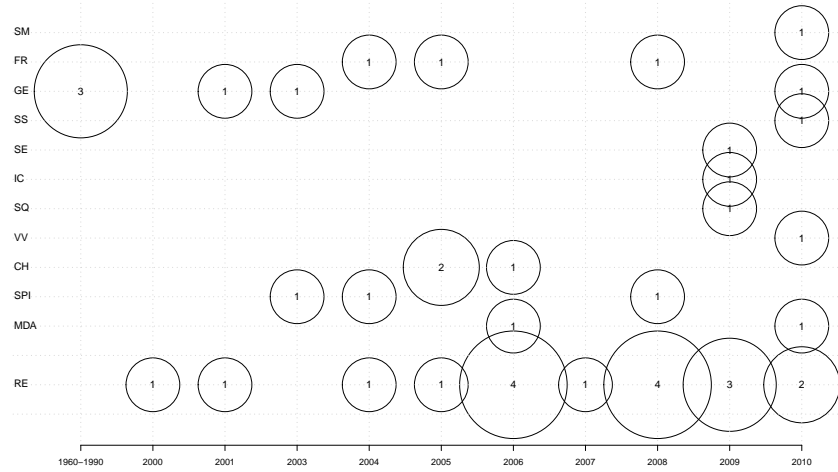


Figure 6: Study quality ratings.

Table 11: Number of papers found in the databases.

database	search strings							unique papers found	primary studies selected
	"100 point method"	"100 dollar method"	"100 dollar test"	"hundred point method"	"hundred dollar method"	"hundred dollar test"	"cumulative voting"		
ACM	2	0	0	1	2	3	31	34	7
IEEE	3	2	0	1	2	6	38	46	11
Inspec/Compendex	1	0	0	1	1	1	22	14	7
ISI web of science	0	0	0	0	1	1	15	16	6
SCOPUS	2	0	0	0	1	2	24	25	9
Springer	2	0	2	0	2	2	89	95	6
unique papers found	6	2	2	1	4	11	165	180	
primary studies selected	1	2	1	1	2	4	18		25



MDA - model driven software development  
 CH - change impact analysis in software engineering  
 RE - requirements engineering and software release planning  
 IC - intellectual capital in software company  
 SPI - software process improvement  
 V&V - software verification and validation  
 FR - forestry  
 GE - government elections  
 SS - software security  
 SQ - software quality  
 SM - software metrics  
 SE - software engineering in general

Figure 7: Distribution of studies over time.

Table 12: Top ranked studies.

	Correctness of research process	Research data availability	Study setting	Number of prioritization items
[40]	36	2	2	17
[17]	41	2	0	29
[41]	40	2	2	5
[8]	31	2	2	27
[42]	34	2	2	14
[43]	22	3	2	30
[44]	34	2	1	14
[45]	24	3	2	8
[31]	21	3	2	91
[46]	34	1	1	7

745 7.2. CV Result Analysis Methods Identified by RQ 1 (RQ 2)

746 The papers identified in the review use various CV result analysis meth-  
747 ods. The main goals for CV result analysis are presented in Table 13 and  
748 a summary of methods used in the primary studies can be found in Section  
749 Appendix B.

750 In order to present prioritization results many studies use charts or tables.  
751 These charts and tables show the average priority of each prioritization item  
752 that is computed from priorities assigned by all stakeholders. In [47] a table  
753 of five items with highest total priority is presented. [48] shows tables with  
754  $min$ ,  $max$ ,  $\tilde{x}$ ,  $\bar{x}$  and  $\sigma$  of priorities assigned by different stakeholders to a  
755 particular prioritization item. Finally, in [48, 49] error bars are added to the  
756 chart of final priorities (denoting  $\sigma$  of priorities).

757 In a few cases final priorities are presented in the form of ranks and  
758 CV results are degraded from ratio to ordinal scale. This is done when the  
759 interest lies only in the order of final priorities.

760 Several papers are interested in the difference between priorities from dif-  
761 ferent prioritization perspectives (e.g. current and ideal situation) or stake-  
762 holder groups (e.g. software developers and management). Pearson or Spear-  
763 man correlation coefficients are commonly used to determine what the level of  
764 similarity is between all priorities from two perspectives. Whereas, Wilcoxon,  
765 Kruskal-Wallis, Nemenyi-Damico-Wolfe-Dunn tests and the  $\chi^2$  statistic are  
766 used to detect if there is a significant difference in the value of one prioritiza-  
767 tion item from two or more perspectives. In addition, PCA is used to detect  
768 if there are distinct groups of stakeholders with common priorities [7, 10, 50].

769 In some cases, a stakeholder may assign equal priority to several prioritiza-  
770 tion items or leave several items unrated, e.g. the stakeholder may not have  
771 carefully considered all prioritization items. Hence, the difference between  
772 the items may have been unnoticed.

773 In [4] the scalability of prioritization is measured using two charts. The  
774 first chart shows the average percentages of items given a non-zero value.  
775 The second chart shows average percentages of divergence of values. If a  
776 stakeholder assigns equal priorities to many prioritization items the diver-  
777 gence of values is low. Unfortunately it is unclear from [4] how the average  
778 percentage of divergence is calculated.

779 In [51] distribution, disagreement, and satisfaction charts are presented.  
780 The distribution chart shows how the final value of a prioritization item  
781 is constructed from priorities assigned by different stakeholders. This chart

Table 13: Goals for CV result analysis.

Purpose of the method	Name
Show the final priority of each prioritization item. Stakeholder priorities are combined into one value.	Chart or table of final priorities
Difference between priorities assigned by different perspectives (status quo, ideal situation) or different stakeholder groups (developers, management) [10]	Bi-plot
detect stakeholder groups with similar priorities [10]	Bi-plot
show the relative number of issues that have received high, medium, or low priority [10]	Ternary plot
detect stakeholder groups with common priorities [10]	PCA
how the final value of prioritization item is constructed from priorities assigned by different stakeholder. This chart shows how much each stakeholder has contributed to the final value of prioritization item [51]	Distribution chart
the level of agreement between different stakeholders on value of particular prioritization item [51]	Disagreement chart
satisfaction of a stakeholder with the prioritization results by the calculating correlation between the final priorities and priorities assigned by a stakeholder [51]	Satisfaction chart
percentage of the divergence of the priorities assigned by a stakeholder [4]	average percentage of divergence
average percentage of items given a non-zero value [4]	
detect equal prioritization items (presented in this paper)	ECV

shows how much each stakeholder has contributed to the final value of a prioritization item. The disagreement chart shows the level of agreement between different stakeholders on the value of a particular prioritization item. The satisfaction chart shows stakeholder satisfaction with prioritization results by calculating the correlation between final priorities and priorities assigned by a stakeholder.

The use of bi-plots and ternary plots are proposed in [10]. A bi-plot shows final priorities and stakeholder viewpoints in a two dimensional plane while a ternary plot shows prioritization items inside a triangle. Ternary plots show how many low, medium or high priorities are assigned to a prioritization item. The corners of the triangle represent high, medium, and low priority, e.g. if a prioritization item has received mostly high priority values then it is shown closer to the high priority corner.

### 7.2.1. Problems with Compositional Data Analysis in Primary Studies

A few primary studies, as revealed by the systematic review, have problems with the analysis of compositional data.

In [7, 50] standard PCA is performed without applying log-ratio transformations to compositional data. According to [52], this is likely to be inadequate and in [53], a more appropriate method for performing PCA on



Table 14: Identified groups of equal items.

Paper identifier & Description	Type of CV	Pairs of equal items	Groups of equal items
[41] Perceived priorities of software product investments in an ideal situation	comp. HCV	(A2, B4) (B4, B5) (B4, C1) (B5, B15) (B6, B7) (B7, B8) (B14, B15) (B14, B18) (B17, B18)	(A2, B4) (B4, C1) (B5, B15) (B6, B7) (B14, B15) (B17, B18)
	uncomp. HCV	(B4, B5) (B4, B8) (B5, B15) (B6, B7) (B7, B12) (B14, B15) (B14, B18) (B16, B17) (B12, B13)	(B4, B5) (B5, B15) (B6, B7) (B14, B15) (B16, B17) (B12, B13)
[17] Software requirements for course management system	uncomp. & comp. HCV	(3:2, 3:3)	(3:2, 3:3)
[46] The view of academia researchers on the requirements understandability criteria	CV	(Development, Verification & Validation) (Development, Product Planning 1)	(Development, Product Planning 1)

801 compositional data is presented.

802 The normality of compositional data is defined in [54]. It is stated that  
803 compositional data must first be transformed using isometric log-ratio trans-  
804 formation before the tests for normality can be applied. In [47] the authors  
805 violate this requirement by applying the Shapiro-Wilk test for normality to  
806 untransformed compositional data.

807 The Kruskal-Wallis test is used in [47] to analyze compositional data.  
808 The test is used to evaluate the difference between three organization levels.  
809 The Kruskal-Wallis test assumes that variables within each sample are in-  
810 dependent [55]. However, values within compositional data vectors are not  
811 independent (as described in Section 2.4). Hence, we claim the Kruskal-  
812 Wallis test to be somewhat misused in [47].

### 813 7.3. Identifying Prioritization Items with Equal Priority Using ECV (RQ 3)

814 This section presents the results of applying ECV to the industrial and  
815 academic CV data as found through the systematic literature review. Six  
816 primary studies included the raw prioritization results in the paper itself or

817 referenced online sources where the data was available. To collect the data  
818 from the remaining 34 papers, the authors of all papers were contacted.

819 First, the email addresses provided in the papers were used. If no answer  
820 was received authors were searched for using Google, Facebook and LinkedIn.  
821 Authors from 11 papers provided us with data to be used in the evaluation  
822 of ECV. However, due to confidentiality reasons we can not publish this data  
823 directly.

824 In short, ECV was applied to 27 CV prioritization cases from 14 studies.  
825 In the cases of HCV, ECV was applied two times to the same data to test both  
826 compensated and uncompensated priorities. Equal items were detected in  
827 three prioritization cases. A summary of the results is presented in Table 14  
828 and below follows a summary of each relevant study.

829 In [46] a prioritization of requirement understandability criteria is pre-  
830 sented. One of the main findings of the paper is that from an academic  
831 viewpoint Development and Verification & Validation are more important  
832 than other criteria. ECV adds new knowledge to these results. It shows that  
833 Development and Verification & Validation are equally important, i.e. it is  
834 not true that either one of the criteria is more important.

835 A prioritization of software requirements for an academic course man-  
836 agement system is presented in [17]. ECV detected that two features—  
837 Assignment Submission and Assignment Feedback—have the same priority.  
838 If the system is developed in several releases Assignment Submission and As-  
839 signment Feedback features can be freely interchanged between the releases  
840 and, hence, in this way ECV simplifies release planning.

841 In [41] software product investments are prioritized with HCV. The re-  
842 sults of ECV was different for uncompensated and compensated HCV. When  
843 compensated HCV was used ECV detected equal items that belonged to dif-  
844 ferent high level prioritization groups (*A*, *B* and *C*) indicating that ECV  
845 provided a more fine-grained view. In the case of uncompensated HCV, on  
846 the other hand, all equal items belonged to one high level prioritization group  
847 (group *B*).

## 848 8. Discussion and Conclusions

849 This section discusses the results of the systematic review and evaluation  
850 of ECV conducted as part of this study.

851 CV has been applied in various areas, but most frequently in requirements  
852 prioritization and release planning, and quite often also as part of research

853 methodologies. A large part of the studies have been conducted in Sweden,  
854 at Ericsson AB. One can see a slight increase in the interest in CV. During  
855 the last five years there have been more studies that use CV than between,  
856 say, 2000–2005.

857 Overall, studies that use CV or analyze the results of CV have a high  
858 quality in terms of correctness of research process and study realism. How-  
859 ever, very few studies present prioritization of more than 30 items and the  
860 availability of research data is somewhat limited. In our particular case we  
861 were able to obtain data from 43% of the primary studies.

### 862 8.1. *Implications for Practitioners*

863 The results of this study provide decision support for industry practition-  
864 ers. We believe that a collection of state of the practice studies help the  
865 adoption of CV prioritization method. (The top studies are summarized in  
866 Table 12.) In addition, a set of CV analysis methods enables comprehen-  
867 sive understanding of the prioritization results. (The analysis methods are  
868 presented in Table 13.) One of the most common goals of CV analysis is to  
869 display the prioritization results and, thus, to show the difference between  
870 several prioritization perspectives.

871 Additionally, we present ECV—a novel method for CV analysis. Priori-  
872 tization often results in the assignment of similar priorities to several prior-  
873 itization items. CV results contain both ‘real priorities’ and random errors.  
874 Due to random errors, equal prioritization items may receive different pri-  
875 orities. ECV identifies such items. It allows stakeholders to disregard the  
876 random part of the CV results. Thus, ECV simplifies the understanding of  
877 the prioritization results.

878 ECV identifies prioritization items with similar priority and tests whether  
879 these items can be considered equal. In this case, ECV can be used in  
880 software release planning. For example, let us suppose that a set of software  
881 requirements are prioritized with regard to the implementation costs. First of  
882 all, ECV can then detect items with equal cost. Second, the equal items can  
883 be freely interchanged between the releases. Finally, the decision to allocate  
884 a requirement to a particular release can be made based on another criteria,  
885 such as risk or business value.

886 ECV has been successfully applied on a considerable amount of CV data  
887 and, additionally, has also detected equal items in different groups of HCV  
888 hierarchies.

889 8.2. *Implications for Academia*

890 In the systematic review 36% of papers were revealed by the snowball  
891 sampling. That is a considerable amount. Several studies do not mention  
892 the name of the prioritization method (i.e. cumulative voting or hundred  
893 dollar test). Others are not available through selected databases because  
894 they are conference publications or theses. It shows, in our opinion, that  
895 snowball sampling ought to be used in all systematic literature reviews.

896 CV results are a special type of data—compositional data. Standard sta-  
897 tistical analysis methods that assume the independence of the samples cannot  
898 be applied to CV results. In [56] methods for the analysis of compositional  
899 data have been presented. The systematic review conducted as a part of this  
900 study revealed that 22 studies analyze CV results; yet, only one study uses  
901 compositional data analysis methods, i.e. [10]. None of the studies, including  
902 [10], present methods for detecting items with equal priority in CV results.  
903 Hence, ECV is, in this respect, a unique method.

904 The small use of compositional data analysis is really not surprising, since  
905 literature describing CV does not state that the results are compositional  
906 data. Standard statistical analysis methods may produce useful results for  
907 compositional data. However, there are cases when they are misleading or  
908 even faulty. Section 7.2.1 contains evidence of inappropriate use of statistical  
909 methods by several papers.

910 This study has collected a set of compositional data analysis methods for  
911 CV analysis (see Table 13). We believe that this could help researchers to  
912 improve the analysis of CV results with appropriate methods.

913 Since CV is associated with compositional data, it might be tempting to  
914 choose another requirements prioritization method. However, it would not  
915 solve the problem *per se*, because any ratio scale prioritization, for instance  
916 AHP, contains compositional data.

917 The principal implications for the academia are mainly the following:

- 918 1. All systematic literature reviews should include snowball sampling.
- 919 2. Researchers can improve their statistical analysis of CV results using  
920 compositional data analysis methods collected and developed by this  
921 study.
- 922 3. When CV or any other ratio scale prioritization method is taught,  
923 compositional data analysis should also be presented as part of the  
924 solution.

### 8.3. *Validity Threats*

The validity of the systematic review is mainly limited by the chosen databases, the design of the review, and human judgement in study selection and data extraction.

To mitigate the threats we use the most popular databases in the field of software engineering. In the beginning of the systematic review a review protocol was developed, peer-reviewed, and revised. Search strategy was validated against a set of previously known papers obtained from other researchers.

One of many terms used to name cumulative voting is ‘\$100 method’. We were not able to search for this term because none of the chosen databases support search for special characters like ‘\$’ and the search string ‘100 method’ yields too many hits. To increase the likelihood of discovering relevant studies snowball sampling was extensively used.

To increase the validity of study selection, all included studies and 20 randomly selected excluded studies were examined by two researchers. There were no disagreement on the inclusion/exclusion of the studies.

The large number of studies identified by snowball sampling (15 out of 40 studies) may be caused by faulty design or by faulty execution of the search in the databases. There are several reasons why the studies revealed by snowball sampling are not revealed by the search in databases. (Reason for each study is given in Table Appendix A.2.) Based on these reasons we argue that snowball sampling does not indicate any problems with the design of the search in the databases.

Four studies were not found because they were not available through databases used in this systematic review. Out of them one is a master thesis, two are conference publications and one is a publication in the area of forestry. Seven studies do not mention the name of the prioritization method (i.e. hundred dollar method or cumulative voting). Only phrases like “distribution of a predefined amount of fictitious money (\$100,000) over the items to be prioritized” or “1,000 points” allowed us to identify that CV was indeed used. One paper used a previously unknown name for CV, i.e. the 100-point technique.

The quality of the data extraction and quality evaluation was validated using inter-rater agreement analysis. In our case, 10% of the studies were rated by two researchers and Krippendorff’s alpha was calculated. The agreement for the data extraction results was 0.86 and the agreement for the quality evaluation was 0.73 (indicating a credible level of quality).

963 There are two main validity threats with ECV itself. First, ECV may not  
964 detect prioritization items with equal priority. Second, ECV may produce a  
965 false positive result, i.e. there may be a real difference between items that  
966 ECV claims as being equal.

967 To mitigate the first threat ECV was applied on artificially created test  
968 data with and without items with similar priority. ECV worked correctly in  
969 both cases.

970 To mitigate the second threat we visually inspected the results of the  
971 application of ECV on the real world data from the primary studies. We  
972 concluded that items identified by ECV can be considered equal.

973 CV results used in the evaluation of ECV were tested for normality. The  
974 tests indicated that CV results are not normally distributed. Therefore, the  
975 design of ECV was based on a non-parametric statistical test.

#### 976 8.4. *Future Research*

977 There are very few studies that apply CV on prioritization sets of more  
978 than 30 items. However, in requirements engineering, industry practitioners  
979 need to prioritize much larger numbers of software requirements. Therefore,  
980 the state of art could benefit from the application of CV and HCV to large  
981 prioritization sets.

982 The proposed method, ECV, has now been evaluated on existing research  
983 data. To further evaluate the ECV, it could be applied in direct industry  
984 practice and in prioritization cases with a larger number of prioritization  
985 items. Additionally, compositional data analysis methods, as the ones iden-  
986 tified by this paper, should be tried with other prioritization methods that  
987 produce ratio scale results.

#### 988 8.5. *Conclusions*

989 CV prioritization results are special type of data – compositional data.  
990 Any analysis of CV results must take into account the compositional nature  
991 of the CV results.

992 This study presents a systematic literature review of the empirical use  
993 of CV. CV has been applied in various areas, but most frequently in re-  
994 quirements prioritization and release planning. The review has resulted in  
995 a collection of state of the practice studies and CV result analysis methods.  
996 We believe that it can help the adoption of CV prioritization method.

In our case, snowball sampling was performed as a part of the review. Since it revealed 36% out of all primary studies, we believe that in future snowball sampling should be used in all systematic reviews.

Additionally, we present ECV—a novel method for CV analysis. As suggested by our evaluation, ECV is able to detect prioritization items with equal priority (i.e. items that have insignificant difference in priority). The evaluation of ECV was based on the data obtained from the authors of the primary studies.

## References

- [1] P. Berander, A. Andrews, Requirements prioritization, in: A. Aurum, C. Wohlin (Eds.), *Engineering and managing software requirements*, Springer-Verlag, Berlin/Heidelberg, 2005, 2005, pp. 69–94.
- [2] D. Leffingwell, D. Widrig, *Managing software requirements: A unified approach*, Addison-Wesley Professional, 1999.
- [3] V. Ahl, An experimental comparison of five prioritization methods, Master’s Thesis, School of Engineering, Blekinge Institute of Technology, Sweden (2005).
- [4] P. Berander, P. Jönsson, Hierarchical cumulative voting (HCV) - Prioritization of requirements in hierarchies, *International Journal of Software Engineering and Knowledge Engineering* 16 (2006) 819–850.
- [5] J. Karlsson, K. Ryan, A cost-value approach for prioritizing requirements, *IEEE Software* 14 (1997) 67–74.
- [6] J. Karlsson, An evaluation of methods for prioritizing software requirements, *Information and Software Technology* 39 (1998) 939–947.
- [7] F. Pettersson, M. Ivarsson, T. Gorschek, P. Öhman, A practitioner’s guide to light weight software process assessment and improvement planning, *Journal of Systems and Software* 81 (2008) 972–995.
- [8] S. Barney, C. Wohlin, Software product quality: Ensuring a common goal, in: Q. Wang, V. Garousi, R. Madachy, D. Pfahl (Eds.), *Trustworthy Software Development Processes*, volume 5543 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2009, 2009, pp. 256–267.

- 1029 [9] P. Jönsson, C. Wohlin, A study on prioritisation of impact analysis  
1030 issues: A comparison between perspectives, *Software Engineering Re-*  
1031 *search and Practice in Sweden* (2005).
- 1032 [10] P. Chatzipetrou, L. Angelis, P. Rovegard, C. Wohlin, Prioritization of  
1033 issues and requirements by cumulative voting: A compositional data  
1034 analysis framework, in: *Proceedings of the 2010 36th EUROMICRO*  
1035 *Conference on Software Engineering and Advanced Applications*, IEEE  
1036 Computer Society, Washington, DC, USA, 2010, pp. 361–370.
- 1037 [11] R. L. Engstrom, D. A. Taebel, R. L. Cole, Cumulative voting as a rem-  
1038 edy for minority vote dilution: The case of Alamogordo, New Mexico,  
1039 *Journal of Law & Politics* 5 (1988) 469.
- 1040 [12] S. Bhagat, J. A. Brickley, Cumulative voting: The value of minority  
1041 shareholder voting rights, *Journal of Law and Economics* 27 (1984)  
1042 339–365.
- 1043 [13] V. Hiltunen, J. Kangas, J. Pykalainen, Voting methods in strategic  
1044 forest planning - Experiences from Metsähallitus, *Forest Policy and*  
1045 *Economics* 10 (2008) 117–127.
- 1046 [14] P. Boldi, F. Bonchi, C. Castillo, S. Vigna, Voting in social networks, in:  
1047 *Proceedings of the 18th ACM conference on Information and knowledge*  
1048 *management*, ACM, New York, NY, USA, 2009, pp. 777–786.
- 1049 [15] H. Ayad, M. Kamel, Cumulative voting consensus method for partitions  
1050 with variable number of clusters, *IEEE Transactions on Pattern Analysis*  
1051 *and Machine Intelligence* 30 (2008) 160–173.
- 1052 [16] M. Svahnberg, A. Karasira, A study on the importance of order in  
1053 requirements prioritisation, in: *Proceedings of the Third International*  
1054 *Workshop on Software Product Management*, IEEE Computer Society,  
1055 Washington, DC, USA, 2009, pp. 35–41.
- 1056 [17] P. Berander, M. Svahnberg, Evaluating two ways of calculating priorities  
1057 in requirements hierarchies - An experiment on hierarchical cumulative  
1058 voting, *Journal of Systems and Software* 82 (2009) 836–850.
- 1059 [18] T. L. Saaty, *The analytic hierarchy process*, McGraw-Hill, New York,  
1060 1980.



- 1061 [19] S. Brenner, J. Schwalbach, Legal institutions, board diligence, and  
1062 top executive pay, *Corporate Governance: An International Review*  
1063 17 (2009) 1–12.
- 1064 [20] V. Pawlowsky-Glahn, J. J. Egozcue, Compositional data and their anal-  
1065 ysis: An introduction, Geological Society, London, Special Publications  
1066 264 (2006) 1–10.
- 1067 [21] J. Martin-Fernandez, C. Barceló-Vidal, V. Pawlowsky-Glahn, Dealing  
1068 with zeros and missing values in compositional data sets using nonpara-  
1069 metric imputation, *Mathematical Geology* 35 (2003) 253–278.
- 1070 [22] P. Filzmoser, K. Hron, Outlier detection for compositional data using  
1071 robust methods, *Mathematical Geosciences* 40 (2008) 233–248.
- 1072 [23] K. Khan, A systematic review of software requirements prioritization,  
1073 Master’s thesis, Blekinge Institute of Technology, Ronneby, Sweden  
1074 (2006).
- 1075 [24] F. Zahedi, The analytic hierarchy process: A survey of the method and  
1076 its applications, *Interfaces* 16 (1986) 96–108.
- 1077 [25] P. Runeson, M. Höst, Guidelines for conducting and reporting case  
1078 study research in software engineering, *Empirical Software Engineering*  
1079 14 (2008) 131–164.
- 1080 [26] L. Goodman, Snowball sampling, *The Annals of Mathematical Statistics*  
1081 32 (1961) 148–170.
- 1082 [27] K. Krippendorff, Bivariate agreement coefficients for reliability of data,  
1083 *Sociological Methodology* 2 (1970) 139–150.
- 1084 [28] K. Krippendorff, Content analysis: An introduction to its methodology,  
1085 Sage Publications, 2nd edition, 2003.
- 1086 [29] B. Kitchenham, S. Charters, Guidelines for performing Systematic Lit-  
1087 erature Reviews in Software Engineering, Technical Report EBSE 2007-  
1088 001, Keele University, 2007.
- 1089 [30] P. Berander, P. Jönsson, A goal question metric based approach for  
1090 efficient measurement framework definition, in: *Proceedings of the 2006*

- 1091 ACM/IEEE international symposium on Empirical software engineer-  
1092 ing, ACM, New York, NY, USA, 2006, pp. 316–325.
- 1093 [31] B. Regnell, M. Höst, J. Natt och Dag, P. Beremark, T. Hjelm, An indus-  
1094 trial case study on distributed prioritisation in market-driven require-  
1095 ments engineering for packaged software, *Requirements Engineering* 6  
1096 (2001) 51–62.
- 1097 [32] B. Kitchenham, Procedures for performing systematic reviews, Techni-  
1098 cal Report TR/SE-0401, Keele University, 2004.
- 1099 [33] M. Ivarsson, T. Gorschek, A method for evaluating rigor and industrial  
1100 relevance of technology evaluations, *Empirical Software Engineering* 16  
1101 (2011) 365–395.
- 1102 [34] C. Wohlin, P. Runeson, M. Höst, Experimentation in software engineer-  
1103 ing: An introduction, Springer Netherlands, 2000.
- 1104 [35] A. Jedlitschka, D. Pfahl, Reporting guidelines for controlled experiments  
1105 in software engineering, in: *Proceedings of the 2005 International Sym-*  
1106 *posium on Empirical Software Engineering*, IEEE Computer Society,  
1107 2005, pp. 10–.
- 1108 [36] K. Rinkevics, R. Torkar, Data extraction and quality evaluation results,  
1109 2011.
- 1110 [37] R. Ihaka, R. Gentleman, R: A language for data analysis and graphics,  
1111 *Journal of computational and graphical statistics* 5 (1996) 299–314.
- 1112 [38] K. Rinkevics, R. Torkar, ECV implementation source code in R, 2011.
- 1113 [39] R. M. Groves, F. J. Fowler, M. P. Couper, J. M. Lepkowski, E. Singer,  
1114 *Survey methodology*, John Wiley and Sons, 2009.
- 1115 [40] S. Barney, A. Aurum, C. Wohlin, The relative importance of aspects  
1116 of intellectual capital for software companies, in: *Proceedings of the*  
1117 *2009 35th Euromicro Conference on Software Engineering and Advanced*  
1118 *Applications*, IEEE Computer Society, 2009, 2009, pp. 313–320.

- 1119 [41] S. Barney, C. Wohlin, Software product quality: Ensuring a common  
1120 goal, in: Proceedings of the International Conference on Software Pro-  
1121 cess: Trustworthy Software Development Processes, Springer-Verlag,  
1122 Berlin, Heidelberg, 2009, pp. 256–267.
- 1123 [42] S. Barney, A. Aurum, C. Wohlin, A product management challenge:  
1124 Creating software product value through requirements selection, *Journal*  
1125 *of Systems Architecture* 54 (2008) 576–593.
- 1126 [43] S. Laukkanen, T. Palander, J. Kangas, A. Kangas, Evaluation of the  
1127 multicriteria approval method for timber-harvesting group decision sup-  
1128 port, *Silva Fennica* 39 (2005) 249–264.
- 1129 [44] G. Hu, A. Aurum, C. Wohlin, Adding value to software requirements:  
1130 An empirical study in the Chinese software industry, in: Proceedings of  
1131 the Seventeenth Australasian Conference on Information Systems, 2006.
- 1132 [45] R. Feldt, R. Torkar, E. Ahmad, B. Raza, Challenges with software ver-  
1133 ification and validation activities in the space industry, in: Proceedings  
1134 of the 2010 Third International Conference on Software Testing, Verifi-  
1135 cation and Validation, IEEE Computer Society, Washington, DC, USA,  
1136 2010, pp. 225–234.
- 1137 [46] M. Svahnberg, T. Gorschek, M. Eriksson, A. Borg, K. Sandahl,  
1138 J. Börster, A. Loconsole, Perspectives on requirements understandabil-  
1139 ity – For whom does the teacher’s bell toll?, in: Proceedings of the 2008  
1140 Requirements Engineering Education and Training, IEEE Computer So-  
1141 ciety, Washington, DC, USA, 2008, pp. 22–29.
- 1142 [47] P. Jönsson, C. Wohlin, Understanding impact analysis: An empiri-  
1143 cal study to capture knowledge on different organisational levels, in:  
1144 Proceedings of International Conference on Software Engineering and  
1145 Knowledge Engineering, IEEE Computer Society, 2005, pp. 707–712.
- 1146 [48] L. Kuzniarz, L. Angelis, Empirical extension of a classification frame-  
1147 work for addressing consistency in model based development, *Informa-*  
1148 *tion and Software Technology* 53 (2011) 214–229.
- 1149 [49] P. Rovegard, L. Angelis, C. Wohlin, An empirical study on views of  
1150 importance of change impact analysis issues, *IEEE Transactions on*  
1151 *Software Engineering* 34 (2008) 516–530.

- 1152 [50] C. Wohlin, A. Aurum, Criteria for selecting software requirements to  
1153 create product value: An industrial empirical study, in: S. Biffl, A. Au-  
1154 rum, B. Boehm, H. Erdogan, P. Grünbacher (Eds.), Value-based soft-  
1155 ware engineering, Springer Verlag, 2006, 2006, pp. 179–200.
- 1156 [51] B. Regnell, M. Höst, J. Natt och Dag, Visualization of agreement and  
1157 satisfaction in distributed prioritization of market requirements, in: Pro-  
1158 ceedings of REFSQ2000, 6th Int. Workshop on Requirements Engineering:  
1159 ing: Foundation for Software Quality, 2000, pp. 1–12.
- 1160 [52] J. Aitchison, Principal component analysis of compositional data,  
1161 Biometrika 70 (1983) 57.
- 1162 [53] P. Filzmoser, K. Hron, C. Reimann, Principal component analysis for  
1163 compositional data with outliers, Environmetrics 20 (2009) 621–632.
- 1164 [54] V. Pawlowsky Glahn, J. Egozcue, R. Tolosana Delgado, Lecture notes on  
1165 compositional data analysis, Technical Report, Universitat de Girona,  
1166 Spain, 2007.
- 1167 [55] W. H. Kruskal, W. A. Wallis, Use of ranks in one-criterion variance  
1168 analysis, Journal of the American Statistical Association 47 (1952) 583–  
1169 621.
- 1170 [56] J. Aitchison, The statistical analysis of compositional data, Chapman  
1171 & Hall, London, 1986.
- 1172 [57] D. Baca, K. Petersen, Prioritizing countermeasures through the coun-  
1173 termeasure method for software security (CM-Sec), in: M. Ali Babar,  
1174 M. Vierimaa, M. Oivo (Eds.), Product-Focused Software Process Im-  
1175 provement, volume 6156 of *Lecture Notes in Computer Science*, Springer  
1176 Berlin/Heidelberg, 2010, 2010, pp. 176–190.
- 1177 [58] S. Bowler, D. Brockington, T. Donovan, Election systems and voter  
1178 turnout: Experiments in the United States, The Journal of Politics 63  
1179 (2001) 902–915.
- 1180 [59] D. Brockington, A low information theory of ballot position effect, Po-  
1181 litical Behavior 25 (2003) 1–27.

- 1182 [60] D. Cooper, A. Zillante, A comparison of cumulative voting and gener-  
1183 alized plurality voting, *Public Choice* (2010).
- 1184 [61] N. Dzamashvili Fogelström, M. Svahnberg, T. Gorschek, Investigating  
1185 impact of business risk on requirements selection decisions, in: *Proceed-*  
1186 *ings of the 2009 35th Euromicro Conference on Software Engineering and*  
1187 *Advanced Applications*, IEEE, 2009, pp. 217–223.
- 1188 [62] S. Hatton, Choosing the right prioritisation method, in: *Proceedings of*  
1189 *the 19th Australian Conference on Software Engineering*, IEEE Com-  
1190 puter Society, Washington, 2008, pp. 517–526.
- 1191 [63] S. Hatton, Early prioritisation of goals, in: *Proceedings of the 2007*  
1192 *Conference on Advances in Conceptual Modeling: Foundations and Ap-*  
1193 *plications*, Springer-Verlag, Berlin, 2007, pp. 235–244.
- 1194 [64] V. Heikkilä, A. Jadallah, K. Rautiainen, G. Ruhe, Rigorous support for  
1195 flexible planning of product releases - A stakeholder-centric approach  
1196 and its initial evaluation, in: *2010 43rd Hawaii International Conference*  
1197 *on System Sciences*, IEEE Computer Society, 2010, pp. 1–10.
- 1198 [65] M. Staron, C. Wohlin, An industrial case study on the choice between  
1199 language customization mechanisms, in: J. Münch, M. Vierimaa (Eds.),  
1200 *Product-Focused Software Process Improvement*, volume 4034 of *Lecture*  
1201 *Notes in Computer Science*, Springer Berlin / Heidelberg, 2006, 2006,  
1202 pp. 177–191.
- 1203 [66] T. Touseef, C. Gencel, A structured goal based measurement framework  
1204 enabling traceability and prioritization, in: *Proceedings of the 2010 6th*  
1205 *International Conference on Emerging Technologies*, 2010, pp. 282–286.
- 1206 [67] P. Berander, C. Wohlin, Differences in views between development roles  
1207 in software process improvement - A quantitative comparison, in: *Pro-*  
1208 *ceedings of the 8th International Conference on Empirical Assessment in*  
1209 *Software Engineering*, 2004.
- 1210 [68] P. Berander, Using students as subjects in requirements prioritization,  
1211 in: *Proceedings of the 2004 International Symposium on Empirical Soft-*  
1212 *ware Engineering*, IEEE Computer Society, 2004, pp. 167–176.

- [69] P. Berander, C. Wohlin, Identification of key factors in software process management - A case study, in: Proceedings of the 2003 International Symposium on Empirical Software Engineering, IEEE Computer Society, Washington, DC, USA, 2003, pp. 316–.
- [70] R. L. Cole, D. A. Taebel, R. L. Engstrom, Cumulative voting in a municipal election: A note on voter reactions and electoral consequences, The Western Political Quarterly 43 (1990) 191.
- [71] J. Kuklinski, Cumulative and plurality voting: An analysis of Illinois' unique electoral system, The Western Political Quarterly 26 (1973) 726–746.
- [72] S. Laukkanen, T. Palander, J. Kangas, Applying voting theory in participatory decision support for sustainable timber harvesting, Canadian Journal of Forest Research 34 (2004) 1511–1524.
- [73] J. Sawyer, D. MacRae, Game theory and cumulative voting in Illinois: 1902-1954, The American Political Science Review 56 (1962) 936–946.

## Appendix A. Primary Studies

### Appendix A.1. Primary studies found in databases.

Title	Reference
Prioritizing countermeasures through the countermeasure method for software security (CM-Sec)	[57]
The relative importance of aspects of intellectual capital for software companies	[40]
Software product quality: Ensuring a common goal	[8]
Balancing software product investments	[41]
Hierarchical cumulative voting (HCV) prioritization of requirements in hierarchies	[4]
A goal question metric based approach for efficient measurement framework definition	[30]
Evaluating two ways of calculating priorities in requirements hierarchies: An experiment on hierarchical cumulative voting	[17]
Election systems and voter turnout: Experiments in the United States	[58]
A low information theory of ballot position effect	[59]
Prioritization of issues and requirements by cumulative Voting: A compositional data analysis framework	[10]
A comparison of cumulative voting and generalized plurality voting	[60]
Challenges with software verification and validation activities in the space industry	[45]
Investigating impact of business risk on requirements selection decisions	[61]
Choosing the right prioritization method	[62]
Early prioritization of goals	[63]
Rigorous support for flexible planning of product releases: A stakeholder-centric approach and its initial evaluation	[64]
Voting methods in strategic forest planning: Experiences from Metsähallitus	[13]
Empirical extension of a classification framework for addressing consistency in model based development	[48]
Evaluation of the multi-criteria approval method for timber-harvesting group decision support	[43]
A practitioner's guide to light weight software process assessment and improvement planning	[7]
An empirical study on views of importance of change impact analysis issues	[49]
An industrial case study on the choice between language customization mechanisms	[65]
Perspectives on requirements understandability—For whom does the teacher's bell toll?	[46]
A study on the importance of order in requirements prioritization	[16]
A structured goal based measurement framework enabling traceability and prioritization	[66]

1231 *Appendix A.2. Primary studies revealed by snowball sampling.*

Reference	Title	Reason why the paper is not revealed by the search in databases
[3]	An experimental comparison of five prioritization methods	Selected databases does not contain the paper, master thesis at BTH
[42]	A product management challenge: Creating software product value through requirements selection	Prioritization method name not mentioned, phrase "1,000 points" used instead.
[67]	Differences in views between development roles in software process improvement—A quantitative comparison	Prioritization method name not mentioned, phrase "100 points" used instead.
[68]	Using students as subjects in requirements prioritization	Unknown CV name: 100-point technique
[69]	Identification of key factors in software process management: A case study	Prioritization method name not mentioned, phrase "100 points" used instead.
[70]	Cumulative voting in a municipal election: A note on voter reactions and electoral consequences	Study published before year 2001.
[44]	Adding value to software requirements: An empirical study in the chinese software industry	Prioritization method name not mentioned, phrase "1,000 points" used instead.
[9]	A study on prioritization of impact analysis issues: A comparison between perspectives	Selected databases does not contain the paper.
[47]	Understanding impact analysis: An empirical study to capture knowledge on different organizational levels	Selected databases does not contain the paper.
[71]	Cumulative and plurality voting: An analysis of Illinois' unique electoral system	Study published before year 2001.
[72]	Applying voting theory in participatory decision support for sustainable timber harvesting	Selected databases does not contain the paper.
[31]	An industrial case study on distributed prioritization in market-driven requirements engineering for packaged software	Prioritization method name not mentioned: "distribution of a predefined amount of fictitious money (\$100,000) over the items to be prioritized."
[51]	Visualization of agreement and satisfaction in distributed prioritization of market requirements	Prioritization method name not mentioned: "distribution of a predefined amount of fictitious money (\$100,000) over the items to be prioritized."
[73]	Game theory and cumulative voting in Illinois: 1902–1954	Study published before year 2001.
[50]	Criteria for selecting software requirements to create product value: An industrial empirical study	Prioritization method name not mentioned: "The subjects had 1,000 points to spend among the 13 criteria."

1233

## Appendix B. CV Result Analysis Methods

	Paper																					
	Svalberg2008	Svalberg2009	Sturon2006	Pettersson2008	Wohlitz2006	Laakkonen2005a	Hu2006	Jonsson2005a	Kuzniarz2010	Rovgaard2008	Berander2006a	Berander2004a	Berander2006	Feldt2010	Barney2009b	Barney2008	Barney2009a	Barney2009	Jonsson2005	Chatzidimitrou2010	Regnell2001	Regnell2000
Analysis method	x			x												x						
Table that shows final priorities	x			x												x						
Chart that shows final priorities	x			x	x	x	x									x						
Table of top-5 prioritization items								x														
min, max, $\bar{x}$ , $\bar{\sigma}$ and $\sigma$ of priorities assigned by different stakeholders										x	x											
Bar chart of prioritization results showing $\bar{x}$ priority and $\sigma$ of priorities										x	x											
Pearson correlation coefficient		x								x	x		x									
Nemenyi Damico Wolfe Dunn														x								
Spearman's $r$															x		x					
Kruskal-Wallis								x														
Wilcoxon							x															
Correlation matrix		x													x		x					
Chart for comparing priorities from two perspectives, priorities are points in two dimensional plane, x- and y-axis represent two different perspectives											x								x			
Difference between priorities assigned by each two stakeholders using $\chi^2$ -statistic										x												
Median ranks		x																				
CV results converted to priority ranks		x											x					x				
PCA				x	x															x		
Percentage of divergence of priorities assigned by a stakeholder											x											
Average percentage of items given non-zero value											x											
Distribution chart																					x	x
Disagreement chart				x																	x	x
Satisfaction chart				x																	x	x
Bi-plot																					x	
Ternary plot																					x	

1234

## Appendix C. Quality Evaluation Checklist

	Item	Question or Description of the Item	Rating
1.	Background, introduction	Introduce research area	
2.	Problem statement, purpose	What is the problem [35]? Where does it occur [35]? Who has observed it [35]? Why is it important to be solved [35]?	
3.	Context, independent variables (aka. environment, setting)	Study location, time constraints, application domain, organization, tools, market, process (e.g. software development methodology), size of project, product that is being developed	
4.	Related work	Other existing work, alternative technologies, solutions, and studies	
5.	Goals and Hypotheses	Null hypothesis and one or more alternative hypotheses for each goal	
6.	Research questions		
7.	Design, Research methods		
7.1.	Design	Description of each step of the study	
7.2.	Control group	If there is a control group, are participants similar to the treatment group participants in terms of variables that may affect study outcomes [29]?	
7.3.	Randomization	Random selection of participants and objects Random assignment of treatment and objects to participants Random order of treatments in case of paired design. If each participant is assigned two treatments A and B, then part of participants perform A first and the other part start with B	
7.4.	Blocking	Group participants of the study into homogeneous groups called blocks (e.g. students in one course, database developers in one company) and implement the study design within each block independently. The idea is that variability of independent variables (e.g. experience and knowledge of subjects) is smaller within a group. That helps measuring changes in dependent variables [32].	
7.5.	Balancing	Equal number of subjects should be assigned to each treatment [32]	
7.6.	Blinding	Automated assignment of treatments to subjects [32] Automated distribution of study materials to subjects [32] Persons who grade the task results should not know which treatment was used [32] Analyst should not know which treatment group is which [32] Automated data collection from subjects [32]	
8.	Subjects (participants)		
8.1.	Population		
8.2.	Sampling	How sampling is performed? What subjects are included and excluded? [29] What is the type of the sampling (e.g. convenience, random)? Is the sample(selected participants) representative of the population?	
8.3.	"Drop outs" and response rate	Are reasons given for refusal to participate [29]?	
8.4.	Subject motivation	E.g. material benefits, course credits for students, etc.	
9.	Objects	E.g. documents and other artifacts	
10.	Measures, Data collection procedures	Who, when, and how to measure [29]? How is the measurement supported? Is it automated [29]? Are the measures used in the study the most relevant ones for answering the research questions [29]?	
11.	Analysis procedure		
11.1.	Data description	Do the numbers add up across different tables and subgroups [29]?	
11.2.	Data types (continuous, ordinal, categorical)		
11.3.	Scoring systems		
11.4.	Data set reduction, outliers		
11.5.	Statistical methods	Are the assumptions of statistical methods met? What statistical programs are used?	
11.6.	Statistical significance	If statistical tests are used to determine differences, is the actual $p$ -value given [29]? If the study is concerned with differences among groups, are confidence limits given describing the magnitude of any observed differences [29]?	
12.	Validity threats	Threats, implications of the threats, and threat mitigation	
12.1.	Side-effects during study execution	Deviations from the plan, solutions for the deviations	
13.	Most important findings	Are all study questions answered [29]? Are negative findings presented [29]?	
14.	Industry impact, inference, generalization	What implications does the report have for practice [29]? How and where the results can be used? Limitations under which findings are relevant [35]?	
15.	Future work		