# On Eliciting Preference and Influence Measures in Goal Models

### Sotirios Liaskos
School of Information Technology
York University, Toronto, Canada
liaskos@yorku.ca

### Rina Jalman
School of Information Technology
York University, Toronto, Canada
rjalman@yorku.ca

### Jorge Aranda
Dept. of Computer Science
University of Victoria, Victoria, Canada
jaranda@uvic.ca

March 6, 2012

### Abstract

Goal models have been found to be useful for supporting the decision making process in the early requirements phase. Through measuring contribution degrees of individual solutions to high-level quality goals and combining them with preference statements, it is possible to compare alternative solutions of the requirements problem against each other. But where do contribution degrees and priority statements come from? In this paper we describe how full application of the Analytic Hierarchy Process (AHP) can be used to quantitatively assess influence relationships and priorities based on stakeholder input and how we can reason about the result in order to make informed decisions. An exploratory experiment shows that the proposed procedure is feasible and offers evidence that the resulting goal model is useful for guiding a decision. It also shows how general and situation-specific knowledge co-exist within goal models, a phenomenon that may need to be studied further in the context of eliciting such models.

# 1   Introduction

Goal models have been recognized as a promising approach for modeling and reasoning about alternative solutions during early requirements [1, 2, 3, 4]. Through refinement hierarchies, such models represent alternative ways by which top-level stakeholder goals can be fulfilled. Each of the identified alternatives is assessed subject to quality requirements, through modeling the *influence* of low-level decisions to each such quality. Based on the relevant significance of these influences, solutions that best match stakeholder priorities are identified and pursued for further analysis. A wealth of such reasoning techniques has being proposed in the literature, employing a variety of measuring scales and reasoning techniques [5].

However, the problem of eliciting/identifying the influence measures, that is, the question *where the numbers come from* [2], has not enjoyed equal attention in the research community. How can real attitudes and influence assessments of stakeholders be elicited and formalized in concrete measures and what factors can affect such assessments? How can these measures be aggregated and used for making a decision and why do we believe that that decision truly represents stakeholder input?

In this paper, we explore the use of the Analytic Hierarchy Process (AHP) [6] for eliciting and aggregating quantitative influence measures within semi-formal goal models, aiming at supporting decision making during early requirements analysis. In particular, we are exploiting the similarity between goal hierarchies which are central in goal models and criteria hierarchies, which is how AHP organizes priority elicitation. Thus, we propose that, for goal models that meet certain structural characteristics, the problem of eliciting influence links and making a decision can be seen as an aggregation of a number of individual standard AHP decision problems. Solving each of these problems is effectively a way to assess the influence measures within the goal model and eventually decide over a solution through a simple reasoning procedure. This way, goal modelers can pose a stronger validity argument both for the resulting representation and for the decisions it yields.

To understand different aspects of this synergy in practice, we also conducted an exploratory experiment. We gave the classic meeting scheduling problem in form of a goal model to a group of participants and asked them to envision themselves in given scheduling scenarios. They followed AHP for eliciting influence measures based on each scenario. We observed the participants for their consistency in their responses within and across scenarios as well as subsequent recognition of their preferences and alternatives of choice. Amongst the findings are that,

for a small problem: (a) the process is applicable and soon converges to consistent or near-consistent results, (b) participants generally recognize the result of their preference input and preferred choice, (c) situational characteristics (i.e. the particular decision making scenario) may influence some but not all of the priorities and (d) quantitative representation of influence does not seem to impair diagrammatic reasoning about influences and preferred decisions compared to qualitative one.

The paper is organized as follows. In Section 2 we provide an overview of goal models and the AHP process and in Section 3 we show how we apply the latter to assess the influence measures of the former. In Section 4 we present the design and results of our exploratory study. Then, in Section 5 we discuss related work and conclude in Section 6.

## 2 Background

### 2.1 Early Analysis using Goal Models

Goal models allow representation and reasoning about how goals of stakeholders relate with and influence each other. Such a model, adapted from [1], can be seen in Figure 1. The model describes goals pertaining to the classic Meeting Scheduling problem, and ways by which such goals can be achieved. The notation makes use of core concepts that can be found, in one form or another, in most dialects of the *i\** family [7, 8].

In the model, high-level hard-goals – the ovals – are recursively decomposed into lower-level ones. The decompositions are modeled through two kinds of links, the *AND-decomposition* links and the *OR-decomposition* links. All children of AND-decompositions need to be fulfilled for the parent goal to be considered fulfilled. Respectively, fulfillment of just one child of an OR-decomposition suffices for us to consider its parent fulfilled. Thanks to the existence of OR-decompositions the goal tree implies a great number of alternative ways by which the root-level goal can be satisfied. These are simply solutions of the AND/OR tree.

Soft-goals – the cloud-shaped elements – also represent goals, but ones for which there is no cut-and-dry satisfaction criterion. As such, their satisfaction is assessed on the basis of satisfaction of other goals. This is traditionally expressed with *contribution links*: a positive (respectively, negative) contribution link drawn from a goal to another means that evidence of satisfaction of the former constitutes

3

evidence for the satisfaction (resp. denial) of the latter. We will call both links with the more general term *influence* links as they show how satisfaction of one goal is understood to influence satisfaction of the other.

Influence links allow us to view soft-goal satisfaction as a *criterion* for assessing the satisfaction of other higher-level goals. In the Figure 1, soft-goal *Minimal Effort* is influenced by goals *Minimal Collection Effort* and *Minimal Matching Effort*. Our knowledge of satisfaction of *Minimal Effort* is thus, assumed to depend exclusively on what we know about the satisfaction of *Minimal Collection Effort* and *Minimal Matching Effort*, and, as such, the latter are the criteria for assessing satisfaction of the former. Thus, soft-goals that vary in their level of specificity, naturally form hierarchical structures in which soft-goals of the lower level serve as criteria for assessing the fulfillment of those at the higher level. At the lowest level, influences to soft-goals originate from hard-goals. Thus, different solutions of the AND/OR decomposition tree imply a different influence to lowest-level soft-goals and, in turn, the entire soft-goal hierarchy. Conversely, different desiderata regarding satisfaction of the soft-goal hierarchy, imply different criteria that need to be met and, in turn, a different goal alternative to be considered.

But what measures can we apply to model influence level and how can they be aggregated when multiple such influences target the same soft-goal? How can the measures be elicited? And how can an informed decision be made based on this? In this paper, we show how we can appeal to the Analytic Hierarchy Process (AHP) to both elicit influence measures and appropriately aggregate them for the purpose of making decisions within goal models. By doing this we both (a) allow goal models to be the basis for structuring and conducting an early requirements decision problem and (b) supply the representational result (i.e. the resulting goal model with all its elicited weights) with a stronger validity argument. Before we see how this is possible, we take a closer look at AHP in the following subsection.

## 2.2 The Analytic Hierarchy Process

The Analytic Hierarch Process (AHP) [6] is a decision support method aiming at quantifying relative priorities for a given set of alternatives based on the subjective judgment of a decision maker. It has been widely used for decision making in many areas like economics, social and management science [12] and in requirements engineering [9]. In AHP complex problems are modeled in a hierarchical structure showing the relationships of the main decision goal, its satisfaction criteria and different levels of sub-criteria thereof. More specifically a sequence of

Figure 1: A goal model

five steps is performed which we describe below.

**Constructing the Criteria Hierarchy.** The first step is the formation of a criteria hierarchy. The top level element of the hierarchy represents the goal or the objective that needs to be met, in form of a problem statement. This is the *decision goal*. For example, if our decision problem is to e.g. purchase a bicycle, the root decision goal would be *Choose Best Bike*. At the same time a set of *alternatives* that can potentially fulfil the objective are identified. In our bicycle example, three such bicycle models may exist, each with different features. Subsequently, one level under the top goal, the major criteria for assessing fulfillment of the decision goal are defined in broad terms. For example, if it is a bike for long commutes we may be interested in its *comfort* and *durability* as top-level criteria. Each criterion may be broken down to lower-level ones depending on how much detail is needed. In our bicycle selection problem, the comfort criterion involves two sub-criteria: *weight* and *shock absorption* (against e.g. pavement bumps). Once the criteria hierarchy is complete, the alternatives are connected to each of the
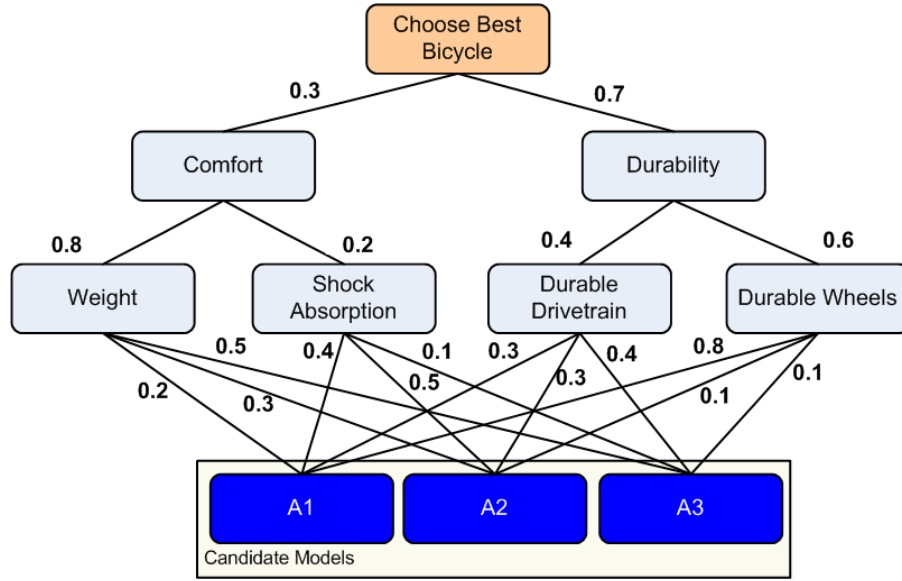
5

Figure 2: An AHP decision hierarchy

leaf-level criteria, forming the bottom level of the hierarchy. Thus, the resulting model, such as that of Figure 2 describing our simple bicycle selection problem has three levels: decision goal, criteria and alternatives. Note that, for simplicity, we focus on acyclic hierarchy structures in this paper, though our proposal applies equally well to hierarchies with undirected cycles.

**Pairwise Comparisons and Comparison Matrices.** At the next stage, comparisons are performed at each level of the hierarchy, in order for the relative importance of the sibling criteria or alternatives to be assessed. In such comparisons the elements are compared *with respect to their parent element* in the hierarchy. Thus, the top level criteria are compared with each other with respect to their relative importance in achieving the top objective. In our bicycle example we compare comfort with durability with respect to how important is each in selecting the best bike. Sub-criteria of every level are likewise compared with respect to their importance for fulfilling the parent criterion. Thus, the durability of the drivetrain is compared with the durability of the wheels with respect to how important each subcriterion is for the overall durability of the bike. At the leaf level, alternatives are compared with respect to their fulfillment of each of the lowest level sub-criteria. Overall, for the decision model shown in Figure 2, seven comparisons need to be made: one for comparing top level criteria with respect to the decision goal, two more for comparing the second-level sub-criteria and

four for assessing the "goodness" of alternatives with respect to each of the four sub-criteria.

The comparisons are performed in a pairwise fashion. An $n \times n$ matrix is constructed, the *comparison matrix*, where each row and column represents each of the $n$ elements to be compared. Each cell of the matrix hosts the result of the pairwise comparison between the corresponding elements. The decision maker fills each cell with a value expressing: (a) in the case of criteria, the relative importance of one sub-criterion (row) over the other (column) with respect to the parent criterion or (b) in case of alternatives, how much better one alternative is judged to satisfy a given leaf-level criterion than the other. In both cases, it is again important to notice that each pairwise comparison is performed with respect to a criterion and the question to the stakeholder must emphasize that. The values are chosen from the set $\{1,3,5,7,9\}$, expressing equal, moderate, strong, very strong, or extreme importance (in case of criteria) or betterness (in case of alternatives) of one element over the other. Thus, back to our bicycle purchase problem, we may say that shock absorption is strongly more important than weight with respect to comfort, hence 5. That number would be different if the parent criterion was different (e.g. performance). Likewise, at the level of alternatives, we may say that we very strongly consider bike model A1 to be more suitable than bike A2 with respect to weight (e.g. the former is aluminium-framed and the latter is steel-framed). Again, with respect to other criteria, such as shock absorption, the preference may be the reverse.

**Calculation of local weights.** In this step, the comparison matrices are transformed into weighted priority profiles amongst the involved items, which we call the *local weights*. Thus, at each level of the hierarchy tree, each of the elements at that level acquires a real number from the interval [0,1] representing its relative importance (for sub-criteria) or relative suitability/"betterness" (for alternatives) of the element compared to its sibling elements and with respect to the parent criterion. Hence local weights represent the influence share of each element to their parent one. The transformation follows the eigenvector method (EVM) – we refer the reader to [6] for details. In Figure 2 such numbers appear as labels on the links that connect alternatives or sub-criteria to higher level criteria.

**Aggregation of local weights into global weights.** Once the local weights of elements are obtained at different levels of the hierarchy, they are aggregated to obtain *global weights* of the decision alternatives (elements at the lowest level). To calculate the global weight $global(a)$ for alternative $a$ we use the formula:

$$global(a) = \sum_{c_l \in C_l} ( \prod_{c_i \in C_l^{root}} local(c_i) \times local(a, c_l))$$

where $C_l$ is the set of criteria $c_l$ subject to which $a$ has been compared, $local(a, c_l)$ the resulting local weight of each such comparison, $C_l^{root}$ is the set of criteria $c_i$ that are ancestors to the criterion $c_l$, $local(c_i)$ being their local weights. Intuitively, the formula indicates that we perform two steps: (a) calculate the global weight of the leaf-level criteria, (b) calculate the global weight of each alternative. The global weight of each leaf level criterion is calculated by multiplying the local weight of each of its ancestors. In the bicycle example of Figure 2 the global weight of the leaf criterion *Durable Drivetrain* is $0.4 \times 0.7 = 0.28$. The global weight of each alternative is calculated as follows. First, collect each of the criteria with respect to which the alternative has been assessed. Second, multiply the global weight of each criterion with the local weight the alternative has with respect to this criterion. Thirdly and finally, add up all the resulting weights. In the bicycle example of Figure 2 the global weight of alternative A1 is $(0.3 \times 0.8) \times 0.2 + (0.3 \times 0.2) \times 0.4 + (0.7 \times 0.4) \times 0.3 + (0.3 \times 0.3) \times 0.8$.

**Decision.** The global weights of the alternatives represent the rating of the alternatives in achieving the decision goal. The result tells us not only which alternative is more important – and we may reasonably want to pursue it as such – but also to what degree.

We next turn our focus on how the AHP process we describe above can be used for assessing influence links and making decisions in goal models.

# 3 Eliciting the Influence Structure and Making Decisions

## 3.1 The Process

Application of AHP to acquisition and aggregation of influence degrees in goal models is based on two principles: (a) every OR-decomposition in the goal model constitutes a separate decision problem and (b) the soft-goal hierarchy plays the role of the AHP criteria hierarchy for making each such individual decision. The aggregated result of solving each and every such decision problem is a complete alternative in the goal tree (i.e. a solution of the AND/OR tree). More specifically let us assume that we are given a typical goal tree such as that of Figure 1. The
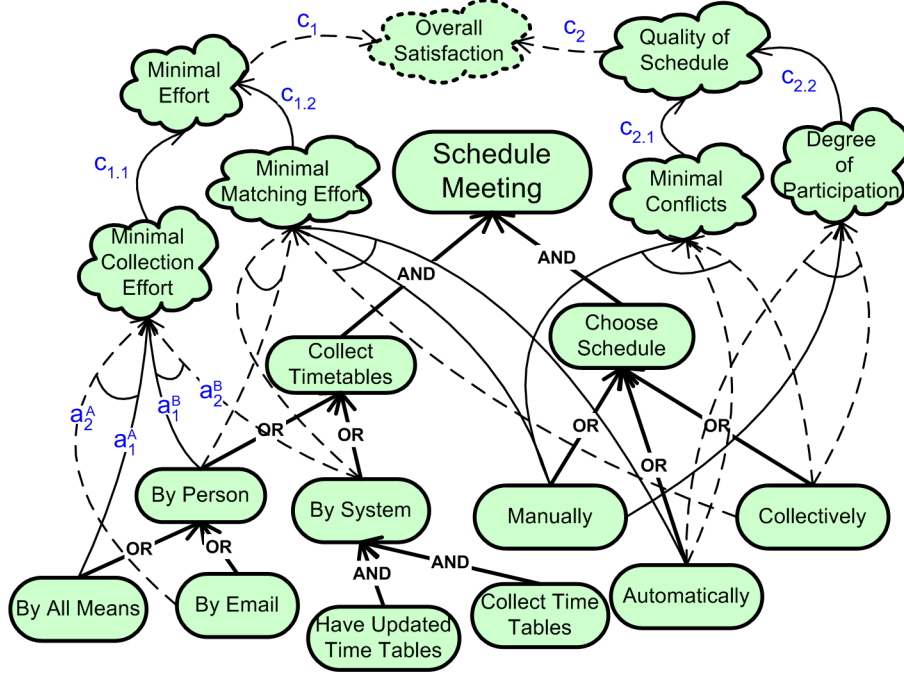
Figure 3: Re-arranging and Enriching Influence Links (added elements appear in dashed lines/borders)

model involves a hierarchy of soft-goals, the AND/OR hierarchy of hard-goals and an initial assumption by the modeller on how these are connected through influence links. We then perform the following steps.

**Re-arranging and Enriching links.** We must first ensure that influence links connecting the hard-goal decomposition with the soft-goals are restricted to ones that connect children of OR-decompositions to leafs of the soft-goal hierarchy. To ensure those two conditions we work as follows. Firstly, for links that originate from a goal that is a child of an AND-decomposition, we move the origin of the link to the closest ancestor in the tree that is a child of an OR-decomposition. We remove the link if such ancestor does not exist, since, in this case, the influence does not serve any purpose for decision or alternatives analysis. Secondly, once this step is done, influence links that point to a non-leaf soft-goal are removed and replaced with links that originate from the original hard-goal to each of the leaf-level descendants of that non-leaf soft-goal. The assumption behind this practice is that an influence to a soft-goal must be explained through influences to one or more of its lower-level soft-goals (i.e. there are no influence means that are not

being modelled).

Once the initial contribution links have been established then they are enriched as follows. For every OR decomposition child, in addition to the influence links that have already been defined for it, connect it also with all the soft-goals which any of its siblings influence. In other words, if a sibling $g'$ of the OR-decomposition child $g$ influences a soft-goal $s$, then $g$ itself must send some influence to $s$. In the end all children of an OR-decomposition influence exactly the same soft-goals.

Finally, if there are multiple soft-goal hierarchies, a soft-goal "Overall Satisfaction" is introduced to the goal graph and placed as the parent of the roots of each of those hierarchies. This way, soft-goals are all organized in a unique tree structure. The result of applying this step to the goal model of Figure 1 can be seen in Figure 3 – ignore the weight labels for the moment, as the influences are unlabeled at this stage. In the figure, the links and elements that have been added are drawn as dashed.

**From hard-goals and soft-goals to criteria and alternatives.** At this stage, the soft-goal hierarchy is isomorphicaly mapped to an AHP criteria hierarchy. Further, each OR-decomposition is considered as a separate decision problem in which each child of the decomposition a distinct alternative. All these decomposition problems however are assumed to share the same set of criteria. Figure 4 illustrates how the goal model of Figures 1 and 3 yields three AHP decision problems sharing the same criteria hierarchy.

**Acquiring Local Weights.** At this step each of the decision problems is solved through AHP pairwise comparisons as described above. Firstly, the local weights in the criteria hierarchy are calculated top-down. Thus, in Figure 4 pairwise comparison between *Minimal Effort* and *Schedule Quality* is performed with respect to *Overall Satisfaction*. This results to local weights $c_1$ and $c_2$. Another two comparisons give us the local weights for the lower level criteria $c_{1.1}, c_{1.2}, c_{2.1}$ and $c_{2.2}$. Secondly, each decision problem is completely solved – following an order we will discuss below. For each problem, each of its alternatives are compared subject to their goodness with respect to each of the criteria they have been associated with. Back to Figure 4, manual timetable selection of alternatives (*By all means* vs. *By email*) are compared with respect to *Minimal Collection Effort*. Likewise, timetable collection alternatives *By Person* vs. *By System* are compared with each other with respect to both *Minimal Collection Effort* and *Minimal Matching Effort*. The global weights of each alternative are calculated as described earlier and, again, the highest score indicates the alternative of choice.

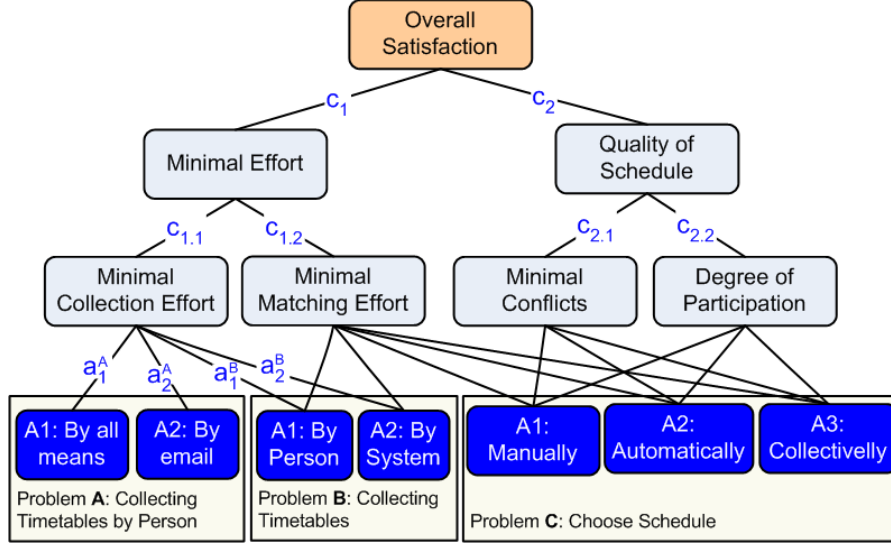**Mapping Result to Goal Model.** The result of AHP analysis has at least two

Figure 4: AHP Decision Model for the Meeting Scheduling Problem

uses with respect to the goal model. Firstly, by mapping the optimal alternative for each individual problem to the OR-decomposition from which the problem was generated, we get a solution to the AND/OR tree that is optimal. Secondly, given the one-to-one correspondence of local weights assessed in AHP and influence links existing in the (enriched) goal model, we can use the former to simply label the latter. Thus, the local weight that expresses the goodness of alternative *By Email* with respect to *Minimal Matching Effort* becomes the label for the influence link from the hard-goal *By Email* to the soft-goal *Minimal Matching Effort*. Some of the weight labels seen in Figure 4 are transferred to Figure 3 to illustrate how the mapping is done.

## 3.2 Discussion

Two aspects of the technique that require some more commenting are (a) the ordering of the decision problems and (b) the reuse of the hierarchy.

**Order.** Solving of each of the AHP problems is ordered based on the hierarchy of the corresponding OR-decompositions in the goal model, starting from the leafs and moving to the top. The reason is that whenever an OR-decomposition is of higher level, each of its alternatives may have descendants that are also OR-decompositions. In such a case, an alternative OR-subgoal may implicitly be a collection of alternatives. In Figure 1, for instance, the alternative to collect

11

timetables *By Person* is in fact two alternatives, one *By Email* and one *By All Means*. Unless we decide on one of those two lower-level alternatives, the comparison between *By Person* and *By System* is problematic. In practice, arbitrary such nestings of OR-decompositions may occur. To address this, we consider solving the low level OR-decompositions first and, when we have a higher level OR-decomposition, we mention the optimal solution of each of its alternatives in our question to the stakeholder. In our example, if *By Email* is found to be the preferred solution for the low level, at the higher level we compare *"By System"* with *"By Person, assuming they do it By Email"*.

**Reuse of Criteria Weights.** Furthermore, notice that the weights that are elicited for the criteria hierarchy are re-used for identifying the optimal alternative for each problem. To see why this choice, which saves significant effort, is justified, consider that, for example, our preference between *Minimal Effort* and *Quality of Schedule* should not depend on the particular aspect of the meeting scheduling solution we are trying to optimize, but rather on circumstances pertaining to the general scheduling problem.

# 4 An exploratory study

## 4.1 Overview

To assess how the above process applies in practice we performed an exploratory experiment. The goals of our exploration are three. Firstly, we want to assess whether the process is at all applicable. AHP has been applied to a wide variety of domains from health care to software specifications. But can we also consider concepts that have traditionally been researched and developed in the *i\** culture (e.g. goals, soft-goals or goal alternatives), and use them as AHP concepts? Secondly, if we populate the model with numbers we acquire through AHP, is the visual model comprehensible in a way that it can, for instance, allow gauging the optimal decision? Finally, what influences the elicitation process? That is, how do stakeholder responses change under different envisioned situations? To study these we subjected a number of experimental participants to a series of tasks, as we describe below.
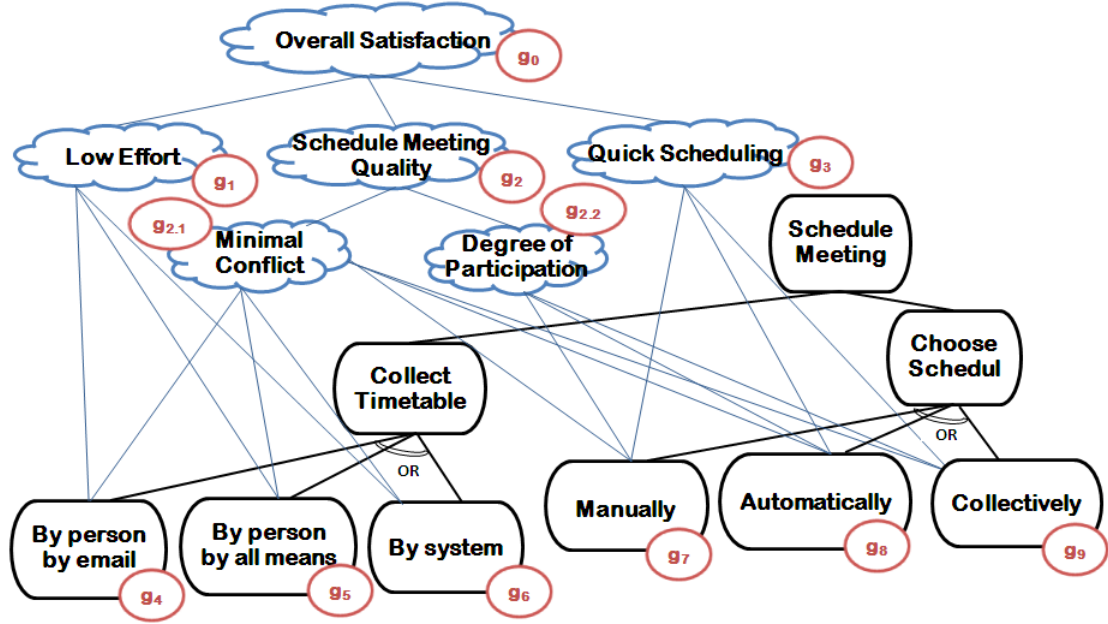
Figure 5: Experimental Goal Model

## 4.2 Experimental Design

Participants in the study are 10 graduate students of Information Technology, recruited from the first author's class. A minor (3%) part of their mark was offered as inducement. They had experience in goal modeling techniques through academic work they performed for the course; some of them had such knowledge from earlier undergraduate courses. The second author conducted the experiment in her office and supervised it the entire time. Participants were asked to perform three different tasks for two different experimental conditions, i.e. in a within-subjects design. The tasks correspond to three different instruments (essentially questionnaires) which we will refer to as Instruments A, B and C and describe in more detail below. Samples of the instruments can be found in the Appendices at the end of the paper.

In each of the two conditions participants were given two different scenarios in which a meeting needs to be scheduled based on the goal model of Figure 5 – a simplified version of the model we saw in our earlier examples. The different scenarios allow us to measure the degree by which information about the situation/context of goal fulfillment affects the elicitation of the influence measures. Thus, in the first scenario (Scenario A), participants envision themselves in a sit-

uation where they work on an academic course project with a group of fellow students, they have a very close deadline, and need to meet in order to resolve an unexpected problem. In the second scenario (Scenario B), they are again students in a class who want to organize an end-of-term social gathering with their fellow students for a date that is one month ahead.

Before presentation of AHP and the scenarios, a preliminary task is performed through Instrument A.

**Step 1 (Instrument A).** This instrument aims at assessing the ability of the participants to look at an arbitrary goal diagram with predefined numeric labels and assess what the optimal goal alternative is for that diagram. More specifically, the instrument contains two pairs of diagrams, each pair displaying a different OR-decomposition of the model of Figure 5. For each pair, one model contains numerical values of the interval [0..1]. The values are crafted by the authors to make visual reasoning about the optimal alternative non-trivial, through maximizing conflicting influences. The other model contains qualitative contribution links "++","+","?", "−" and "−−". The quantitative labels are derived from the quantitative ones in the first model of the pair through fragmenting the continuous space [0..1] into 0.2-step intervals ([0,0.2),[0.2,0.4),..., [0.8, 1.0]). Therefore, the qualitative model is roughly a discretization of the quantitative one and, as such, we assumed it to have the same optimal goal alternative. For each model in each quantitative-qualitative pair the participants are asked to choose the optimal alternative based on the given influence measures, through visual reasoning and without performing any pen and paper calculations.

**Step 2 (AHP Training).** The participants are then given a brief introduction to the AHP method and detailed instructions for filling out the comparison tables; an example with diagrams and a short training exercise are employed to ensure understanding. This training was administered by the second author throughout.

The following tasks are performed for each of the two scenarios:

**Step 3 (Fill-in Comparison Matrices).** The participants are given the scenario and are asked to imagine themselves as part of the described situation. Based on that situation, they then fill in the AHP comparison matrices. In total, seven (7) comparison matrices are filled for each scenario, as needed for the goal model of Figure 5 (two for the soft-goal hierarchy, two for the first OR-decomposition and three for the second OR-decomposition).

**Step 4 (Calculation).** For each resulting comparison table the local weights are calculated by following the eigenvalue method. The optimal alternative is calculated based on the aggregation rule we discussed above. Two instruments, *Instrument B* and *Instrument C* are then given to the participants.

**Step 5 (Instrument B).** In this instrument we want to measure to what degree participants are able to recognise the local weight profile that best matches their comparison matrix input, according to the eigenvector calculation. The participants are given four (4) comparisons, that are randomly selected from the seven they completed before (Step 3). For each comparison, a number of different options for local weights are provided to the participant. The options are generated as follows: one of them is the one that results from applying the AHP transformations of Step 4 (i.e. the "correct" one) and the rest are all possible permutations of that one. The participant is asked to select the one that best describes their own inputs in the comparison matrix.

**Step 6 (Instrument C).** In this instrument, we aim at measuring how well the participants' perception of the optimal alternative of the entire goal model matches the AHP-based aggregation of their inputs. Thus, the participants are given four (4) different alternatives for the model of Figure 5 and are asked to select the one that they think best matches their overall priorities that they have been specifying earlier in Step 3. One of the alternatives is the optimal based on participant's input and the AHP aggregation procedure we described, another two are partially optimal (they miss one of the two OR-decompositions – Figure 5) and the fourth is totally incorrect (misses both decompositions).

Steps 2-6 above are performed in two repetitions (for both scenarios) with distance of about 1 month from each other. In the second session, when Step 1 also takes place for the first time, the administrator is actively detecting inconsistent responses and requests the participants to revise their preferences into more consistent ones (without, of course, dictating the preference per se). This practice of resolving inconsistencies with participants is legitimate and recommended in AHP, as it results in more valid inputs. Thus, all the data we report on come from the second session.

## 4.3 Results

We now take a look at the main results obtained from the experimentation and assess them with respect to the original goals of our exploration.

### 4.3.1 General Applicability

Our first evaluation goal is to assess whether application of AHP's pairwise comparisons is relevant and doable in goal models. We use two criteria to asses this:
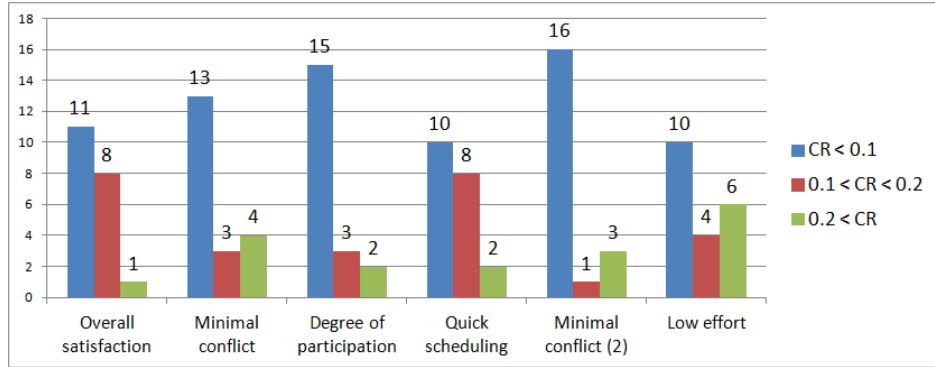
Figure 6: Frequencies of Consistency Ratio Levels

firstly, the degree of internal consistency of each individual response, and, secondly, the degree by which the participants are able to recognize their own preferences by looking at the output.

An initial observation is that all participants were able to complete the process without explicating any reservations about the logic or the procedure they followed. Further, the *consistency ratio (CR)* is calculated for each response. The CR tells us to what extend the participant has entered conflicting priorities in the comparison matrix. Calculation details can be found in e.g. [11]. According to Saaty [6], a value of CR more than 0.1 indicates that judgments should be elicited once again from the participant until she gives more consistent judgments – which is why two sessions were performed. Yet, participants may refuse to completely lift inconsistencies even if they acknowledge their existence, which happened in our case and is probably why CRs greater than 0.1 seem to be frequent in practice [11].

Thus, after the second session, the CR values that are greater than 0.1 are 25.7% for scenario A and 38.6% for scenario B, mostly by small amounts. The bar-chart of Figure 6, shows the frequency of CR values that were below 0.1, from 0.1 to 0.2, and above 0.2, for each comparison exercise, for both scenarios. In all cases, more than half of the CRs are below 0.1. The highest CR observed was 0.52.

To assess our second goal, i.e. whether participants are able to recognize their output we look at the results from Instruments B and C. In Instrument B, in which the participants were asked to select the local weights that best match their input in the comparison matrix, in 87.50% and 80% of the responses, for scenario A and B respectively, the participants have correctly selected the local weight pro-

16

files that correspond to their pairwise input. In other words, participants were largely able to recognize the numeric priority model that corresponds to their own judgment of the situation as guided by AHP. In Instrument C, where participants were asked to select the overall goal alternative that they think best matches their priorities, in scenario A 80% of them have selected the alternatives that the AHP aggregation process indicates as preferred. From the remaining, 20% selected one that was partially correct (i.e. got one of the two OR-decompositions right) and none (0%) selected one that was totally incorrect. The corresponding percentages for scenario B are 70% (totally correct) 20% (half-correct) and 10% (totally incorrect).

Do these successful responses occur by chance? To investigate this, the binomial test is applied. In Instrument B, the result of the two-tailed binomial test for each scenario was statistically significant ($p < 0.01$ for questions with six options; $p < 0.02$ for questions with two options). The test is statistically significant for Instrument C as well ($p < 0.05$ for four-option questions). Hence, it is highly unlikely that these responses were successful by chance.

*Discussion.* In analyzing the results we first find the fact that we had consistent or near-consistent inputs as indicative both of relatively reliable data (participants don't answer randomly) and of some basic sanity of the elicitation approach. In Instrument B, the fact that participants are able to recognize *their own* preferences in the profile seems to support the appropriateness of the pair-wise comparison method for the purpose we are using it. In other words, what the participants observe in the resulting local priority weights is in agreement with their pair-wise input, which seems to validate use of the latter (pair-wise comparisons) to produce the former (local weights as measures of influence share). Furthermore, the ability to recognize not just the local weight but the entire goal alternative (Instrument C) offers evidence on the intuitiveness of AHP-based approach for aggregating local weights: the participants' intuition of what the preferred alternative is coincides to a great degree with the AHP-dictated aggregation of individual local weights, which, in turn, as we saw, are remarkably consistent with the participants' pair-wise inputs. In other words, the way the parts are aggregated yields a whole that is consistent with the participants intuition.

### 4.3.2 The role of representation of weights: quantitative vs. qualitative

At the second stage we study the comprehensibility of the weighted measures when placed on the goal model. As we saw, Instrument A offers participants partial goal models with the influence measures completed by the researchers, in two

versions: a numeric, having the exact results of the AHP process, and a qualitative, having a discretization of the numbers. Of the three options they are given, the participants are asked to check the one that is optimal. The result is compared with what the AHP-based aggregation decides as optimal. The participants are also asked to rate their confidence in their decisions in a 1 to 10 scale.

The results show that, in the quantitative models, in 95% of the responses, participants have selected the alternatives correctly with an average confidence of 83%. In the qualitative models, 75% of the responses are correct and with 74% confidence. To investigate whether the participants are just guessing, the binomial test was again applied: the p-values are $< 0.01$ for both cases, which are both statistically significant for $\alpha = 0.05$

*Discussion.* In this exercise we focus on one of the possible *uses* of the goal model, namely its ability to support a decision. In this context, we investigate the correspondence between, on one hand, the guidance that the visual aspect of the goal model gives to the user regarding how influences are aggregated and, on the other hand, the mathematical result of the AHP-based aggregation approach. The result suggests a strong connection between the two for both kinds of labels. This allows us to hypothesise that, for a simple model, if we assume validity in the AHP decision making approach, then if its result is placed in the goal model it may also allow visual reasoning over the resulting graphical representation. Another important observation is that the result does not offer evidence that quantitative labelling impairs the effectiveness of the visual reasoning or even the confidence of the respondents. In other words, our exploratory experiment suggests that adding weights derived from an AHP exercise into a goal model, either as-is or after turning them into qualitative labels, may aid visual reasoning.

Note that, from a validity standpoint, we must keep in mind that in the instrument the measures were constructed by the researchers to make the process more challenging – this may have introduced a bias to the opposite direction in an unknown way. Nevertheless, the qualitative vs. quantitative comparison seems to be less exposed to this threat.

### 4.3.3 The influence of the scenario

As a final step, we investigate the role of the scenarios to participant input. One should expect that goals are more or less important from each other depending on the situation in which the goal model is used. But is the participants' input in fact influenced by the scenario? If yes, which part of her input is influenced and which part is not? Is the influence justified or is it a result of an unwanted cognitive bias?
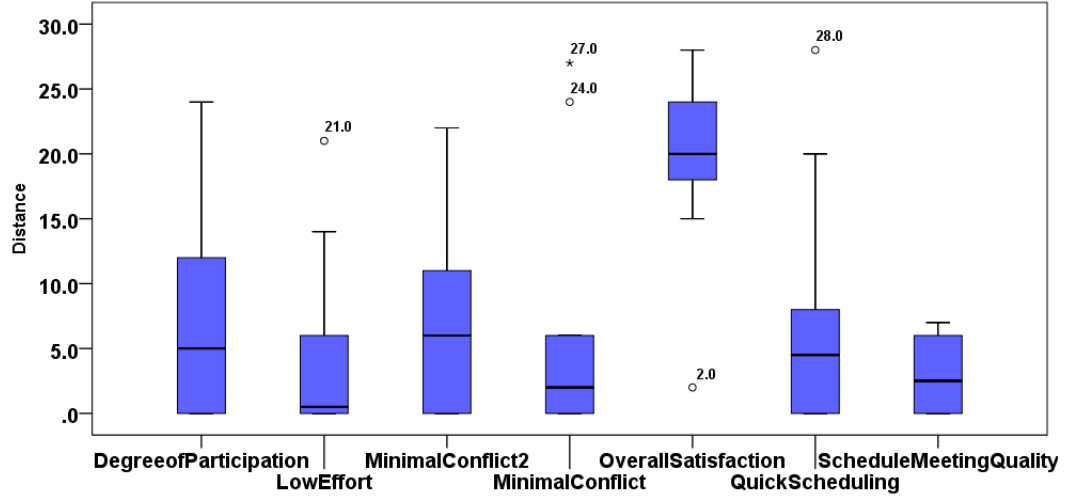
Figure 7: Similarity between comparison matrices across scenarios

To examine this we used two measures. Firstly, we devised a distance measure between comparison matrices (as we could not find such in the literature). The measure simply calculates the distance between each of the individual inputs in the comparison matrix and sums them up. Thus, for a particular comparison, the more the influence of the scenario the more the distance between the two corresponding inputs. Figure 7 shows a box-plot in which the distances between responses are summarized per goal/comparison for all participants. It is notable that the comparison corresponding to the highest level goal appears to be more dependent on the scenario than the other ones.

We then moved on to assess the statistical significance of the observed difference. Recall that each comparison results in a profile of two or three local weights (depending on the number of items under comparison). For each such individual weight, we tested for significant differences with respect to the scenario change. To avoid Normal distribution assumptions, we used the Wilcoxon test. Twenty such two-tailed tests are performed: out of the seven (7) total comparisons, six (6) are of three (3) choices and one (1) is of two (2) choices. In Table 1 the Wilcoxon value for each test can be seen – values that are below our significance threshold $\alpha = 0.05$ can be considered significant.

*Discussion.* Both the descriptive analysis and the Wilcoxon test give us a strong impression that some influence measures in goal models depend on the perceived situation in which goals need to be fulfilled. What is interesting is

19

that not all such influences are volatile, although they all have the same status in the goal modeling language. While the concept of (weighted) *preferences* has been studied as an exogenous component in goal modeling languages [4], our hypothesis is that the distinction between what constitutes a preference, influenced by the situation, and common knowledge, that applies in all situations, may also exist within goal models of the *i\** family. In our result, we found that the top level goal is influenced by the scenario the most, indicating that higher-level influences may have an inherently more preferential/subjective nature compared to low level ones that describe commonly accepted causal relationships.

Nevertheless, our observation from interacting with the participants, shows another issue pertaining to elicitation: some participants seem to be influenced by the scenario even for aspects where this should not have been the case. For example, when asked to compare timetable collection means with respect to *Low Effort* it became obvious that in scenario A (which involves urgent scheduling), the due date constraint influenced some participants to pick the alternative which allows scheduling as soon as possible regardless of the effort level. This seems to suggest that analysts should be aware that cognitive biases may accompany weight elicitation: participants focus on a piece of information (scenario in our case) that influences their decisions in ways that it should not. In the particular case, the element with respect to which a comparison is supposed to be made (*Low Effort*) is replaced by dominant characteristics of the situation (e.g. urgency). When such biases were detected in the experiment, the intervention of the second author was often required to carefully (i.e. without dictating a response) remind the participants that the comparison is performed with respect to a very specific criterion. Hence, although we believe that the end-result is not dominated by such biases, the phenomenon is present and needs to be studied more.

## 4.4   Discussion and Validity Threats

While we already mentioned some specific validity threats we now turn our focus on the general validity of the experiment. The fronts in which *external validity* can be challenged are several. Firstly, we are using only one domain (meeting scheduler) and a particular goal model. Would different goal models for the same domain provide different findings? How about different domains? Further, graduate students of Information Technology seem to be a population that can potentially represent analysts, but cannot possibly represent arbitrary stakeholders, who have a different intuition of numbers and box-and-line models. Furthermore, the small number of comparison tables used in the experiment to limit the fatigue effect and

| $g_0$ | | | $g_{2.1}$ | | | $g_{2.2}$ | | |
|---|---|---|---|---|---|---|---|---|
| $g_1$(*) | $g_2$(*) | $g_3$(*) | $g_7$ | $g_8$ | $g_9$ | $g_7$ | $g_8$ | $g_9$ |
| $< 0.01$ | $< 0.01$ | $0.01$ | 0.89 | 0.50 | 0.89 | 0.92 | 0.46 | 0.92 |

| $g_3$ | | | $g_{2.1}$ | | |
|---|---|---|---|---|---|
| $g_7$ | $g_8$ | $g_9$ | $g_4$ | $g_5$ | $g_6$ |
| 0.89 | 0.89 | 0.92 | 0.60 | 0.09 | 0.40 |

| $g_1$ | | | $g_2$ | |
|---|---|---|---|---|
| $g_4$ | $g_5$ | $g_6$ | $g_{2.1}$ | $g_{2.2}$ |
| 1.00 | 0.68 | 0.89 | 0.60 | 0.60 |

Table 1: Wilcoxon tests per element per comparison

the small size of each (three by three maximum), prevents generalization to situations where a larger set of (larger) comparison tables are needed. Finally, the model per se is small. If it were larger, both Instruments A and C would perhaps yield significantly different results. To that end, though, we believe that the idea of dealing with each OR-decomposition as a separate problem makes size less problematic.

Standard measures were adopted to address *internal validity* threats. Thus, as we follow a within-subjects design, counterbalancing was applied to eliminate ordering effect. The comparison tables are given with different, random orders to each participant to eliminate biases in that regard too. Sessions are designed to not exceed 1.5 hours in duration and breaks are taken to avoid fatigue effects. Furthermore, despite our recruitment method (volunteering) we see the threat of self-selection bias irrelevant in our case. Nevertheless, we consider 10 participants to be a small sample despite the statistical significance that emerges by having each complete many exercises. As such, a larger sample size would offer us more confidence in future experimentation.

# 5 Related Work

The Analytic Hierarchy Process is a well-established multi-criteria decision support method. Several papers have compiled the AHP success stories in very different fields (e.g. [12, 13]). In the literature, there is also considerable commentary describing AHP as a broadly accepted method, based on firm theoretical foundation and, for many, being the most reliable approach to prioritization [14, 15].

In Requirements Engineering the method has been extensively studied as a tool for prioritizing specifications. Karlsson et al. [16] experimentally evaluated six different methods of requirements prioritization and found that AHP is the most promising technique in terms of providing trustworthy results, being fault tolerant, and including consistency measures – Karlsson and Ryan describe the application of AHP in requirements prioritization elsewhere [11]. Nevertheless, we could not find work in Requirements Engineering that makes full use of the hierarchical style of AHP criteria and applies it to goal hierarchies.

In goal modeling there is substantial work on reasoning about goal satisfaction and influence thereof ([17, 18, 4] – [19] for a survey). Such techniques typically model satisfaction of denial of soft-goals using qualitative or quantitative labels. Influence/contribution operators of various types model how satisfaction of one goal influences that of others and how groups of such, potentially conflicting, influences are aggregated. Despite the rigorous theoretical foundation and tool support these proposals offer, they leave the elicitation problem – how we elicit the influence labels and why we aggregate multiple and conflicting such the way we do – out of their scope.

To address the problem of number elicitation, the literature seems to pursue two directions. Letier and van Lamsweerde [2], propose a probabilistic interpretation of numbers in order to reduce the problem to one of assessing probabilities. This approach, however, cannot be assumed to be applicable to all domains and problems, it may sometimes require too detailed analysis for early requirements and, most importantly, it may not accurately capture stakeholder attitudes, preferences and likings. The latter can be addressed though the introduction of utility measures. We however could not find an extension to that line of work that makes effective use of utility functions/values – let alone address the problem of elicitation thereof. The second approach to validating influence input and propagation is provided by Horkoff and Yu [5], who propose an interactive conflict resolution technique to address the problem of aggregating influences. The approach seems to be geared towards scalable exploration and comprehension of the model. In comparison, our approach makes fewer presumptions about stakeholder attitudes, requiring the stakeholder to offer their input for every aggregation problem in the goal model using an established elicitation and local weight aggregation procedure. We find the possibility of combining the two methods very promising.

Finally, the problem of the effect of the visual/diagrammatic notation in comprehending goal structures has been addressed by Moody et al. [20] through appeal to current theories of perception and cognition. We believe this is a crucial line of investigation which, nevertheless, needs insights from the empirical front

in order to shed light on how people understand and use goal diagrams.

# 6   Conclusions

We presented an approach for applying the analytic hierarchy process in order to elicit influence measures in goal models. Soft-goal hierarchies of the goal models are treated as AHP criteria hierarchies and each OR-decomposition of hard-goals is treated as a separate AHP decision problem. The results are plugged back in the goal model, where the optimal goal alternative can be found following the AHP weight aggregation approach. An exploratory experiment offers evidence that AHP's pairwise elicitation is applicable to goals and that both the numeric/visual results and their AHP-based aggregation approach are comprehended by participants.

For the future we wish to tackle certain technical challenges, such as the presence of directed cycles or optional goals in the goal model. A more fundamental issue is that of the *semantics* of numeric measures as shares of influence and their comparison to the existing tradition of influence measures as absolute values of contribution. We believe that more extended empirical work will reveal what representation is more visually natural and how it is influenced by factors such as complexity and size of the model. Finally, we wish to understand the nature of influence and the kinds of theoretical constructs we need in order to distinguish better, for instance, influences that are attitudes versus influences that are "pure" domain assumptions.

# References

[1] J. Mylopoulos, L. Chung, S. Liao, H. Wang, and E. Yu, "Exploring alternatives during requirements analysis," *IEEE Software*, vol. 18, no. 1, pp. 92–96, 2001.

[2] E. Letier and A. van Lamsweerde, "Reasoning about partial goal satisfaction for requirements and design engineering," in *Proceedings of the 12th International Symposium on the Foundation of Software Engineering FSE-04*. Newport Beach, CA: ACM Press, November 2004, pp. 53–62.

[3] P. Giorgini, J. Mylopoulos, E. Nicchiarelli, and R. Sebastiani, "Reasoning with goal models," in *Proceedings of the 21st International Conference on Conceptual Modeling (ER'02)*, London, UK, 2002, pp. 167–181.

[4] S. Liaskos, S. McIlraith, S. Sohrabi, and J. Mylopoulos, "Representing and reasoning about preferences in requirements engineering," *Requirements Engineering Journal (REJ)*, vol. 16, pp. 227–249, 2011.

[5] J. Horkoff and E. Yu, "Comparison and evaluation of goal-oriented satisfaction analysis techniques," *Requirements Engineering (REJ)*, pp. 1–24, 2011.

[6] T. L. Saaty, *The Analytic Hierarchy Process*. McGraw-Hill International, New York, 1980.

[7] E. S. K. Yu, "Towards modelling and reasoning support for early-phase requirements engineering," in *Proceedings of the 3rd IEEE Int. Symposium on Requirements Engineering (RE'97)*, Washington D.C., USA, January 1997.

[8] D. Amyot and G. Mussbacher, "User requirements notation: The first ten years, the next ten years," *Journal of Software (JSW)*, vol. 6, no. 5, pp. 747–768, 2011.

[9] J. Karlsson, *Software requirements prioritizing*. Proceedings of the 2nd IEEE International Conference on Requirements Engineering (ICRE1996) Colorado, USA, 1996.

[10] S. Liaskos, R. Jalman, J. Aranda, "On eliciting preference and influence measure in goal models," School of IT, York University, http://www.yorku.ca/liaskos/Docs/AHPGoals.pdf, Tech. Rep., 2012.

[11] J. Karlsson and K. Ryan, "A cost-value approach for prioritizing requirements," *IEEE Software*, vol. 14, no. 5, 1997.

[12] O. S. Vaidya and S. Kumar, "Analytic hierarchy process: An overview of applications," *European Journal of Operational Research*, vol. 169, no. 1, pp. 1–29, 2006.

[13] W. Ho, "Integrated analytic hierarchy process and its applications – a literature review," *European Journal Of Operational Research*, vol. 186, no. 1, pp. 211–228, 2008.

[14] M. Bernasconi, C. Choirat, and R. Seri, "The analytic hierarchy process and the theory of measurement," *Management Science*, vol. 56, no. 4, pp. 699–711, 2010.

[15] M. S. H. Triantaphyllou, E., "Using the analytic hierarchy process for decision making in engineering applications: some challenges," *International Journal of Industrial Engineering: Applications and Practice*, vol. 2, no. 1, pp. 35–44, 1995.

[16] J. Karlsson, C. Wohlin, and B. Regnell, "An evaluation of methods for prioritizing software requirements," *Information & Software Technology*, vol. 39, no. 14-15,1998.
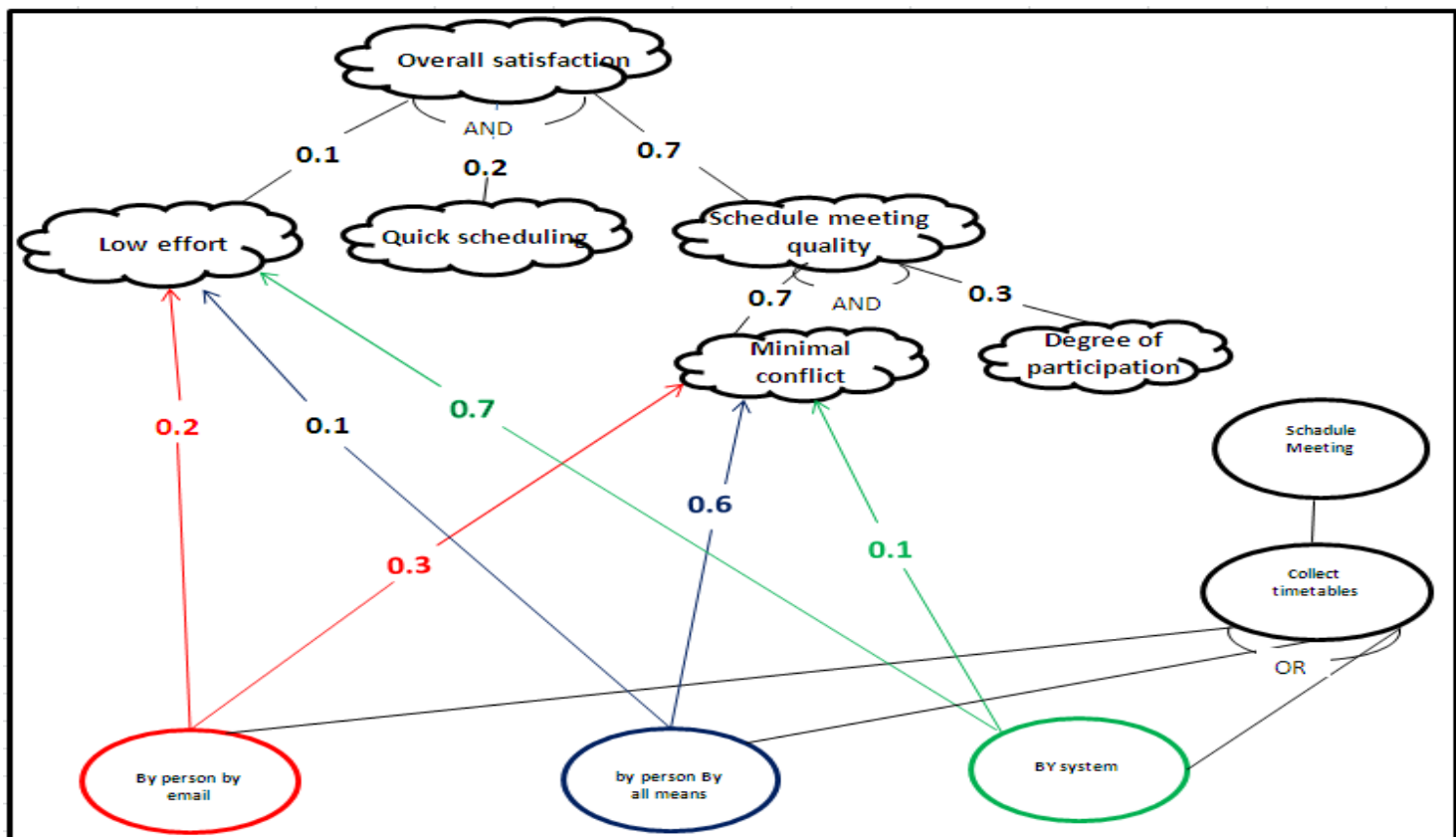
[17] D. Amyot, S. Ghanavati, J. Horkoff, G. Mussbacher, L. Peyton, and E. Yu, "Evaluating goal models within the goal-oriented requirement language," *International Journal Intelligent Systems*, vol. 25, no. 8, pp. 841–877, 2010.

[18] P. Giorgini, J. Mylopoulos, and R. Sebastiani, "Goal-oriented requirements analysis and reasoning in the Tropos methodology," *Engineering Applications of Artificial Intelligence*, vol. 18, no. 2, pp. 159–171, Mar. 2005.

[19] J. Horkoff and E. Yu, "Finding solutions in goal models: an interactive backward reasoning approach," in *Proceedings of the 29th International Conference on Conceptual modeling (ER'10)*, ser. ER'10.   Berlin, Heidelberg, 2010, pp. 59–75.

[20] D. L. Moody, P. Heymans, and R. Matulevičius, "Visual syntax does matter: improving the cognitive effectiveness of the i* visual notation," *Requirements Engineering Journal (REJ)*, vol. 15, no. 2, pp. 141–175, 2010.

# A   Instrument A (Step 1)

The quantitative models are constructed first. The numbers are "random" in a sense that they were put together arbitrarily in a way that it is not immediately obvious (based on our judgement) which alternative is the optimal. Each quantitative model comes with its qualitative counterpart, through the following mapping:

| | |
|---|---|
| [0.0,0.2) | $--$ |
| [0.2,0.4) | $-$ |
| [0.4,0.6) | ? |
| [0.6,0.8) | $+$ |
| [0.8,1.0] | $++$ |

How confident do you have that you chose the correct alternative?

Not at all confident  |  |  |  |  |  |  |  |  |  |  Very confident
                      1  2  3  4  5  6  7  8  9  10

Overall satisfaction

AND

--     -     +

Low effort     Quick scheduling     Schedule meeting quality

+    AND    -

Minimal conflict     Degree of participation

+     --     +

-     --     +     --

Schadule Meeting

Collect timetables

OR

By person by email     by person By all means     BY system

**How confident do you have that you chose the correct alternative?**

Not at all confident    |   |   |   |   |   |   |   |   |   |   Very confident

1   2   3   4   5   6   7   8   9   10

**Makes (++)**
**Helps (+)**
**Unknown (?)**
**Hurts (-)**
**Breaks (--)**

**How confident do you have that you chose the correct alternative?**

Not at all confident | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | Very confident

**How confident do you have that you chose the correct alternative?**

Not at all    |    |    |    |    |    |    |    |    |    |    Very
confident    1    2    3    4    5    6    7    8    9    10    confident

Makes (++)
Helps (+)
Unknown (?)
Hurts (-)
Breaks (--)

# B    Comparison Matrices (Step 3)

To facilitate consistency a graph-based representation was followed, where participants were labelling each edge with (a) the direction of the preference (by drawing the tip of the arrow) and (b) the number expressing relative importance of one element over the other (e.g. 3, 5, 9 etc.).

## Scenario

Suppose that the due date for the group project for ITEC6970 is tomorrow, and the project is still not finished .There are some unclear points, and you need to schedule a meeting with your group to finalize these points.
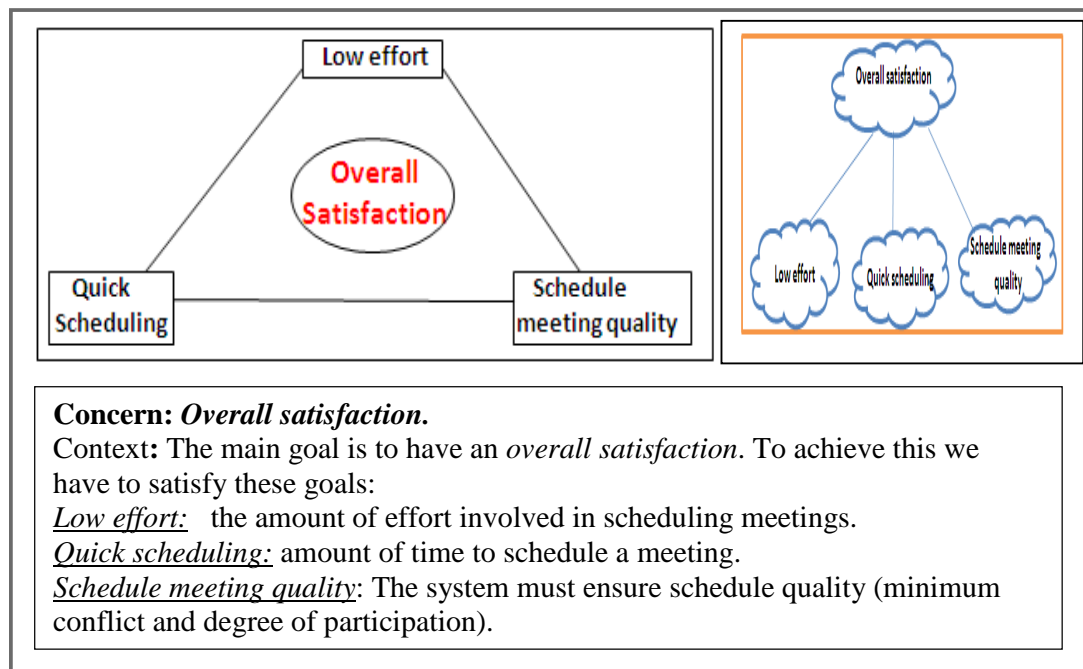
## Scenario [To be replaced with the above]

Suppose you want to arrange a memorable end-of-year party before the end of this academic year (end of May) that will gather most of your colleagues from Information Systems and Technology class 2010/2011.

# Task

Please do the following comparisons using the instructions for *Task 1.* Remember that the direction of all arrows should not create a cycle.

**Table 1:**



**Concern:** *Overall satisfaction.*
Context**:** The main goal is to have an *overall satisfaction*. To achieve this we have to satisfy these goals:
*Low effort:*   the amount of effort involved in scheduling meetings.
*Quick scheduling:* amount of time to schedule a meeting.
*Schedule meeting quality*: The system must ensure schedule quality (minimum conflict and degree of participation).

**Table 2:**

Choose schedule manually

Minimal conflict

Choose schedule automatically

Choose schedule collectively

Minimal conflict

Choose schedule manually

Choose schedule automatically

Choose schedule collectively

**Concern:** *Minimal conflict.*
Context: The main goal is to avoid all possible conflicts. To achieve this we have to satisfy these goals:
*Choose schedule manually*: The schedule will be selected manually.
*Choose schedule automatically:* The schedule will be selected automatically.
*Choose schedule collectively:* The schedule will be selected collectively.

**Table 3:**

Choose schedule manually

Degree of participation

Choose schedule automatically

Choose schedule collectively

Degree of participation

Choose schedule manually

Choose schedule automatically

Choose schedule collectively

**Concern:** *Degree of participation.*
Context: The main goal is to have a good number of participants. To achieve this we have to satisfy these goals:
*Choose schedule manually*: The schedule will be selected manually.
*Choose schedule automatically:* The schedule will be selected automatically.
*Choose schedule collectively:* The schedule will be selected collectively.

**Table 4:**



**Concern:** *Quick scheduling.*
Context: The main goal is to schedule a meeting in the shortest time possible. To achieve this we have to satisfy these goals:
*Choose schedule manually*: The schedule will be selected manually.
*Choose schedule automatically:* The schedule will be selected automatically.
*Choose schedule collectively:* The schedule will be selected collectively.


**Table 5:**



**Concern:** *Minimal conflict.*
Context: The main goal is to avoid all possible conflicts. To achieve this we have to satisfy these goals:
*Collect time tables by person by email:* Participants timetable information will be collected via email by a human (e.g. secretary).
*Collect time tables by person by all means:* A person will collect participants timetable information by using all possible means (telephone, regular mail,……..).
*Collect time tables by system*: Participants timetable information will be collected by a system.

**Table 6:**



**Concern:** *schedule meeting quality.*
Context: The main goal is to improve *schedule meeting quality*. To achieve this we have to satisfy these goals:
*Minimal conflict:* Minimize scheduling conflicts among participants.
*Degree of participation*: Number of participants that actually show-up.

**Table 7:**



**Concern:** *Low effort.*
Context: The main goal is to minimize the amount of effort involved in scheduling meetings. To achieve this we have to satisfy these goals:
*Collect time tables by person by email:* Participants timetable information will be collected via email by a human (e.g. secretary).
*Collect time tables by person by all means:* A person will collect participants timetable information by using all possible means (telephone, regular mail,……..).
*Collect time tables by system*: Participants timetable information will be collected by a system.

# C   Instrument B (Step 5)

Step 5 makes use of the results of the calculation of Step 4. The header reminds of the scenario. For each of the 5 comparisons (picked randomly from the previous questionnaire, the administrator picks a random line and fills in the result of the calculation (the "correct result"). She then generates all possible permutations thereof to complete the other rows (again randomly).

# Task 2/3 - SCENARIO [X]

## Scenario
[the scenario repeated]

# Task
Please select the priority profile that best matches your preferences for this scenario, as described in instructions for *Task 2.*

## Question 1:

**Quick scheduling**

| | | | Option 1 |
| | | | Option 2 |
| | | | Option 3 |
| | | | Option 4 |
| | | | Option 5 |
| | | | Option 6 |

Choose schedule manually

Choose schedule automatically

Choose schedule collectively

**Question 3:**

**Question 4:**



Low effort

Option 1

Option 2

Option 3

Option 4

Option 5

Option 6

Collect time tables by person by email

Collect time tables by person by all means

Collect time tables by system

# D   Instrument C (Step 6)

Step 6 also makes use of the results of the calculation of Step 4. The header reminds of the scenario. The administrator calculates the optimal ("correct") alternative and then picks another three as follows: two by changing one of the OR-decompositions of the correct alternative (so two partially "correct" alternatives) and one more by changing both OR-decompositions (yielding a totally "incorrect" alternative).

**Scenario**
[the scenario repeated]

**Task**
Please select your preferred solution for the scenario as per instructions for *Task 3*.

**[note: of the alternatives below, four (4) are chosen each time according to the description given in the main text]**
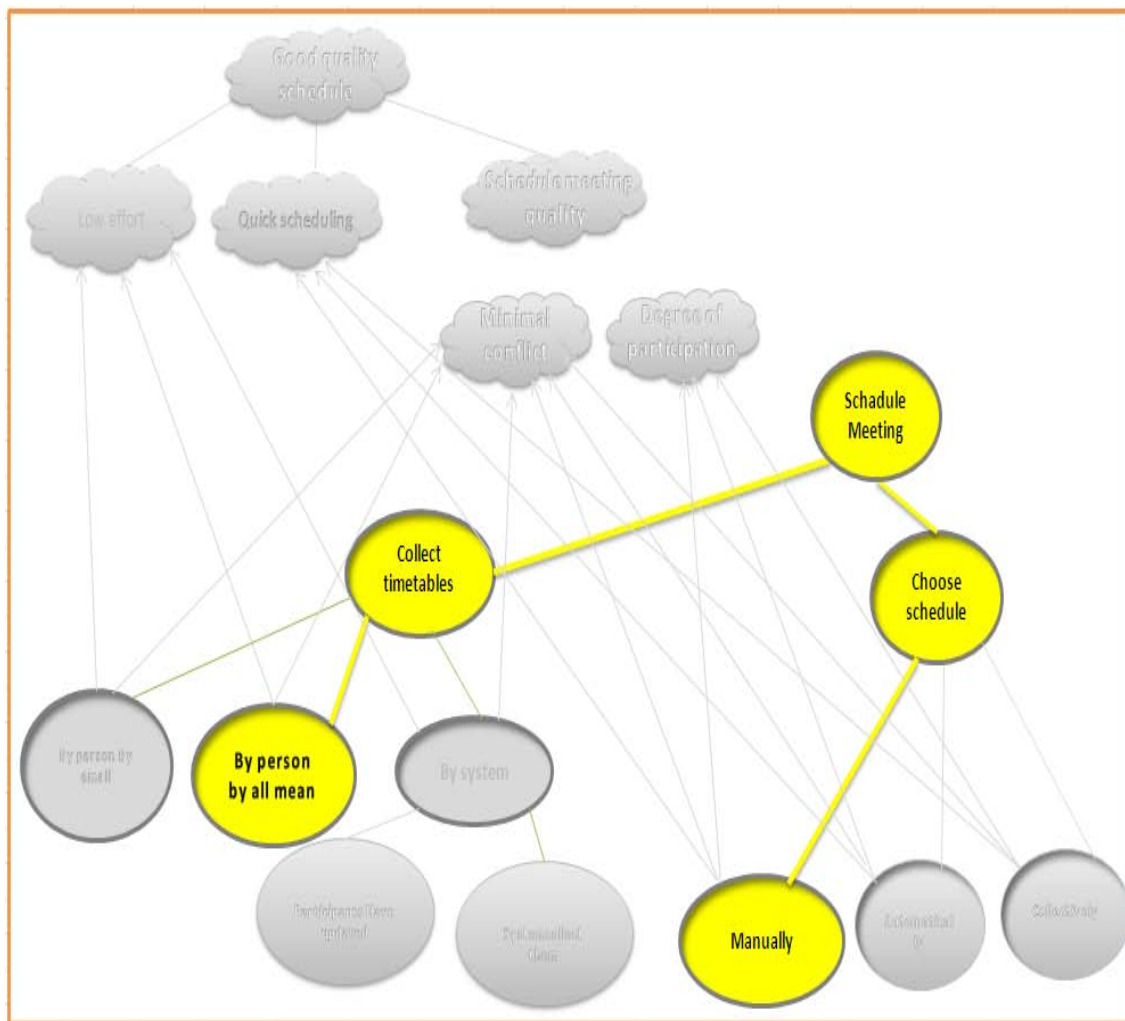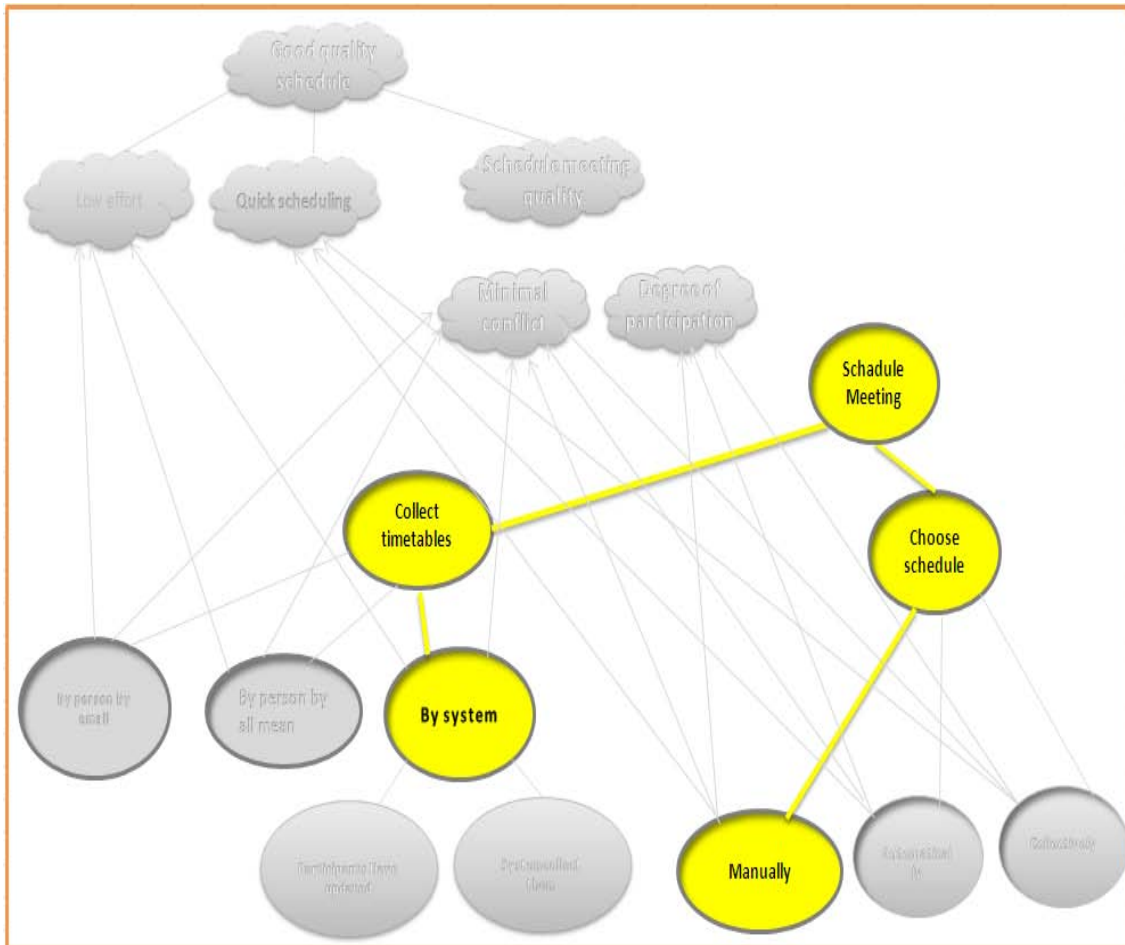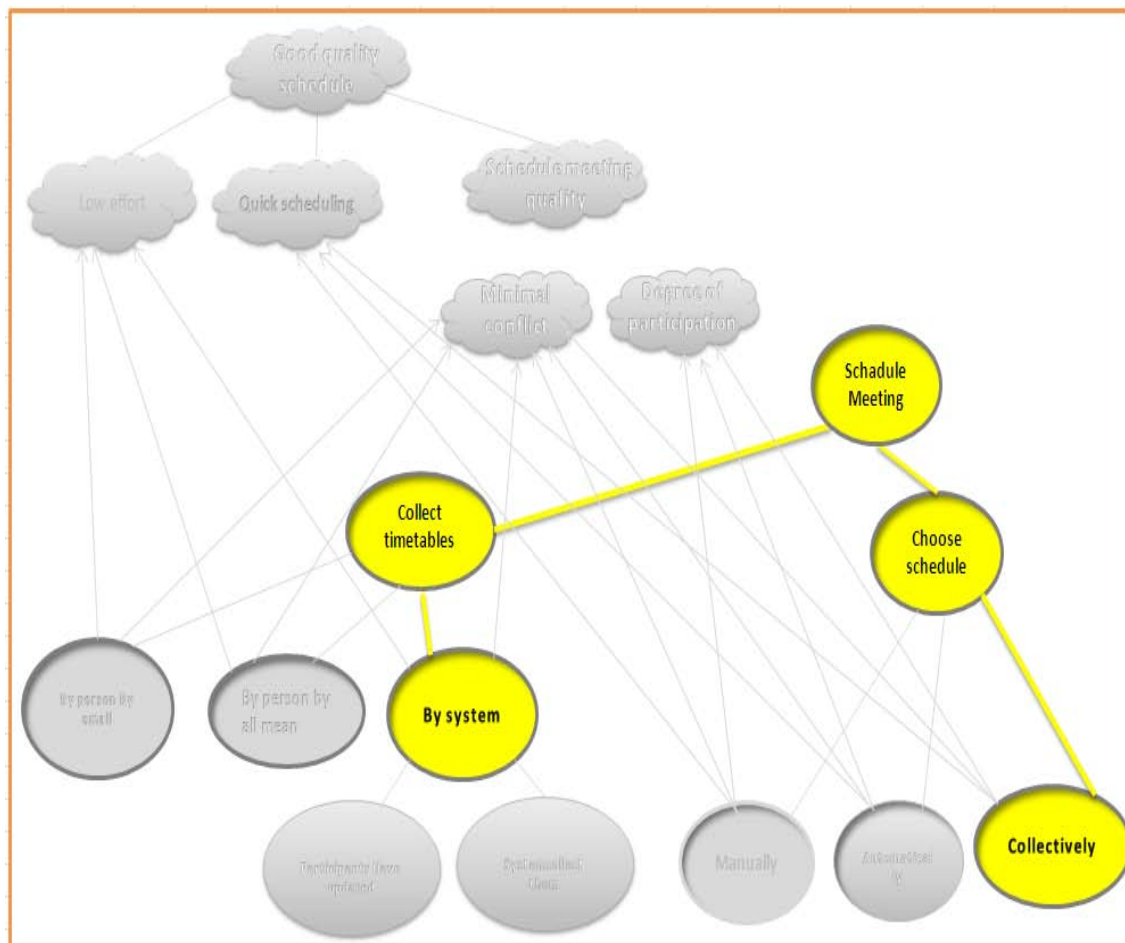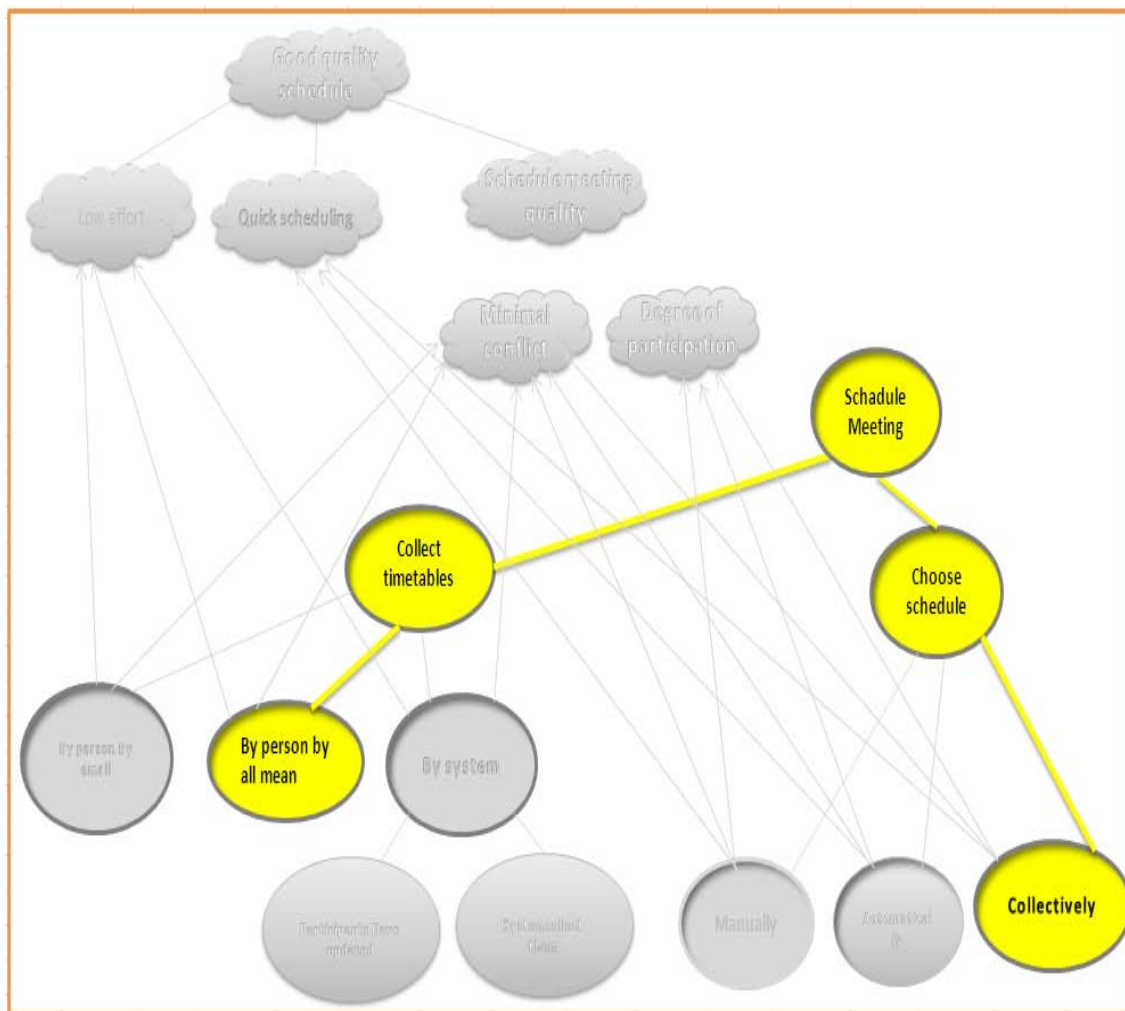
**Option :**

**Option :**

**Option :**

**Option :**

**Option :**

**Option :**

**Option :**

**Option :**

**Option :**