

# Forging High-Quality User Stories: Towards a Discipline for Agile Requirements

Garm Lucassen, Fabiano Dalpiaz, Jan Martijn E.M. van der Werf, and Sjaak Brinkkemper

Department of Information and Computing Sciences

Utrecht University, The Netherlands

{g.lucassen, f.dalpiaz, j.m.e.m.vanderwerf, s.brinkkemper}@uu.nl

**Abstract**—User stories are a widely used notation for formulating requirements in agile development projects. Despite their popularity in industry, little to no academic work is available on assessing their quality. The few existing approaches are too generic or employ highly qualitative metrics. We propose the Quality User Story Framework, consisting of 14 quality criteria that user story writers should strive to conform to. Additionally, we introduce the conceptual model of a user story, which we rely on to design the AQUASA software tool. AQUASA aids requirements engineers in turning raw user stories into higher-quality ones by exposing defects and deviations from good practice in user stories. We evaluate our work by applying the framework and a prototype implementation to three user story sets from industry.

**Index Terms**—User stories, requirements quality, AQUASA, QUS Framework, requirements engineering, NLP

## I. INTRODUCTION

As practitioners transition to agile development, requirements are increasingly expressed as user stories [1]. Invented by Connextra in the United Kingdom and popularized by Mike Cohn [2], user stories only capture the essential elements of a requirement: *who* it is for, *what* it expects from the system, and, optionally, *why* it is important (a key aspect in RE [3]). The most well-known format, popularized by Mike Cohn [2] is: “As a *<type of user>*, I want *<goal>*, [so that *<some reason>*]”. For example: “As an Administrator, I want to receive an email when a contact form is submitted, so that I can respond to it”.

In a 2014 survey among requirements analysts in an agile environment, user stories were the most used requirements documentation method [1]. Despite this popularity, the number of methods to assess and improve user story *quality* is limited. Existing approaches to user story quality employ highly qualitative metrics, such as the heuristics of the INVEST (Independent-Negotiable-Valuable-Estimable-Scalable-Testable) framework [4], and the generic guidelines for ensuring quality in agile RE proposed by Heck and Zaidman [5].

The goal of this paper is to introduce a comprehensive approach to assessing and enhancing user story quality. To achieve this goal, we take advantage of the potential offered by natural language processing (NLP) techniques, while taking into account the reservations of Daniel Berry and colleagues [6]. Existing state-of-the-art NLP tools for RE such as QuARS [7], Dowser [8], Poirot [9] and RAI [10] are unable to transcend from academia into practice. The ambitious

objectives of these tools necessitate a deep understanding of the requirements’ contents [6]. This necessity is currently unachievable and will remain impossible to achieve in the foreseeable future [11].

Instead, to be effective, tools that want to harness NLP should focus on the *clerical* part of RE that software can perform with 100% recall and high precision, leaving thinking-required work to human requirements engineers [6]. Additionally, they should conform to what practitioners actually do, instead of what the published methods and processes advise them to do [12]. User stories’ popularity among practitioners and their simple yet strict structure make them ideal candidates for applying NLP tools and techniques.

Throughout the remainder of this paper, we make five concrete contributions that pave the way for the creation of these types of tools:

- Sec. II formulates a revised notion of user story quality. The Quality User Story Framework separates the algorithmic aspects that NLP can automatically process from the thinking-required concerns which necessitate involving human requirements engineers. We illustrate each quality criteria with a real-world example to demonstrate that the quality defect exists in practice;
- Sec. III proposes the conceptual model of a user story, detailing the decomposition of a single user story. We use this conceptual model throughout the paper;
- Sec. IV explores the cross-story relationships within sets of user stories, enabling the identification of possible semantic quality improvements;
- Sec. V presents the design of a prototype tool (an example of a *dumb tool* [6]) that restricts itself to correcting the algorithmically determinable errors of user stories;
- Sec. VI empirically evaluates the feasibility of our approach by applying the framework and the prototype tool to multiple real-world user story sets.

We review related literature in Sec. VII, and conclude with discussion and future research directions in Sec. VIII.

## II. WHAT IS USER STORY QUALITY?

Many different perspectives on requirements quality exist. The IEEE Recommended Practice for Software Requirements Specifications defines eight quality characteristics for a requirement [13]: correct, unambiguous, complete, consistent, ranked for importance/stability, verifiable, modifiable

TABLE I  
QUALITY USER STORY FRAMEWORK

Criteria	Description
<b>Syntactic</b> - Atomic - Minimal - Well-formed	A user story expresses a requirement for exactly one feature A user story contains nothing more than role, means and ends A user story includes at least a role and a means
<b>Semantic</b> - Conflict-free - Conceptually sound - Problem-oriented - Unambiguous	A user story should not be inconsistent with any other user story The means expresses a feature and the ends expresses a rationale, not something else A user story only specifies the problem, not the solution to it A user story avoids terms or abstractions that may lead to multiple interpretations
<b>Pragmatic</b> - Complete - Explicit dependencies - Full sentence - Independent - Scalable - Uniform - Unique	Implementing a set of user stories creates a feature-complete application, no steps are missing Link all unavoidable, non-obvious dependencies on user stories A user story is a well-formed full sentence The user story is self-contained, avoiding inherent dependencies on other user stories User stories do not denote too coarse-grained requirements that are difficult to plan and prioritize All user stories follow roughly the same template Every user story is unique, duplicates are avoided

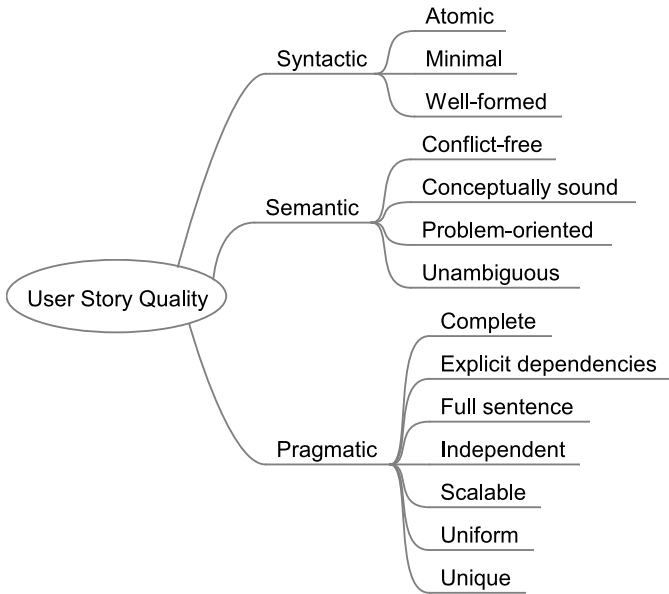


Fig. 1. Quality User Story Framework

and traceable. However, most requirements specifications are unable to adhere to these in practice [14]. On top of this, these quality characteristics were not developed with user stories nor agile development in mind.

The Agile Requirements Verification Framework [5] defines three high-level verification criteria for requirements for agile development: completeness, uniformity, and consistency & correctness. The framework proposes specific criteria to be able to apply the quality framework to both feature requests and user stories. Many of these criteria, however, require supplementary, unstructured information that is not captured in the primary user story text, making them inadequate for dumb RE tools and our perspective.

With this in mind, we take inspiration from the verification framework [5] to define a new Quality User Story (QUS)

Framework (Figure 1 and Table I). The QUS Framework only focuses on the information that is derivable from user story texts themselves, disregarding all requirements management concerns such as effort estimation and additional information sources such as descriptions or comments. The QUS Framework comprises 14 criteria that influence the quality of a user story or set of user stories. Because user stories are entirely textual, we classify each quality criteria according to three concepts borrowed from linguistics, similar to Lindland [15]:

**Syntactic** quality, concerning the textual structure of a user story without considering its meaning;

**Semantic** quality, concerning the relations and meaning of (parts of) the user story text;

**Pragmatic** quality, regarding choosing the most effective alternatives for communicating a given set of requirements.

In the next subsections, we introduce each criterion by presenting: (1) a comprehensive explanation of the criterion, (2) an example user story that violates the specific criterion, and (3) why the example violates the specific criterion. The example user stories originate from two real-world user story databases of software companies in the Netherlands. One contains 98 stories that specify the development of a tailor-made web information system. The other consists of 26 user stories from an advanced healthcare software product for home care professionals. We refrain from disclosing additional application details due to confidentiality constraints.

#### A. Syntax

1) *Atomic*: A user story should concern only one feature. It is tempting to combine multiple features in one user story when they are related or similar. Not doing this, however, improves the accuracy of the expected effort estimation. According to practitioners, the combined effort estimation of two small, clear-cut user stories is more accurate than the estimation of one larger, more opaque user story [16]. The user story US<sub>1</sub> in Table II consists of two separate requirements: the act of clicking on a location, and the display of associated

TABLE II  
USER STORIES THAT BREACH QUALITY CRITERIA FROM TWO REAL-WORLD CASES

ID	Description	Issues
US <sub>1</sub>	As a User, I'm able to click a particular location from the map and thereby perform a search of landmarks associated with that latitude longitude combination	Not atomic: two stories in one
US <sub>2</sub>	As a care professional I want to see the registered hours of this week (split into products and activities). See: Mockup from Alice NOTE: - First create the overview screen - Then add validations	Not minimal, due to additional note about the mockup
US <sub>3</sub>	Add static pages controller to application and define static pages	Missing role
US <sub>4</sub>	As a User, I want to open the interactive map, so that I can see the location of landmarks	Conceptual issue: the end is in fact a reference to another story
US <sub>5</sub> US <sub>6</sub>	As a User, I'm able to edit any landmark As a User, I'm able to delete a landmark which I added	Conflict: US <sub>5</sub> refers to any landmark, while US <sub>6</sub> only to those that user has added
US <sub>7</sub>	As a care professional I want to save a reimbursement. - Add save button on top right (never grayed out)	Hints at the solution
US <sub>8</sub>	As a User, I am able to edit the content that I added to a person's profile page	Unclear: what is content here?
US <sub>9</sub>	As an Administrator, I am able to view content that needs to be reviewed	The type of content is not specified
US <sub>10</sub>	Server configuration	In addition to being syntactically incorrect, this is not even a full sentence
US <sub>11</sub> US <sub>12</sub>	As an Administrator, I am able to add a new person to the database followed by As a Visitor, I am able to view a person's profile	Viewing relies on first adding a person to the database
US <sub>13</sub>	As a care professional I want to see my route list for next/future days, so that I can prepare myself (for example I can see at what time I should start traveling)	Difficult to estimate because it is unclear what see my route list implies
US <sub>14</sub>	As an Administrator, I receive an email notification when a new user is registered	Deviates from the template, no "wish" in the means
EP <sub>A</sub> US <sub>15</sub>	As a Visitor, I'm able to see a list of news items, so that I can stay up to date on news As a Visitor, I'm able to see a list of news items, so that I can stay up to date on news	The same requirement is repeated both in epic EP <sub>A</sub> , and in a user story US <sub>14</sub>

landmarks. The requirements engineer should split this user story into two autonomous user stories.

2) *Minimal*: User stories should contain a role, means and, optionally, ends. Any additional information such as comments, descriptions of the expected behavior or testing hints are to be captured as additional notes. Take, for example, US<sub>2</sub>; aside from a role and means, this user story includes a reference to an undefined mockup and a note on how to approach the implementation. The requirements engineer should move both to a comments or description section.

3) *Well-formed*: Before it can be considered a user story, the core requirements text needs to define a role and the expected functionality: the *means*. US<sub>3</sub> is one of the eight user stories in our sample user story database that do not adhere to this basic syntax, forgoing the role. This is most likely done because the requirement is quite technical, defining static pages and database issues. Nevertheless, the requirements engineer should fix this issue by introducing the relevant role.

## B. Semantic

1) *Conceptually Sound*: The means and ends parts of a user story play a specific role. The means should capture a concrete feature, while the ends express the rationale for that feature. Consider US<sub>4</sub>: the end is actually a dependency on another (hidden) functionality, which is required in order for the means to be realized, implying the existence of a landmark database which is not mentioned in any of the other stories. A significant additional feature that is mistakenly represented as an end, but should be a means in a separate user story.

2) *Conflict-free*: To prevent implementation errors and rework, a user story should not conflict with any of the other user stories in the database. Requirements conflict occur when

two or more requirements cause an inconsistency [17], [18]. Although a comprehensive taxonomy of all types of conflicts is beyond the scope of this work, two major families of conflicts concern *activities* or *resources* [19]. Story US<sub>6</sub> contradicts the requirement that a user can edit any landmark (US<sub>5</sub>). The new information that users are only allowed to delete content that they added themselves, raises the question whether this constraint also counts for the first requirement. The requirements engineer should modify the first user story and explicitly include whether a user can modify all landmarks.

3) *Problem-oriented*: A user story should only specify the problem, not the solution. If absolutely necessary, include implementation hints as a comment or description of the user story. Aside from breaking the minimal quality criteria, US<sub>7</sub> includes an implementation specification within the core user story text. The requirements engineer should remove this section and, if essential, post it as a comment.

4) *Unambiguous*: Ambiguity is inherent to natural language requirements, but the requirements engineer writing user stories should avoid it as much as possible. Not only should a user story be internally unambiguous, but it should also be clear in relationship to all other user stories. The Taxonomy of Ambiguity Types [20] is a comprehensive overview of the kinds of ambiguity that can be encountered in a systematic requirements specification. For user stories, we explore some possibilities in Sec. IV. Examples include registering alternative means with the same purpose or linking superclass terms to its respective elements. In US<sub>8</sub>, "content" is a superclass referring to audio, video and textual media uploaded to the profile page. The requirements engineer should explicitly mention which media are editable.

### C. Pragmatic

1) *Complete*: Implementing a set of user stories should lead to a feature-complete application. While user stories should not thrive to cover 100% of the application’s functionality preemptively, the requirements engineer must take care not to forget or neglect crucial user stories that cause a show stopping feature-gap. Unfortunately, it is impossible to include a missing user story in a table, but as an example consider that to be able to delete an item you first need to create it.

2) *Explicit Dependencies*: Whenever a user story has a non-obvious dependency, it should explicitly link to the user story tag of the user story it depends on. For example, US<sub>9</sub> does not explicate what types of text or media “content” references. To fix this, the requirement engineer should add the user story tags of the user stories that capture creating the relevant types of consistent.

3) *Full Sentence*: A user story should read like a full sentence, without typos or grammatical errors. US<sub>10</sub>, for example, is not expressed as a full sentence (in addition to not complying with syntactic quality). By reformulating the feature as a full sentence user story, it will automatically specify what exactly needs to be configured.

4) *Independent*: User stories should not overlap in concept and should be schedulable and implementable in any order [4]. Note that this quality criterion is more of a ground rule that engineers try to follow to the best of their possibilities. Much like in programming loosely coupled systems, it is impossible to never breach this quality criterion. For example, US<sub>12</sub> is dependent on US<sub>11</sub>, because it is impossible to view a person’s profile without first laying the foundation for creating a person. Although it is contradictory that a quality criteria is unattainable, requirements independence is sufficiently important to warrant inclusion in the QUS Framework.

5) *Scalable*: As user stories grow in size, it becomes more difficult to accurately estimate the effort required for the entire set of user stories. Therefore, each user story should not become so large as to avoid their estimation and planning with reasonable certainty [4]. For example, US<sub>13</sub> requests a route list so that care professionals can prepare themselves. While this might be just an unordered list of places to go to during a workday, it is likely that the feature includes ordering the routes algorithmically to minimize distance traveled and/or showing the route on a map. These many functionalities makes accurate estimation difficult. The requirements engineer should split the user story into multiple, more specific user stories so that effort estimation is more accurate.

6) *Uniform*: All user stories should follow the same, agreed upon template. Minimal deviations are allowed when this better suits the narrative structure of the user story. For instance, using “I am able” instead of “I want to” is an acceptable deviation. However, US<sub>14</sub> is a bigger deviation, along with 9 other stories in the first user story set (10.9% of the total). For uniformity, the requirements engineer should redefine these user stories so that they express a wish.

7) *Unique*: A user story should not be a duplicate of another user story nor epic (large user story). For example,

both epic EP<sub>A</sub> and user story US<sub>15</sub> (expressed as part of epic EP<sub>A</sub>) share the same text. In this case, the epic contains seven other user stories that also concern other artifacts than news items. To resolve this issue, an epic can be formulated that captures the importance of both news items and events.

### III. A CONCEPTUAL MODEL OF USER STORIES

There are more than 80 syntactic variants of user stories [21]. Although originally an unstructured written description similar to use cases [22] but restricted in size [2], nowadays user stories follow a strict, compact template that captures *who* it is for, *what* it expects from the system, and (optionally) *why* it is important in a simple manner [21].

When used in Scrum, two other artifacts are relevant: epics and themes. An epic is a label for a large user story, which is broken down into smaller, implementable user stories. A theme is a collection of user stories grouped according to a given criterion [2]. For simplicity, and due to their greater popularity, we only include epics in our conceptual model.

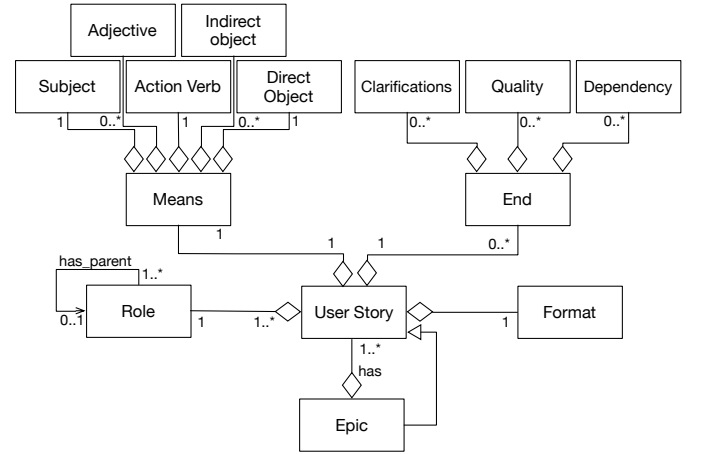


Fig. 2. Conceptual model of user stories

Figure 2 shows an UML class diagram for user stories. A user story itself consists of four parts: one role, one means, optionally one or more ends and a format. In the following subsections, we elaborate on how to decompose each of these. Note that we deviate from Cohn’s terminology as presented in the introduction, using the more abstract means-end instead of goal-reason. We do so because alternative user story templates talk of desire instead of goal, or benefit as opposed to reason, and we do not want to position Cohn’s terminology as better than the others.

#### A. Format

A user story should follow some pre-defined template as agreed upon among the stakeholders, such as the one proposed in [2]. The pre-defined text or skeleton of the template is what we depict as *format* in the conceptual model. This format is the foundation, in between which the role, means and optional end(s) are interspersed to forge a user story.

## B. Role

A user story always defines one relevant role. Roles hierarchically relate to each other through the *has\_parent* relationship. This can be used, for example, to determine that “*Editor*” is a specific type of (has parent) “*User*”.

## C. Means

Means can have different structures, for they can be used to represent different types of requirements. However, grammatically speaking, all means have three things in common: (1) they contain a *subject* with an intent such as “*want*” or “*am able*”, (2) followed by an *action verb*<sup>1</sup> that expresses the action related to the feature being requested, and (3) a *direct object* on which the subject executes the action. For example: “*I want to open the interactive map*”. Aside from this basic requirement, means are essentially free form text which allow for an infinite number of constructions. Two common additions are an adjective or an indirect object, which results in the following example: “*I want to open a larger (adjective) view of the interactive map from the person’s profile page (indirect object)*”. We included these interesting cases in the conceptual model, but left out all other variations, which will be studied in future research. Moreover, we find them appropriate for inclusion in domain-specific conceptual models such as those for Web Applications Requirements Engineering [23] or Embedded Real-Time Software [24].

## D. End

The end part provides the reason for the means [2]. However, user stories frequently end up including other information. By analyzing the ends available in our learning set of user stories, we define three possible variants of a well-formed end:

1. **Clarification of means.** The end explains the reason of the means. Example: “*As a User, I want to edit a record, so that I can correct any mistakes*”.
2. **Dependency on another functionality.** The end (implicitly) references a functionality which is required for the means to be realized. Although dependency is an indicator of a bad quality criteria, having no dependency at all between requirements is impossible [4]. There is no size limit to this dependency on the (hidden) functionality. Small example: “*As a Visitor, I want to view the homepage, so that I can learn about the project*”. The end implies the homepage also has relevant content, which requires extra input. Larger example: “*As a User, I want to open the interactive map, so that I can see the location of landmarks*”. The end implies the existence of a landmark database, a significant additional functionality of the same requirement.
3. **Qualitative requirement.** The end communicates the intended qualitative effect of the means. For example: “*As a User, I want to sort the results, so that I can more easily*

<sup>1</sup>While other types of verbs are in principle admitted, in this paper we focus on action verbs, which are the most used in user stories requesting features

*review the results*” indicates that the means contributes maximizing easiness.

Note that these three are not mutually exclusive, but can occur simultaneously such as in “*As a User, I want to open the landmark, so that I can more easily view the landmark’s location*”. The means only specifies that the user wishes to view a landmark’s page. The end, however, contains elements of all three types: (1) a clarification that you want to open the landmark to view its location, (2) additional functionality of the landmark, and (3) the qualitative requirement that it should be easier than an alternative.

## IV. IDENTIFYING CROSS-STORY RELATIONSHIPS

There are two dimensions of user story quality: individual and cross-story. While several criteria in Table I concern individual user stories, some require taking into account all other user stories for the project at hand. For example, whether a single user story is complete, independent, uniform and unique depends on the entire set of user stories.

In this section, we explore the relationships *between* user stories. By characterize each relationship, formalize them via first-order logic predicates to enable their concrete identification, and associate the relevant quality criteria. Our aim is to provide an initial set of relevant relationships, which is far from being complete.

*Notation.* Lowercase identifiers refer to single elements (e.g., one user story), and uppercase identifiers denote sets (e.g., a set of user stories). We use the following notation:

- $U$  for the set of user stories (in a given project);
- $r_1, r_2, \dots$  for role identifiers;
- $m_1, m_2, \dots$  for means identifiers, where  $m = \langle s, av, do, io, adj \rangle$  with  $s$  being a subject,  $av$  an action verb,  $do$  a direct object,  $io$  an indirect object, and  $adj$  an adjective ( $do$  and  $io$  may be null, see Fig. 2);
- $e_1, e_2, \dots$  for end identifiers;
- $E_1, E_2, \dots$  for identifiers of sets of ends;
- $f_1, f_2, \dots$  for user story format;
- $\mu_1, \mu_2, \dots$  for user stories, where  $\mu = \langle r, m, E, f \rangle$ , or, expanding  $m$ ,  $\mu = \langle r, \langle s, av, do, io, adj \rangle, E, f \rangle$

Furthermore, we assume that the equality, intersection, etc. operators are semantic, i.e., they look at the meaning of an entity (e.g., they account for synonyms, etc.), in addition to the syntax. To denote that an operator is merely syntactic, we add the “*syn*” subscript; so, for instance,  $=_{syn}$  is syntactic equivalence. In this section, we employ theoretical semantic operators; Sec. V will discuss the feasibility of some of them. Finally, the function  $depends(av, av')$  denotes that executing the action  $av$  on a specific object requires first executing  $av'$  on that very object (e.g., “*delete*” depends on “*create*”).

### A. Complete

Some user stories imply the necessity of other functionality not yet captured in any of the user stories written so far. A simple example is user stories with action verbs that refer to a non-existent direct object: to read, update or delete an item

you first need to create it. In practice, determining completeness is a context-dependent issue that is hard to generalize. Thus, we define this relationship on a high abstraction level, but focusing on dependencies concerning the means' direct object. Formally, the predicate  $missesDep(\mu)$  holds when a dependency for  $\mu$ 's direct object is missing:

$$missesDep(\mu) \leftrightarrow depends(av, av') \wedge \nexists \mu' \in U. do' = do$$

### B. Independent

Many different types of dependency exist, and our intent is not to list them here. For illustration, we explain two cases.

1) *Causality*: In some cases, it is necessary that one user story  $\mu_1$  is completed before the developer can start on another user story  $\mu_2$  (US<sub>11</sub> and US<sub>12</sub> in Table II). Such a causal dependency foremost impacts the *independent* quality criterion. When this dependency is non-obvious and implicit, the requirements engineer needs to consider adding an *explicit dependency*. Formally, the predicate  $hasDep(\mu_1, \mu_2)$  holds when  $\mu_1$  causally depends on  $\mu_2$ :

$$hasDep(\mu_1, \mu_2) \leftrightarrow depends(av_1, av_2) \wedge do_1 = do_2$$

2) *Superclasses*: An object of one user story  $\mu_1$  can refer to multiple other objects of a set of user stories  $\{\mu_2, \mu_3, \dots\}$ , indicating that the object of  $\mu_1$  is a parent or *superclass* of the other objects. "Content" for example can refer to different types of multimedia, as exemplified in US<sub>8</sub>. This relationship has an impact on the *unambiguous*, *independent* and *explicit dependencies* quality criteria. One can detect ambiguity or missing dependencies by scanning for Create, Read, Update and Delete (CRUD) actions which have a unique direct object. Formally, predicate  $hasIsaDep(\mu, do')$  is true when  $\mu$  has a direct object superclass dependency based on the sub-class  $do'$  of  $do$ .

$$hasIsaDep(\mu, do') \leftrightarrow \exists \mu' \in U. is-a(do', do)$$

### C. Uniformity

Uniformity in the context of user stories means that a user story has a format that is consistent with the format of all other user stories. To test this, the requirements engineer needs to determine the most frequently occurring format, typically the format agreed upon with the team. The format of an individual user story  $\mu = \langle r, m, E, f \rangle$  is syntactically compared to the most common format  $f_{std}$  to determine whether it adheres with the *uniformity* quality criterion. US<sub>14</sub> in Table II was an example of a non-uniform user story. Formally, predicate  $isNotUniform(\mu, f_{std})$  is true if the format of  $\mu$  deviates from the standard:

$$isNotUniform(\mu, f_{std}) \leftrightarrow f \neq_{syn} f_{std}$$

### D. Unique

A user story is unique when no other user story is (semantically) the same or too similar. There are many different ways in which two user stories can be similar. For example, all user stories should follow a similar user story format. The type of similarity we are interested in, however, is a potential

indicator of duplicate or conflicting user stories such as user stories US<sub>5</sub> and US<sub>6</sub>, or user story US<sub>15</sub> and epic EP<sub>A</sub> in Table II. We discuss five similarity relationships that help detecting similarity among user stories.

To detect these types of relationships, each user story part needs to be compared with the parts of other user stories, using a combination of similarity measures that are either syntactic (e.g., Levenshtein's distance) and semantic (e.g., employing an ontology to determine synonyms). When similarity exceeds a certain threshold, a human analyst is required to examine the user stories for potential conflict and/or duplication.

1) *Full Duplicate*: A user story  $\mu_1$  is an exact duplicate of another user story  $\mu_2$  when the stories are identical. This impacts the *unique* quality criterion. Formally,

$$isFullDuplicate(\mu_1, \mu_2) \leftrightarrow \mu_1 =_{syn} \mu_2$$

2) *Semantic Duplicate*: A user story  $\mu_1$  that duplicates the request of  $\mu_2$ , while using a different text; this has an impact on the *unique* quality criterion. Formally,

$$isSemDuplicate(\mu_1, \mu_2) \leftrightarrow \mu_1 = \mu_2 \wedge \mu_1 \neq_{syn} \mu_2$$

3) *Different Means, Same End*: Two or more user stories that have the same end, but achieve this using different means. This relationship potentially impacts two quality criteria, as it may indicate: (i) a feature variation that should be explicitly noted in the user story to maintain an *unambiguous* set of user stories, or (ii) a conflict in how to achieve this end, meaning one of the user stories should be dropped to ensure *conflict-free* user stories. Formally, for user stories  $\mu_1$  and  $\mu_2$ :

$$diffMeansSameEnd(\mu_1, \mu_2) \leftrightarrow m_1 = m_2 \wedge E_1 \cap E_2 \neq \emptyset$$

4) *Same Means, Different End*: Two or more user stories that use the same means to reach different ends. This relationship affects the qualities of user stories to be *unique* or *independent* of each other. If the ends are not conflicting, they could be combined into a single larger user story; otherwise, they are multiple viewpoints that should be resolved. Formally,

$$sameMeansDiffEnd(\mu_1, \mu_2) \leftrightarrow m_1 \neq m_2 \wedge (E_1 \setminus E_2 \neq \emptyset \vee E_2 \setminus E_1 \neq \emptyset)$$

5) *Different Role, Same Means and/or Same End*: Two or more user stories with different roles, but same means and/or ends indicates a strong relationship. Although this relationship has an impact on the *unique* and *independent* quality criteria, it is considered good practice to have separate user stories for the same functionality for different roles. As such, requirements engineers could choose to ignore this impact. Formally,

$$diffRoleSameStory(\mu_1, \mu_2) \leftrightarrow r_1 \neq r_2 \wedge (m_1 = m_2 \vee E_1 \cap E_2 \neq \emptyset)$$

6) *Purpose = Means*: The end of one user story  $\mu_1$  is identical to the means of another user story  $\mu_2$ . Indeed, the same piece of text can be used to express both a wish, and a reason for another wish. When there is this strong a semantic relationship between two user stories, it is important

to add *explicit dependencies* to the user stories, which is a trade-off with the *independent* quality criterion. Formally,  $purposeMeans(\mu_1, \mu_2, x)$  is true if  $x$  is an end in  $\mu_1$  and a means in  $\mu_2$

$$purposeMeans(\mu_1, \mu_2, x) \leftrightarrow \exists e \in E_1 \text{ s.t. } e = m_2$$

## V. A DUMB TOOL FOR IMPROVING USER STORY QUALITY

The Quality User Story (QUS) Framework provides structured guidelines for improving the quality of (a set of) user stories. To support the framework, we propose the Automatic Quality User Story Artisan (AQUSA) tool, which exposes defects and deviations from good practice in user stories.

Unlike most NLP tools for RE, and in line with Berry's notion of a *dumb tool* [6], we require our tool to detect defects with close to 100% recall. Not fulfilling this requirement means that a human requirements engineer has to double check the entire requirements document for missed defects, which we want to avoid [25]. On the other hand, *precision*, the number of false positives in proportion to the detected defects, should not be so high that the user perceives AQUSA to report useless errors. Thus, AQUSA is designed as a tool that focuses on easily describable, algorithmically determinable defects. This allows the requirements engineer to focus on thinking-required defects for which 100% recall with high precision is impossible [25]. Consequently, AQUSA supports only certain QUS criteria:

- Syntactical criteria are detectable with 100% recall; AQUSA can report need-to-improve defects with high precision.
- Semantic criteria are impossible to detect with 100% recall and are thus out of scope for AQUSA.
- Some aspects of the pragmatic criteria are detectable with 100% recall, while others are not. For AQUSA, we select a number of algorithmically determinable subparts.

Next, we present the selected quality criteria, discuss their theoretical implementation and provide accompanying example input and output user stories. This is followed by AQUSA's architecture and implementation status.

### A. Syntax

1) *Well-formed*: One of the essential aspects of verifying whether a string of text is a requirement, is splitting it into role, means and end(s). This process consists of two steps: (1) *chunking* on commas and common indicator texts such as *As a*, *I want to*, *I am able to* and *so that*; (2) verifying that each chunk contains their relevant part. A grammatical tagger assigns a word category to each of the words in the chunk. For each chunk, AQUSA tests the following rules:

**Role:** Is the last word a noun? Do the words before the noun match a known role format?

**Means:** Is the first word I? Can we identify a known means format? Does the part include two verbs and a noun?

**End:** Is the end present? Does it start with a known end format?

When AQUSA encounters the user story "*Add static pages controller to application and define static pages*", it includes the user story in its report, highlighting that it is not well-formed because it does not explicitly contain a role nor means. The well-formed user story "*As a Visitor, I want to register at the site, so that I can contribute*", however, would be verified and separated into the following chunks for later use:

**Role:** As a Visitor

**Means:** I want to register at the site

**End:** so that I can contribute

2) *Atomic*: To audit that the means of the user story concerns only one feature, AQUSA parses the means for occurrences of "*and*", "*&*", "*+*" to include any double feature requests in its report. Additionally, AQUSA suggests the reader to split the user story into multiple user stories.  $US_1$  from Table II, would generate a suggestion to be split into two user stories: (1) "*As a User, I want to click a location from the map*" and (2) "*As a User, I want to search landmarks associated with the lat long combination of a location*".

3) *Minimal*: To test this quality criterion, AQUSA relies on the results of chunking and verification of the *role and means* quality criterion. When this process has been successfully completed, AQUSA reports any user story that contains additional text after a dot, hyphen, semicolon or other separating punctuation marks. For example the "*See: Mockup from Alice*" text from  $US_2$  is sufficient for AQUSA to report that the user story is not minimal.

### B. Pragmatic

1) *Explicit Dependencies*: Whenever a user story includes an explicit dependency on another user story, it should include a navigable link to the dependency. Because the popular issue trackers Jira and Pivotal Tracker use numbers for dependencies, AQUSA checks for numbers in user stories and checks whether the number is contained within a link. The example "*As a care professional, I want to edit the planned task I selected - see 908*." would prompt the user to change the isolated number to "*See PAW-908*". On the end of the issue tracker, this should automatically change to "*see PAW-908*" (<http://company.issue tracker.org/PAW-908>)

2) *Uniform*: Aside from chunking, AQUSA extracts the user story format parts out of each chunk and counts their occurrences throughout the set of user stories. The most commonly occurring format is used as the standard user story format. All other user stories are marked as non-compliant to the standard and included in the error report. For example, AQUSA reports that "*As a User, I am able to delete a landmark*" deviates from the standard "*I want to*".

3) *Unique*: AQUSA implements each of the similarity measures that we outlined in Sec. IV using the WordNet lexical database [26] to detect semantic similarity. For each verb and object in a means or end, AQUSA runs a WordNet::Similarity calculation with the verbs or objects of all other means or ends. Combining the calculations results in one similarity degree for two user stories. When this metric is bigger than 90%, AQUSA reports the user stories as potential duplicates.

### C. Architecture and Implementation

AQUSA is designed as a simple, stand-alone, deployable as-a-service application that analyzes a set of user stories regardless of its source of origin. AQUSA exposes an API for importing user stories, meaning that AQUSA can easily integrate with any requirements management tool including MS Excel spreadsheets by developing middleware. By retaining its independence from other tools, AQUSA is capable of easily adapting to future technology changes. Aside from importing user stories, AQUSA consists of five main architectural components (Figure 3):

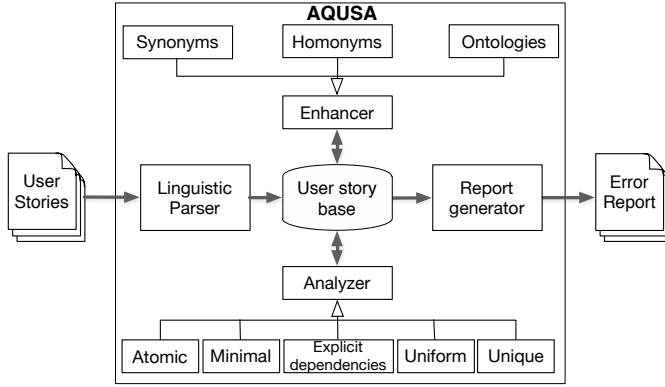


Fig. 3. Architecture of AQUSA

**Linguistic parser:** the first step for every user story is validating that it is well-formed. This takes place in the linguistic parser, which chunks the text according to common indicator texts and assigns word categories based on the NLTK<sup>2</sup> grammatical tagger.

**User story base:** a parsed user story is captured as an object—aligned with the conceptual model in Fig. 2—in the user story base, ready to be further processed.

**Enhancer:** AQUSA enhances user stories by adding possible synonyms, homonyms and relevant semantic information—extracted from an ontology—to the relevant words in each chunk.

**Analyzer:** AQUSA analyzes user stories by running methods that verify the selected syntactic and pragmatic quality criteria: atomic, minimal, explicit dependencies, uniform and unique.

**Report generator:** AQUSA captures the results of the parser, enhancer and analyzer in a comprehensive report.

We developed an AQUSA proof of concept with a limited scope that is still under active development<sup>3</sup>. Although currently only the syntactical quality criteria are implemented that has only been tested with a single user story set, these preliminary results are promising. Out of 96 user stories, this proof of concepts finds 7 that are not well-formed and 8 that are not atomic. In this circumstance, AQUSA led to 100% recall of the manual analysis that will be presented in Sec. VI.

<sup>2</sup>Natural Language Toolkit: <http://www.nltk.org/>

<sup>3</sup>Code available at <http://github.com/gglucass/aqusa>

## VI. EMPIRICAL EVALUATION

We evaluate the QUS Framework in two ways. In Sec. VI-A, we manually apply all the quality criteria to a third set of mature user stories to validate their occurrence besides the two user story sets we used while constructing our approach. Additionally, in Sec. VI-B, we apply the algorithmically determinable quality criteria of AQUSA to all three user story sets to show which errors exist in real-world user story requirements, and to evaluate the effectiveness of our identification mechanisms.

### A. Applying the Quality User Stories Framework

We validate the QUS Framework by applying its criteria to a third set of *intuitively high-quality* user stories, with the aim to assess if its perceived quality is confirmed by our quality framework. These user stories belong to a point of sales software product with customers across the world. A requirements engineering team in Belgium creates the user stories and outsources their development to a near-shore team in Romania. Over the span of a year, both sides made significant investments to be able to collaborate in an effective manner. User stories are central to this collaboration. In the experience of the Belgian team, user story quality directly influences the quality of the software product. Despite their high quality, however, these user stories still contain flaws. In the following paragraphs, we discuss how and why some criteria are breached, emphasizing three interesting patterns that arose.

1) *The Pragmatic Requirements Engineer:* The data set contains a number of quality errors that are easily avoidable. They request more than one feature, do not comply to the user story format or specify a solution instead of a problem. Aside from a few exceptions, these errors occur because they are technical requirements that are difficult to capture from the end-user's perspective.

2) *Keep it Minimal:* It is particularly difficult to adhere to one quality criteria: minimal. More than 10% of the user stories include additional information, distracting the reader from the user story itself.

3) *Context Is King:* The examined user stories concern a mature software product for a specific business domain with its own particular jargon. As a consequence, it is difficult for a non-informed outsider to accurately test some quality criteria. Not knowing the technical implications of user stories makes it difficult to estimate their scalability, an effect that jargon amplifies by causing lexical ambiguity. Moreover, without intimate knowledge of the application, it is nearly impossible to detect dependencies or conflicts among user stories of which one cannot be sure that they are conceptually sound. This clearly calls for reliance on domain-specific ontologies.

The overall quality is high: no dependencies or conflicts were immediately obvious. Nevertheless, structured approaches can detect some defects, e.g., by collecting all user stories that contain a frequently occurring phrase we did find some dependencies. For instance: by grouping nine user stories



concerning “*price*”, we can identify a causal dependency that was hidden among the other user stories.

### B. Applying the Automatic QUS Artisan

To test whether a full AQUUSA implementation is effective to determine user story quality, we manually apply the rules of Sec. V to our three user story sets. The resulting quality criteria breaches in Table III show promising results that indicate high potential for successful further development. For each set, at least 25% of the processed user stories violate one or more quality criteria that the AQUUSA algorithms can detect. AQUUSA detected the total amount of user stories with errors with 71% precision. Furthermore, the quality criteria breaches significantly vary per user story set. User story set #1 is very minimal, but breaks uniformity in 15% of all user stories. For user story sets #2 and #3 the inverse is true; although they are very uniform, they are not minimal.

TABLE III  
NUMBER OF IDENTIFIED VIOLATIONS (V) AND FALSE POSITIVES (FP) PER QUALITY CRITERIA IN THE THREE USER STORY SETS

	Set 1 (n=96)		Set 2 (n=24)		Set 3 (n=124)	
	V	FP	V	FP	V	FP
<i>Atomic</i>	7	5	10	3	17	12
<i>Minimal</i>	0	-	17	-	16	-
<i>Well-formed</i>	8	-	2	-	6	-
<i>Explicit dependencies</i>	0	-	1	-	0	-
<i>Uniform</i>	14	4	2	-	1	-
<i>Unique</i>	2	-	0	-	0	-
<i>Total US with errors</i>	27	9	19	3	37	12

Using these results, we compare the main issues between the user story sets. By looking at AQUUSA’s grouping of all non-uniform user stories of user story set #1, we recognize its primary issue: 10 user stories omit “*want to*”, directly expressing the functionality. Moreover, 5 out of 7 reports of atomic errors are false positives, while the not well-formed user stories are the consequence of technical requirements that are difficult to capture in the user story format.

This is in stark contrast with user story set #2, which contains errors in all AQUUSA’s criteria but one (unique). It is immediately apparent that the majority of these user stories are not minimal. More than 2/3rds of these user stories contain notes, feedback, testing hints, todo’s and/or solution specifications. Although user story set #3 has 17 potential atomic defects according to AQUUSA, 12 are false positives. The real primary issue is minimality, albeit less frequent and less severe than for #2. In 7 cases the requirements engineer has added a reference to a specific document, the remaining 9 breaches contain no relevant patterns.

## VII. RELATED LITERATURE

Despite their popularity among practitioners [1], academic research on user stories is few and far between. The little work that is available, concerns a diverse list of topics. The connection of user stories to code was studied to retrieve reusable test steps [27]. A conceptual method for identifying

dependencies between User Stories [28] was proposed by Gomez and colleagues, relying on the data entities that stories refer to. Along the same lines, a basic tool was developed for writing consistent user stories [29]. Plank et al. reported on the potential of applying NLP to user stories [30]. They proposed that by analyzing source code, comments, bug reports one can establish links between user stories and their implementation progress. Unfortunately, in private communication, the first author indicated that they chose not to pursue this research line any further. In future work, however, we intend to pursue a similar goal by building on the conceptual model presented in this paper.

Applying natural language processing to RE has historically been heralded as the final frontier of requirements engineering. Nowadays, this ambitious objective is understood to be unattainable in the foreseeable future—at least not without a significant, fundamental breakthrough [11]. Nevertheless, a wide variety of contemporary research in RE applies NLP for specific uses: automatically identifying security requirements hidden in other requirements [31], detecting uncertainty in NL requirements [32] or improving NL requirements quality by semi-automatically detecting a range of bad practices [7]. Tools like these are interesting research artifacts, but still far from becoming mainstream in practice.

Arguing that these tools deliver the opposite effect of what they intend, Berry [6] calls for NLP-supported tools that support 100% recall. AQUUSA aims to satisfy this constraint to obtain quality requirements. However, we will investigate the techniques other tools rely upon to determine if some of them have the potential for improving user story quality.

Multiple frameworks exist for characterizing requirements quality, a very vague concept in general. The IEEE Recommended Practice for Software Requirements Specifications is the standard body of work on this subject, defining eight quality characteristics [13]. Unfortunately, most requirements specifications are unable to adhere to them in practice [14], although evidence shows a correlation between high-quality requirements and project success [33].

## VIII. CONCLUSION AND FUTURE RESEARCH

In this paper, we have argued for user stories as an ideal candidate for improving requirements quality using Natural Language Processing (NLP) techniques. They conform to what practitioners in agile development actually do, and detection of errors is possible with 100% recall and high precision. This paper lays the theoretical foundations for such a tool by making three contributions:

- 1) A revised notion of user story quality in the QUS Framework which supports requirements engineers to forge higher quality user stories.
- 2) A conceptual model of user stories that both computers and humans can use to identify improvement points.
- 3) An initial set of properties, with preliminary identification techniques, that pinpoint low-quality user stories.

Based on these theoretical contributions, we design the Automatic Quality User Story Artisan (AQUUSA), a prototype

tool which exposes defects and deviations from good practice in user stories. The promising results of our application of both the QUS Framework and AQUASA to industrial user story sets demonstrates the feasibility and relevance of this work.

This paper paves the way for future work. By studying how requirements engineers apply and experience the QUS Framework in practice, we intend to validate and extend it further. Additionally, we intend to create a fully-functional, robust implementation of AQUASA to verify the expected effectiveness according to our manual analysis. Additionally, this enables quantitative analysis of user story databases and demonstrating the impact of using the tool for longer periods of time on initial user story quality. A key challenge will be to reduce the number of false positives, while maintaining 100% recall. For example, we are investigating whether applying domain and foundational ontologies will improve the tool. Another direction concerns the tool-supported resolution of errors and conflicts: how to assist requirements engineers by suggesting automatic fixes to existing user stories?

#### ACKNOWLEDGEMENTS

The authors would like to thank Floris Vlasveld, Erik Jagroep, Jozua Velle and Frieda Naaijer for providing real-world user story data. Additionally, we would like to thank Leo Pruijt for his comments on an earlier draft of this paper.

#### REFERENCES

- [1] X. Wang, L. Zhao, Y. Wang, and J. Sun, "The Role of Requirements Engineering Practices in Agile Development: An Empirical Study," in *Proc. of the Asia Pacific Requirements Engineering Symposium*, ser. CCIS. Springer, 2014, vol. 432, pp. 195–209.
- [2] M. Cohn, *User Stories Applied: for Agile Software Development*. Redwood City, CA, USA: Addison Wesley, 2004.
- [3] E. S. K. Yu and J. Mylopoulos, "Understanding "Why" in Software Process Modelling, Analysis, and Design," in *Proc. of the International Conference on Software Engineering*. IEEE, 1994, pp. 159–168.
- [4] B. Wake, "INVEST in Good Stories, and SMART Tasks," <http://xp123.com/articles/invest-in-good-stories-and-smart-tasks/>, 2003, accessed: 2015-02-18.
- [5] P. Heck and A. Zaidman, "A Quality Framework for Agile Requirements: A Practitioner's Perspective," *CoRR*, vol. abs/1406.4692, 2014. [Online]. Available: <http://arxiv.org/abs/1406.4692>
- [6] D. Berry, R. Gacitua, P. Sawyer, and S. Tjong, "The Case for Dumb Requirements Engineering Tools," in *Proc. of Requirements Engineering: Foundation for Software Quality*, ser. LNCS. Springer, 2012, vol. 7195, pp. 211–217.
- [7] A. Bucchiarone, S. Gnesi, and P. Pierini, "Quality Analysis of NL Requirements: An Industrial Case Study," in *Proc. of the IEEE International Conference on Requirements Engineering*, 2005, pp. 390–394.
- [8] D. Popescu, S. Rugaber, N. Medvidovic, and D. M. Berry, "Reducing Ambiguities in Requirements Specifications Via Automatically Created Object-Oriented Models," in *Innovations for Requirement Analysis. From Stakeholders' Needs to Formal Designs*, ser. LNCS. Springer, 2008, vol. 5320, pp. 103–124.
- [9] J. Cleland-Huang, B. Berenbach, S. Clark, R. Settini, and E. Romanova, "Best Practices for Automated Traceability," *Computer*, vol. 40, no. 6, pp. 27–35, 2007.
- [10] R. Gacitua, P. Sawyer, and V. Gervasi, "On the Effectiveness of Abstraction Identification in Requirements Engineering," in *Proc. of the IEEE International Requirements Engineering Conference*. IEEE, 2010, pp. 5–14.
- [11] K. Ryan, "The Role of Natural Language in Requirements Engineering," in *Proc. of the IEEE International Symposium on Requirements Engineering*. IEEE, 1993, pp. 240–242.
- [12] N. Maiden, "Exactly How Are Requirements Written?" *IEEE Software*, vol. 29, no. 1, pp. 26–27, 2012.
- [13] IEEE Computer Society, "IEEE Recommended Practice for Software Requirements Specifications," *IEEE Std 830-1993*, 1994.
- [14] M. Glinz, "Improving the Quality of Requirements with Scenarios," in *Proc. of the World Congress on Software Quality*, 2000, pp. 55–60.
- [15] O. I. Lindland, G. Sindre, and A. Sølvberg, "Understanding Quality in Conceptual Modeling," *IEEE Software*, vol. 11, no. 2, pp. 42–49, 1994.
- [16] O. Liskin, R. Pham, S. Kiesling, and K. Schneider, "Why We Need a Granularity Concept for User Stories," in *Agile Processes in Software Engineering and Extreme Programming*, ser. LNBIP. Springer, 2014, vol. 179, pp. 110–125.
- [17] W. N. Robinson, "Integrating Multiple Specifications Using Domain Goals," *SIGSOFT Software Engineering Notes*, vol. 14, no. 3, pp. 219–226, 1989.
- [18] E. Paja, F. Dalpiaz, and P. Giorgini, "Managing Security Requirements Conflicts in Socio-Technical Systems," in *Proc. of the International Conference on Conceptual Modeling*, ser. LNCS, vol. 8217, 2013, pp. 270–283.
- [19] M. Kim, S. Park, V. Sugumaran, and H. Yang, "Managing Requirements Conflicts in Software Product Lines: A Goal and Scenario Based Approach," *Data & Knowledge Engineering*, vol. 61, no. 3, pp. 417–432, 2007.
- [20] "Ambiguity in Requirements Specification," in *Perspectives on Software Requirements*, ser. International Series in Engineering and Computer Science. Springer, 2004, vol. 753.
- [21] Y. Wautelet, S. Heng, M. Kolp, and I. Mirbel, "Unifying and Extending User Story Models," in *Proc. of the International Conference on Advanced Information Systems Engineering*, ser. LNCS. Springer, 2014, vol. 8484, pp. 211–225.
- [22] K. Beck, *Extreme Programming Explained: Embrace Change*. Boston, MA, USA: Addison-Wesley, 2000.
- [23] M. Escalona and N. Koch, "Metamodeling the Requirements of Web Systems," in *Web Information Systems and Technologies*, ser. LNBIP. Springer, 2007, vol. 1, pp. 267–280.
- [24] P.-A. Hsiung, S.-W. Lin, C.-H. Tseng, T.-Y. Lee, J.-M. Fu, and W.-B. See, "VERTAF: An Application Framework for the Design and Verification of Embedded Real-Time Software," *IEEE Transactions on Software Engineering*, vol. 30, no. 10, pp. 656–674, 2004.
- [25] S. F. Tjong and D. M. Berry, "The Design of SREE: A Prototype Potential Ambiguity Finder for Requirements Specifications and Lessons Learned," in *Proc. of Requirements Engineering: Foundation for Software Quality*, ser. LNCS. Springer, 2013, vol. 7830, pp. 80–95.
- [26] G. A. Miller, "WordNet: A Lexical Database for English," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, 1995.
- [27] M. Landhäuser and A. Genaid, "Connecting User Stories and Code for Test Development," in *Proc. of the International Workshop on Recommendation Systems for Software Engineering*. IEEE, 2012, pp. 33–37.
- [28] A. Gomez, G. Rueda, and P. Alarcn, "A Systematic and Lightweight Method to Identify Dependencies between User Stories," in *Proc. of the International Conference on Agile Software Development*, ser. LNBIP. Springer, 2010, vol. 48, pp. 190–195.
- [29] M. Śmiałek, J. Bojarski, W. Nowakowski, and T. Straszak, "Writing Coherent User Stories with Tool Support," in *Proc. of the International Conference on Agile Software Development*, ser. LNCS. Springer, 2005, vol. 3556, pp. 247–250.
- [30] B. Plank, T. Sauer, and I. Schaefer, "Supporting Agile Software Development by Natural Language Processing," in *Trustworthy Eternal Systems via Evolving Software, Data and Knowledge*, ser. CCIS. Springer, 2013, vol. 379, pp. 91–102.
- [31] M. Riaz, J. King, J. Slankas, and L. Williams, "Hidden in Plain Sight: Automatically Identifying Security Requirements from Natural Language Artifacts," in *Proc. of the IEEE International Requirements Engineering Conference*. IEEE, 2014, pp. 183–192.
- [32] H. Yang, A. De Roeck, V. Gervasi, A. Willis, and B. Nuseibeh, "Speculative Requirements: Automatic Detection of Uncertainty in Natural Language Requirements," in *Proc. of the IEEE International Requirements Engineering Conference*. IEEE, 2012, pp. 11–20.
- [33] M. I. Kamata and T. Tamai, "How Does Requirements Quality Relate to Project Success or Failure?" in *Proc. of the IEEE International Requirements Engineering Conference*. IEEE, 2007, pp. 69–78.