

Leveraging teenagers feedback in the development of a Domain-Specific Language – the case of programming low-cost robots

Abstract—Domain Specific Languages (DSLs) empower end-users to express software tasks that were traditionally developed by software engineers. DSLs allow users to express themselves in terms closer to the way they think about their problems, rather than in computational terms. However, conceiving a DSL with an adequate user experience for its end-users is not a trivial task, and the process of engineering that adequacy tends to be performed *ad-hoc*.

VL¹ is an open-source DSL for controlling low-cost rover-like Arduino robots, designed for being used by teenagers with no previous computing skills, so they can be introduced to programming in a fun way. In this paper, we discuss an iterative process building on teenagers’ early feedback, collected in a series of empirical evaluations with 128 teenagers, and how this has helped us driving VL to a competitive level in terms of usability, when compared to well-established alternatives such as Lego, or Scratch.

I. INTRODUCTION

The increasing software pervasiveness fosters a growing concern for making some of its development accessible to end-users with no formal training in programming. Creating Domain-Specific Languages (DSLs) for empowering end-users is challenging, as we need to bridge the gap between computation concepts and the concepts mastered by the end-users. Two complementary ways of bridging this gap are (i) nurturing ‘computational thinking’ [1] skills in end-users, and (ii) devising adequate metaphors that hide the unnecessary complexity of computational concepts from those end-users. We were contacted by a company interested in developing an open source web-based DSL, called VL, targeted to teenagers, to control low-cost rover-like Arduino robots. The challenge was how to assess the DSL in a timely way, so that end-user feedback could lead to a competitive product. In particular, the company was concerned with the user-friendliness of the DSL.

However, the development process of DSLs lacks a systematic and iterative approach to evaluating and detecting usability issues since the early stages of the DSL construction.

Therefore, as the goal is to empower teenagers with an adequate DSL for programming low-cost rover-like robots, we add to the common iterative life cycle of DSL development an evaluation task involving the end-users, to be performed in each iteration. The teenagers’ feedback is used to help to steer the VL’s evolution through the identification of several improvement opportunities in the language. This evaluation stage in each iteration is often not reported in the context of

developing DSLs but is key for our development effort. In this assessment, we contrast VL with two popular DSLs that are used to control rover-like robots: a commercial competitor (Lego) and an open source initiative (Scratch).

We report on the design and results of the empirical studies used in this evaluation that helped us identifying the language’s strengths and weaknesses. We discuss how this led to the improvement of the VL language. To help to achieve this higher level goal, we answer two more detailed research questions:

- **RQ1:** *How does the current VL (VL2) compare to baselines (a previous version of VL (VL1), Lego and Scratch) regarding the **Effectiveness** of the teenagers when programming a robot?*
- **RQ2:** *How does the current VL (VL2) compare to baselines (a previous version of VL (VL1), Lego and Scratch) regarding the **Satisfaction** of the teenagers when programming a robot?*

The chosen baselines are aimed at providing a comparison basis with (1) a previous version of VL (VL1), so that we can assess the extent to which the feedback collected with that previous version has helped in its evolution and (2) two popular competitor languages, so that we can assess how competitive the VL language can be, when contrasted with those languages.

II. RELATED WORK

A. Language Usability

Language usability can be defined as the degree to which a language can be used by specific users to meet their needs to achieve specific goals with effectiveness, efficiency and satisfaction in a specific context of use (adapted for the particular case of languages from [2]).

User-centered design (UCD) [3], [4] can contribute to more usable DSLs. For example, [5] presented an innovative visualisation environment, which eases and makes more effective the experimental evaluation process, implemented with the help of UCD. A visual query system was also designed and implemented following the UCD approach [6]. Although there is a lack of general guidelines and best practices to conduct language usability evaluations, they are slowly being recognised as an important step in the Language Engineering life-cycle [7]. An iterative approach allows us to trace usability requirements and the impact of usability recommendations throughout the DSL development process [8].

¹renamed here for double blind review purposes

B. User characterization and evaluation

Teenagers who are familiar with computers and technologies tend to be more successful in new computer-related tasks [9]. Also, they are likely to work and play in groups, (e.g. sharing a single computer) [10]. Involving teenagers as subjects in empirical studies is a valid option to evaluate usability on software products targeted to teenagers [11] but also challenging, as teenagers have a high variability in cognitive and development abilities at a given age [12]. Nevertheless, it has been shown that teenagers can identify and report usability problems using methods like ‘think aloud’, interviews or questionnaires [13].

Existing studies (e.g. [14], [15], [16]) concentrate on issues related to the curricula of K12 courses, the motivation of students to engineering, and effective education of computational teaching. Teenagers have been used in the past as subjects for language-related studies. For example, a hybrid approach combining text and visual notations offers the best compromise to increase efficiency and effectiveness of teenagers while using that language [17]. However, to the best of our knowledge, there are no formal studies concerned about improving software language usability, involving teenagers as subjects.

C. Programming languages for teenagers

Current robot platforms offer Application Programming Interfaces (APIs) for programming with languages such as *Java* (e.g. [18]), [19]) or *.NET Framework* languages (e.g. [20]). These languages may exclude many teenagers who are not (yet) proficient in programming. Without proper training, it is hard for them to program and master a textual programming language with a complex syntax full of technical concepts. This is true even with the help of powerful Integrated Development Environments (IDEs). We argue it is far from trivial to design a DSL for teenagers (and even more so for younger children).

There are several programming languages designed specifically for children (including teenagers). Examples include *Alice* [21], *Blockly* [22], *Lego* [23] and *Scratch* [24]. We used for comparison purposes with VL *Lego*, one of the most widely used languages for programming toy robots, and *Scratch*, a popular visual programming language for teaching programming concepts to children which can also be used for controlling rover-like robots.

1) *Lego*: Inspired by Papert’s seminal work [25], *Lego Mindstorms (LegoMS)* is one of the most well-known technologies, in educational robotics for children. Several case studies have been done for this technology in order to understand which are the main individual needs of the children, when working with robots. *LegoMS*’s technology combines hardware and software, so that users can develop and deploy programs specifying behavior to be executed by the *Lego* robot.

Lego Mindstorms NXT 2.0, used in our usability trials, presents building blocks (like bricks) as elements to build a program. Each block represents a programming concept, such as an execution control element, e.g. loops, conditions, arithmetic, or an actuate block that interacts with the robot

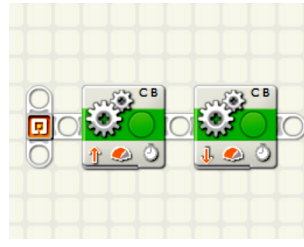


Fig. 1. LegoMS

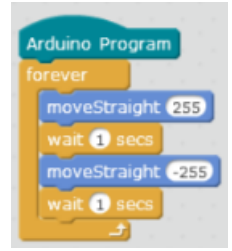


Fig. 2. Scratch

components, e.g. motors. The sequence of blocks is constructed by a behavior flowchart, structuring the program’s blocks.

LegoMS has an appealing notation and is used only with relatively expensive *Lego* robots. However, the development of complex behaviors may become difficult [26]. The increasing number of blocks and the size of the diagram make it difficult to analyze and read the program solution. Non-trivial behaviors are difficult to implement in a visual programming language like this. Fig. 1 exemplifies the move back and forth example. The diagram contains two composition elements describing how the robot should move.

2) *Scratch*: Scratch [24] relies on actions (Blocks) to operate specific objects (Sprites). They can be seen as a visual abstraction to the Object Oriented programming paradigm with some restrictions (no support for custom objects and the dynamic generation of sprites). Open source tools like *mBlock* [27] and *Enchanting* [28], are built upon *Scratch*, and meant to be simple and used for robot programming. *mBlock* is solution which compile to the *Arduino* open source hardware platform, which makes it suitable for the our assessment studies, as it use the same hardware as VL.

Fig. 2 illustrates how to program the robot to move forward during one second and to go backwards during one second. However, as this language is general purpose, it does not have abstractions of operations like move forward or backwards. Therefore, the user needs to detail the movement operator. This includes to describe the pins of each motor on the *Arduino* board and deal with concepts like motor rotation directions, rotation speed, angular velocity.

III. CONTEXT

VL is a visual language that allows the user to implement programs for robot behavior, through the manipulation of visual elements, or objects. This manipulation is expected to help the user to understand quickly the programming mechanisms. Visual languages are thought useful for introducing programming concepts to children [29], while robots perform the developed code in the real world. The observation of the program running on a physical robot provides an engaging feedback on the implications of changes in the program.

The robot is an *Arduino* board with a set of sensors and actuators. *Arduino* is an open-source prototyping platform [30]. The low-cost robot (see Fig. 3) which works with VL has the following configuration: (i) An infra-red distance sensor; (ii) Two bumpers (e.g. collision detectors sensors); (iii) One motor

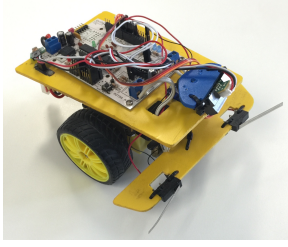


Fig. 3. Arduino Robot

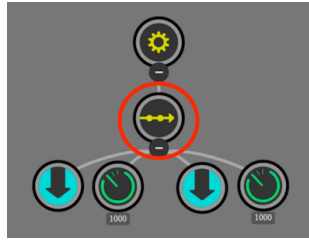


Fig. 4. VL1 - Back and Forward

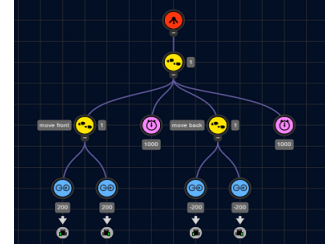


Fig. 5. VL2 - Back and Forward with bumpers

for each wheel; (iv) A Servo motor that actuates under the distance sensor; (v) A simple LED light.

Programming the *Arduino* textual code requires technical skills not owned by most teenagers. To mitigate this problem, the robot may be programmed using *VL*, which is then automatically translated into the *Arduino* textual code.

VL's visual syntax is based on the behavior tree paradigm [31], a mathematical model of plan execution, used for a diversity of areas including robotics, control systems and software games. A complex behavior is mapped into smaller and simpler behaviors through its branches. This descending order of complexity provides a structured way of defining complex behaviors (which are used to define objectives) through simple tasks defined hierarchically. Each node may have a specification that determines how the actions of its children will be executed (in parallel, or sequentially). The child node returns its status to the parent node, and this successively happens until the root of the tree.

Fig. 4 illustrates how to program the robot to move back and forth. The root node is decomposed into a single sequential node (highlighted with a red circle). The sequential node runs all its children (in this case, the leaf nodes) in depth first traversal order. Leaf nodes represent the most primitive actions that could be taken by an agent. In Fig. 4 the nodes represented by arrows pointing downwards are outputs actions. In this case they are left and right motors commands. The outputs can also be visual using LEDs or audio using a buzzer. Both nodes displayed by 2 round timer shapes are wait commands in which the execution of the next node is delayed by the duration determined by the end user for each node.

1) *From VL0 to VL1*: The first development iteration of *VL* VL0 started with a domain analysis, followed by the design, implementation, and evaluation, in order to quickly deploy an early prototype. The last phase involved 22 (10+12) children, with the age range from 8 to 12, as subjects included in an exploratory study [X]². Children under 10 had a hard time learning *VL*. Further development resulted in the next release VL1 (Fig. 4), which improved the interaction model and provided a web based solution.

2) *From VL1 to VL2*: VL1 was evaluated experimentally and compared to the one of the most popular commercial competitors in the market, *LegoMS*. Based on the results of this empirical study, some improvements were suggested and

applied to *VL*. The focus was on improving the user interface providing better readability of the programs being developed and improving error prevention. The suggestions were:

- Improve error feedback and suggestions to solve problems.
- Add auto-alignment of new nodes added to a sequence node, preventing situations where the visual order of nodes from left to right did not match the order of their execution.
- Highlight when a user selects a new node to add to the diagram the nodes that are available to form a connection. This provides better user assistance and error prevention.
- Allow the users to create blocks of tree structures. This feature promotes reuse of previously developed structures and allows to share and slowly introduce more complex behaviours to novice users.
- Show the Icon's label with mouse-over events. This helps the user to recognize the available options.
- Introduce a different method of collapsing and expanding tree structures. The use of the keys "+" and "-" for these operations was not clear to the user.
- Introduce new zooming options allow an easier navigation of the tree structure by introducing a fit to screen operation of the entire tree structure or only the selected nodes.

These suggestions lead to the release of a new improved version of the language *VL2* (Fig. 5). Finally, *VL2* was compared with *Scratch*.

IV. EXPERIMENT PREPARATION

In this section, we describe the experiment protocol. Additional materials can be found in this paper's companion site [32]. We had two runs of the same experiment, with the second prototype release (VL1) which we compared to *Lego*; and second *VL* version (VL2) which was compared to *Scratch*. The experiment was announced as a robot programming challenge to engage teenagers to participate.

During the first run, VL1 was compared to the best product in the area but appropriate only for costful *Lego* robots. Therefore, the robots used in this run were different. However, during the second experiment run, both VL2 and *Scratch* programs were running on the same robot (see Figure 3). The two different experiment runs were designed as two parts of the same experiment. Both participant's profile and the tasks were similar, so that, in practice, we can think of them as a single experiment with four different languages (two of which are different versions of the same language).

²Hidden for double blind review purposes

A. Experiment objectives

The high-level objective of our study is to evaluate the usability of VL, “*Usability of Programming the robot*”, as it is expected to be used by a broad group of people that are not expected to have previous programming skills.

As we already have different generations of a functional prototype of VL, we measure VL’s *effectiveness* and *satisfaction*. The metrics for these usability requirements are defined in Table I. We are not particularly concerned with other requirements such as efficiency or learnability at this point because we are at a stage of the language development where we want to know if the teenagers can program the robot, which is a prerequisite for measuring other characteristics. The experiment objectives follow the GQM template [33] and are defined as follows:

G1 - Analyse the effect of VL2, **for the purpose** of evaluation, **with respect to** its impact on the *effectiveness* in programming a robot **when compared to** three *baselines*, namely VL1, *Lego* and *Scratch*, **from the point of view** of researchers, **in the context of** an experiment conducted with secondary school subjects.

G2 - Analyse the effect of VL2, **for the purpose** of evaluation, **with respect to** its impact on the *satisfaction* in programming a robot **when compared to** three *baselines*, namely VL1, *Lego* and *Scratch*, **from the point of view** of researchers, **in the context of** an experiment conducted with secondary school subjects.

TABLE I
USABILITY REQUIREMENTS

Requirement	Metric
<i>Effectiveness</i> : Is the <User> able to correctly implement a given <Use Case>?	PCorrUCInst - percentage of correctly implemented concepts in <Use Case>
<i>Satisfaction</i> How much is the user satisfied with VL?	ConfLevel - self rated confidence score in a Likert scale LikeLevel - self rated likeability score in a Likert scale LearnLevel - self rated learnability score in a Likert scale

In particular we test the following (null) hypotheses:

- $H1_0$: Using VL2 has no influence in the *effectiveness* of programming the robot when compared to programming the robot with the *baselines* VL1, *Lego* and *Scratch*.
- $H2_0$: Using VL2 has no influence in the *satisfaction* of programming the robot when compared to programming the robot with the *baselines* VL1, *Lego* and *Scratch*.

B. Experiment design

We used a *between groups* design where participants were randomly assigned to the task of either programming a robot using VL1 or *Lego*, in the first run, or VL2 or *Scratch*, in the second run. Each participant only participated in one of the four alternatives to avoid learning effects. The experimental process was similar for all runs and subject groups, starting with a learning session that lasted 30 minutes, during which the participants filled in a background questionnaire and learned about programming robots using their assigned language. Then, the participants had 15 minutes to solve a programming

challenge. The contents of the computer screen during the session were recorded. Finally, there was a feedback session, taking up to 5 minutes, to collect the teenager’s subjective opinions about using the language they were assigned to.

1) *Participants, teams and groups*: The participants were high school students recruited through convenience sampling, among the visitors of two different open doors days in the same university. In each day, participants were randomly assigned to one of the two groups: in the first day, there was a group with VL1 and another one with *Lego*; in the second day, there was a group with VL2 and another one with *Scratch*. Students and their teachers were aware of this study and volunteered to participate. Two lecturers guided each language’s groups. The teenagers were requested to participate, in teams of up to three elements, using their assigned language and corresponding robots. Each group had a maximum of 5 teams participating in the same session. Undergraduates helped in the experiment, while a researcher monitored the data collection process.

2) *Technical, social and physical environment*: Each team worked on a desktop computer with OS Windows 7 Professional (Athlon 64x2 Dual Core 5000 2.6 GHZ processor, 4 GB RAM and a 17-inch monitor with a screen resolution configured at 1280x720) to implement the applications. The interaction between user and computer was achieved by the use of keyboard, mouse, and screen, and was captured by *Debut* video recorder³.

The atmosphere was set up to be challenging and educative, but also playful and entertaining, to keep the teenagers still interested to participate and to reduce the sensation of failure. Team members were allowed to talk to each other. They could seek help regarding the solution of the exercises presented or if they had technical issues with the robots.

3) *Exercises and Challenge*: During the learning sessions, participants solved three basic robot programming exercises. They tested each of those exercises on the corresponding robot. Furthermore, they were encouraged to ask questions and to ask for help if necessary.

The first exercise was to program the robot to *move forward* and then *move back* (see Fig. 6). It gave the participants a notion of how to program the robot to move in a basic way. The second exercise was to program the robot to move along a path similar to a ‘5’ (see Fig. 8). In each turn, the robot needed to execute the move operation using the same amount of time and then to make a 90° angle turn. The third exercise was to program a robot to *move forward* until it *bumped* into some object, then it would *move back* and *stop* (see Fig. 7). During this exercise, participants learned how to use the detector sensors and to compose their behaviour in a sequence with other modelling elements.

The challenge in the first experiment was to program the robot to make a shape of ‘5’ but the robot would only turn when it would hit the bumper (see Fig. 9). It was composed by basic user story actions of the exercises. The challenge in second experiment was a slightly different (see Fig. 10). The

³<http://www.nchsoftware.com/capture/>

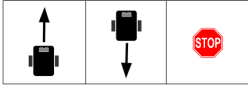


Fig. 6. Exercise 1

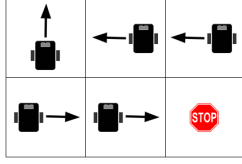


Fig. 8. Exercise 2

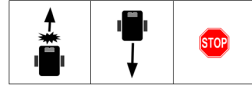


Fig. 7. Exercise 3

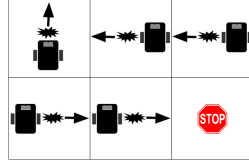


Fig. 9. Challenge1

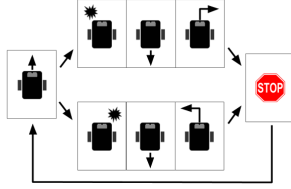


Fig. 10. Challenge2

robot was expected to detect if it hit the obstacles with the left or the right bumper, and was supposed to turn accordingly to the opposite direction. When the team thought they had a solution, they were invited to physically test it in the arena with a given robot. If they were not happy with the result of the test, they could go back and try to fix whatever was wrong with their solution.

C. Experimental instruments and measurements

During the experiment, we used survey forms, video recordings and a competition arena to collect data for further analysis. The survey forms were composed of “Smileyometers” which are found to be appropriate for teenagers questionnaires [34]. While answering to forms, teenagers were assisted by an adult (one of the experiment assistants), to ensure that there were no misinterpretations of questions and answers and to confirm that participants did not experience reading problems. As we grouped the participants into teams, the participants’ individual answers to questionnaires were merged. We computed the mean response within each group, for each answer.

TABLE II
INSTRUMENTS AND SCALES

	Instrument	Value
Profile	Background Questionnaire	0 / 0.5 / 1
Effectiveness	Video recording	0 / 1
Satisfaction	Satisfaction Questionnaire	-1 / 0 / 1

Profile is a measure influenced by *Experience* factors on Computer Games, Programming or Programming a robot, and *Tendency* factors reflecting the teenager’s tendency to Mathematics, Physics or to Learn programming. The data for calculating the *Profile* was collected through the background questionnaire, that consisted of questions designed to assess those pre-defined factors. Each answer was encoded in three

possible answers: ‘Yes’ (with a score of 1), ‘Intermediate’ (0.5), and ‘No’ (0).

We used the recorded videos to evaluate the **Effectiveness**. The challenge could be solved by composing elements of the training exercises, which are marked as Success (S) or Failure (F). Effectiveness measures the percentage of modelling elements correctly built and composed to achieve the solution.

Satisfaction is characterized by: a *Confidence* factor, reflecting ConfLevel metric that tells how confident teenagers were about their solution; *Likeability* factor, reflecting LikeLevel metric that tells how interesting and enjoyable they found the challenge itself; and, *Learnability* factor, reflecting LearnLevel metric that tells how useful they found what was taught during the learning session which helped them to face the final challenge. The data was collected through a satisfaction questionnaire, that consisted of questions designed to assess the defined factors. A ‘Yes’ scored 1 point, ‘Intermediate’ scored 0, and ‘No’ scored -1.

V. RESULTS

A. Descriptive statistics

Table III presents the descriptive statistics for the metrics collected in our data analysis. In the *Characteristic* column we present the property under scrutiny. For each of these characteristics, we show four rows, one for each *Language* (in the second column). We further detail the *mean*, *standard deviation*, *skewness*, *kurtosis*, and the *p-value* for the Shapiro-Wilk normality test. The number of participant teams is not always the same for all languages. The shape of the distributions concerning most of the variables suggests that, in general, normality is not a reasonable assumption ($p\text{-value} < 0.05$). The exceptions are metrics concerning *profile*, *experience* and *tendency*, where, for most languages, normality is a reasonable assumption ($p\text{-value} \geq 0.05$). The visual inspection of boxplot diagrams, Q-Q plots and kernel density plots (omitted here for the sake of brevity) further reinforced our assessment concerning data normality. As several of the variables have a non-normal distribution, we assume this for the remainder of our analysis. Note also that all teams expressed the maximum *confidence* level on their solution for VL2, as well as the maximum *likeability* level for VL2 and Scratch.

The variance of the distributions is not similar, when comparing the characteristics metrics distributions for the different languages, as shown in table IV.

B. Dataset preparation

The video recordings were analysed following a protocol previously established by the research team. Apart from (i) assessing the success, we checked if the team (ii) reused the concepts from previous exercises, (iii) experienced technical problems or functional errors, (iv) had interaction difficulties (e.g. using copy/paste for visual objects or connecting the same objects in sequence), (v) reused previously constructed sequences (within the same exercise), or (vi) used any other additional language features (e.g. zooming). The remaining data was extracted from the questionnaires [32].

TABLE III
DESCRIPTIVE STATISTICS

Characteristic	Language	N	Mean	Std. Dev.	S-W
Age	Scratch	8	16.7062	1.23782	0.023
	VL2	9	16.3689	0.91408	0.180
	VL1	17	16.6059	1.31694	0.018
	Lego	14	16.4429	1.28465	0.039
Profile	Scratch	8	0.6875	0.17107	0.542
	VL2	9	0.6019	0.13029	0.213
	VL1	17	0.6225	0.16435	0.455
	Lego	14	0.6250	0.22349	0.201
Experience	Scratch	8	0.5625	0.23465	0.241
	VL2	9	0.4259	0.14699	0.338
	VL1	17	0.5294	0.24463	0.214
	Lego	14	0.4643	0.31473	0.042
Tendency	Scratch	8	0.8125	0.18767	0.197
	VL2	9	0.7778	0.20412	0.122
	VL1	17	0.7157	0.20211	0.059
	Lego	14	0.7857	0.2210	0.004
Effectiveness	Scratch	8	0.8000	0.23905	0.041
	VL2	9	0.9556	0.08819	0.000
	VL1	17	0.3765	0.40548	0.001
	Lego	14	0.7857	0.32783	0.000
Satisfaction	Scratch	8	0.9583	0.05893	0.002
	VL2	9	0.9630	0.07349	0.000
	VL1	17	0.6275	0.23221	0.085
	Lego	14	0.8571	0.22746	0.000
Confidence	Scratch	8	0.9688	0.05786	0.000
	VL2	9	1.0000	0.00000	-
	VL1	17	0.5000	0.30619	0.234
	Lego	14	0.8571	0.25409	0.000
Learnability	Scratch	8	0.9063	0.12939	0.000
	VL2	9	0.8889	0.22048	0.000
	VL1	17	0.4412	0.42875	0.023
	Lego	14	0.7500	0.37978	0.000
Likeability	Scratch	8	1.0000	0.00000	-
	VL2	9	1.0000	0.00000	-
	VL1	17	0.9412	0.24254	0.000
	Lego	14	0.9643	0.13363	0.000

TABLE IV
TEST OF HOMOGENEITY OF VARIANCES

Characteristic	Levene Statistic	df1	df2	Sig.
Age	0.353	3	44	0.787
Profile	1.758	3	44	0.169
Experience	3.588	3	44	0.021
Tendency	0.537	3	44	0.660
Effectiveness	6.975	3	44	0.001
Satisfaction	2.892	3	44	0.046
Confidence	8.079	3	44	0.000
Learnability	2.828	3	44	0.049
Likeability	1.711	3	44	0.178

C. Hypotheses testing

For testing our hypotheses, we used the Welch t test, as it is robust to deviations from the normal distribution, different sample sizes and different variance in the samples, thus following the recent recommendations on data analysis for Software Engineering empirical evaluations (which summarises best practices in statistical analysis on other domains) [35]. Table V summarizes the results of these tests. Note that, because *Confidence* and *Likeability* had a constant value (the

top possible score for each of them, the corresponding lines are not filled in table V.

TABLE V
WELCH T TEST SCORES

	Statistic	df1	df2	Sig.
Age	.170	3	21.109	.915
Profile	.423	3	20.891	.739
Experience	.921	3	21.381	.448
Tendency	.527	3	20.597	.669
Satisfaction	10.723	3	24.108	.000
Confidence	-	-	-	-
Learnability	5.678	3	23.334	.005
Likeability	-	-	-	-
Effectiveness	10.886	3	20.483	.000

We used a Games-Howell post-hoc test to determine which languages were significantly different from which languages, according to our set of characteristics under scrutiny. Table VI summarises this test's result, for the comparisons involving either VL1 or VL2, or both. We observe that VL1 lead to a significantly lower *Satisfaction* and *Effectiveness* when compared to Scratch, VL2 and Lego. VL1 was also significantly harder to learn than Scratch and VL2, but not significantly different when compared to Lego. In contrast, VL2 was consistently as good as Lego and Scratch, while superior to VL1 in terms of *Satisfaction*, *Learnability* and *Effectiveness*.

TABLE VI
GAMES-HOWELL TEST

Characteristic	(I)Tool	(J)Tool	MD(I-J)	Std. Err.	Sig.
Satisfaction	VL2	Lego	.10582	.06554	.397
	VL2	Scratch	.00463	.03216	.999
	VL2	VL1	.33551	.06142	.000
	VL1	Lego	-.22969	.08287	.046
	VL1	Scratch	-.33088	.06005	.000
Learnability	VL2	Lego	.13889	.12531	.689
	VL2	Scratch	-.01736	.08657	.997
	VL2	VL1	.44771	.12734	.009
	VL1	Lego	-.30882	.14531	.169
	VL1	Scratch	-.46507	.11360	.003
Effectiveness	VL2	Lego	.16984	.09242	.293
	VL2	Scratch	.15556	.08948	.362
	VL2	VL1	.57908	.10264	.000
	VL1	Scratch	-.42353	.12967	.018
	VL1	Lego	-.40924	.13171	.021

RQ1: How does the current version of VL (VL2) compare to the used baselines (VL1, Lego and Scratch) regarding the Effectiveness of the children when programming a robot?

As seen in table V, there was a statistically significant difference among languages, with respect to the overall *Effectiveness*, with VL2 (M=.9556; SD=.08819), Scratch (M=.8; SD=.23905) and Lego (M=.7857; SD=.32783) clearly allowing the teenagers using it to outperform those using VL1 (M=.3765; SD=.40548), as detailed in table VI. This suggests that **VL2 has achieved a similar level of Efficiency compared to the two commercial baseline languages, and has significantly improved when compared to its previous iteration, VL1,**

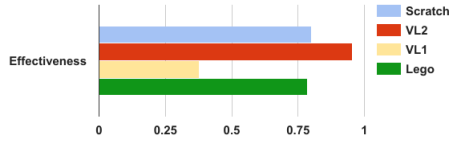


Fig. 11. Effectiveness results

in this aspect. Figure 11 illustrates this. As such, we can reject our null hypothesis H_01 .

RQ2: How does the current version of VL (VL2) compare to the used baselines (VL1, Lego and Scratch) regarding the Satisfaction of the teenagers when programming a robot?

As seen in table V, there was a statistically significant difference among languages, on the overall satisfaction, with Scratch ($M=.9583$; $SD=.05893$), VL2 ($M=.963$; $SD=.07349$) and Lego ($M=.85710$; $SD=.22746$) providing a higher satisfaction than VL1 ($M=.62750$; $SD=.23221$). This suggests that **VL2 has achieved a similar level of Satisfaction compared to the two commercial baseline languages, and has significantly improved when compared to its previous iteration, VL1, in this aspect.** Figure 12 illustrates this. As such, we can reject our null hypothesis H_02 .

We can further break down this observation with a closer look to the components of *Satisfaction*: *Confidence*, *Learnability* and *Likeability*. *Confidence* obtained a perfect score ($M=1.0$; $SD=.0$) for VL2, closely followed by Scratch ($M=.96880$; $SD=.05786$) and Lego ($M=.8571$; $SD=.25409$), but contrasting to VL1 ($M=.5$; $SD=.30619$). Concerning *Learnability*, Scratch ($M=.90630$; $SD=.12939$) and VL2 ($M=.8889$; $SD=.22048$) had a statistically significantly higher score than VL1 ($M=.4412$; $SD=.42875$), but not significantly higher than Lego ($M=.75$; $SD=.37978$). Finally, the perfect scores of Scratch and VL2 concerning *Likeability* ($M=1.0$; $SD=.0$) were not statistically significantly different from those of Lego ($M=.9643$; $SD=.13363$) and VL1 ($M=.9412$; $SD=.24254$).

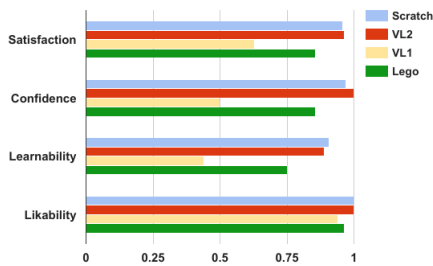


Fig. 12. Satisfaction feedback

Participants Background:

We also need to stress that the observed differences are **not** attributable to different participant backgrounds. Indeed, the results from the background questionnaires indicate a comparable profile of participants for all languages. This is visible from table V, where none of the properties *Age*, *Profile*,

Experience and *Tendency* are statistically significantly different for any of the languages, i.e., they are essentially comparable.

Most participants had some experience in playing computer games, but very few of them had previously programmed a robot (see Figure 13). Some of the participants also had some knowledge of programming. These three factors gave an average Experience score for teams. The motivation to participate in the challenge was high. Most subjects expressed their tendency to learn to program, mathematics and physics, leading to a high Tendency score for all groups. The Profile score, calculated as an average of Experience and Tendency, indicates balanced teams regarding their background.

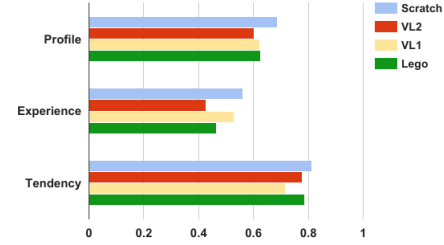


Fig. 13. Profile analysis

VI. DISCUSSION

We presented a systematic usability evaluation of the VL language. In this section, we discuss the contributions to the software language engineering process, evaluation results and finally the implications for the further development.

A. Contributions to the development cycle

By designing the systematic experimental study, it was necessary to describe the context of use of VL during the experiment and its goals explicitly. We have identified “Who will use the language?” and characterised the intended group of end users by the Profile, where each of the relevant characteristics was measured with a particular set of questions and can be further reused and analysed for selection of experimental subjects and specification of VL’s audience. We explicitly tackle the question “Where will the language be used?” by defining its technical, social and physical environment. Also, we systematically analysed the working environment of the VL’s competitors’ tools to further identified both limitations and advantages of VL. Further, we explicitly expressed “How will the language be used?” through user stories. Developers performed the functional testing of these user stories, and the teaching artefacts (i.e. presentation documents and videos) were produced, which can be now be reused as a part of language documentation. The readability and understandability of these presentation materials were validated during pilot experiment sessions and materials can be reused in further evaluations.

Additionally, we address the question “Why will the language be used?” by defining the usability requirements and their GQM definition. Our high-level goal to evaluate usability is context-dependent, and therefore we see that its success is linearly

related to the validated requirements (e.g. effectiveness and satisfaction) and a given context (e.g. evaluating a Challenge scenario in a similar experiment with LegoMS, in the first run of the comparative evaluation, using VL1, and Scratch in the second run of the experiment, this time competing with VL2). The trial design was successful and can be reused and save time in further evaluation sessions. The experiment results helped in making the further design decisions and identifying the pros and cons of the previous implementation.

B. Evaluation results

RQ1: How does VL2 compare to the baselines (VL1, Lego and Scratch) in terms of the Effectiveness of the teenagers when programming a robot? Based on the results obtained in the first run we conclude that programming the robots with VL1 resulted in lower Effectiveness scores when compared to programming the robot with Lego. However, thanks to the evaluation with the VL1 prototype, we found that VL was already usable to some extent by teenagers. The feedback and observations of that initial run gave us insights into particular features that should be implemented to improve VL's usability. This helped to steer the development of VL2, that was then tested in the second run of the experiment, this time using Scratch as a competitor (although, in practice, both VL1 and Lego can be seen as competitors). The feedback collected in the previous iteration proved extremely valuable, leading to a significant improvement of the *effectiveness* from VL1 to VL2, which is now on par to Lego and Scratch.

RQ2: How does VL2 compare to baselines (VL1, Lego and Scratch) in terms of the Satisfaction of the teenagers when programming a robot? Results showed that both VL1 and Lego were rated as satisfactory, with Lego providing a more intense satisfaction level. This shows that, although not being as successful as with Lego, the teams solving the challenge using VL1 still had an entertaining and motivating experience. Concerning *Satisfaction*, VL2 is significantly more competitive than VL1 and on par to Scratch and Lego.

C. Implications for the development of VL

As a consequence of the analysis of all the data collected and the observation of the interaction of the experiment's participants with the VL language, the VL development team identified that new features are needed, and they will be included in future releases.

The focus is on improving the following areas:

Error Prevention and program readability needs: 1) Improve errors' feedback (including suggestions on how to solve the problems found); 2) When a user selects a new node to add to the diagram, the nodes that are possible to connect to the new node should become highlighted; 3) The users should be able to create blocks of tree concepts (this will promote reuse and facilitate the sharing and introduction of complex behaviours to novice users); 4) To prevent misinterpretations of the execution of a sequence node, an auto-alignment of new nodes added to a sequence node feature has to be introduced; and, 5) Icon

Labels should be presented on mouse-over events (this should improve the user recognition of the available options).

Diagram navigation: A *fit to screen* of the entire program or only the selected nodes should be available; an improved tree collapse and expand feature (the former approach, using keys "+" and "-" was not clear to the user) should be added.

VII. THREATS TO VALIDITY

Usually, it is not safe to rely on teenagers self-rating questions. To maximise the reliability of the responses, we had adult helpers interviewing them. This helped to ensure the validity and integrity of results and gave strength to design recommendations or decisions.

We did not have participants using both languages. Although this design prevents learning effects, it does create the risk of, by accident, having more "competent" teams using one language, or the other. However, by obtaining similar background scores (Profile and Age) in the four groups, we are confident that this threat was mitigated.

We compared two approaches with different robots, in the first experiment run: one using *Arduino*, and the other using *Lego* hardware. This might have introduced a bias in the results if one of the robots was easier to program than the other for some reason not directly associated with the programming language. To mitigate this threat, the second run of the evaluation contrasted VL2 with the Scratch for *Arduino* language. So, in this second run, the robot was the same for both languages (and also the same used with VL1).

The choice of *LegoMS NXT* as the platform to test against VL1 can also be regarded as a potential threat. Lego recently released the *Mindstorms EV3* platform that introduces improvements in their development software and robots. However, in the second run of the experiment, we used a different, but also quite popular language (although not EV3), so VL2 is progressively being compared to other alternative languages, rather than just to *LegoMS NXT*. This diversity is expected to mitigate the effects of using a single comparison point.

VIII. CONCLUSION

We reported how we involved teenagers, the end users, in several iterations of the engineering process of a programming language for low-cost robots and performed usability studies. This language, VL, was contrasted with *LegoMS* and *Scratch*.

The evaluation described in this paper, thanks to the involvement of the end-users (teenagers) since early stages of the development process of the language, was helpful to timely detect, prioritise, and improve crucial usability aspects of VL by identifying its strengths and weaknesses. We observed the convergence of the visual language to the degree of usability (regarding satisfaction and effectiveness) achieved by existing mature and commercial languages.

As future work, we intend to apply the same described technique in the development of DSLs for other purposes and for different end-user profiles. Also, we foresee the need for developing tools to support software language developers to deal with the significant overhead of the assessments and track the improvements on the different features of the language.

REFERENCES

- [1] J. M. Wing, "Computational thinking," *Commun. ACM*, vol. 49, no. 3, Mar. 2006.
- [2] International Standard Organization, "ISO/IEC FDIS 25010:2011 systems and software engineering – systems and software quality requirements and evaluation (SQuaRE) – system and software quality models," March 2011. [Online]. Available: http://www.iso.org/iso/catalogue/catalogue_tc/catalogue_detail.htm?csnumber=35733
- [3] D. A. Norman and S. W. Draper, "User centered system design," *Hillsdale, NJ*, 1986.
- [4] K. Vredenburg, J.-Y. Mao, P. W. Smith, and T. Carey, "A survey of user-centered design practice," in *SIGCHI conference on Human factors in computing systems*. ACM, 2002, pp. 471–478.
- [5] M. Angelini, N. Ferro, G. Santucci, and G. Silvello, "VIRTUE: A visual tool for information retrieval performance evaluation and failure analysis," *Journal of Visual Languages & Computing*, vol. 25, no. 4, pp. 394–413, 2014.
- [6] E. Bauleo, S. Carnevale, T. Catarci, S. Kimani, M. Leva, and M. Mecella, "Design, realization and user evaluation of the SmartVortex Visual Query System for accessing data streams in industrial engineering applications," *Journal of Visual Languages & Computing*, vol. 25, no. 5, pp. 577–601, 2014.
- [7] T. Kosar, M. Mernik, and J. Carver, "Program comprehension of domain-specific and general-purpose languages: comparison using a family of experiments," *Empirical Software Engineering*, vol. 17, no. 3, pp. 276–304, 2012.
- [8] A. Barišić, P. Monteiro, V. Amaral, M. Goulão, and M. Monteiro, "Patterns for evaluating usability of domain-specific languages," *Proc. 19th Conference on Pattern Languages of Programs (PLoP), SPLASH 2012*, October 2012.
- [9] M. Hatzigianni and K. Margetts, "I am very good at computers: young children's computer use and their computer self-esteem," *European Early Childhood Education Research Journal*, vol. 20, no. 1, pp. 3–20, 2012.
- [10] A. Bruckman and A. Bandlow, "Human-computer interaction for kids," in *The Human-Computer Interaction Handbook: Fundamentals, Evolving Technologies, and Emerging Applications*, J. Jacko and A. Sears, Eds. Lawrence Erlbaum and Associates, 2002.
- [11] G. Sim, B. Cassidy, and J. C. Read, "Understanding the fidelity effect when evaluating games with children," in *Proc. 12th International Conference on Interaction Design and Children*, ser. IDC '13. New York, NY, USA: ACM, 2013, pp. 193–200.
- [12] N. Borgers, E. De Leeuw, and J. Hox, "Children as respondents in survey research: Cognitive development and response quality," *Bulletin de methodologie Sociologique*, vol. 66, no. 1, pp. 60–75, 2000.
- [13] P. Markopoulos and M. Bekker, "On the assessment of usability testing methods for children," *Interacting with computers*, vol. 15, no. 2, pp. 227–243, 2003.
- [14] J. Johnson, "Children, robotics, and education," *Artificial Life and Robotics*, vol. 7, no. 1, pp. 16–21, 2003.
- [15] F. Klassner and B. Schafer, "Using the new lego mindstorms ev3 robotics platform in cs courses (abstract only)," in *Proceedings of the 45th ACM Technical Symposium on Computer Science Education*, ser. SIGCSE '14, 2014.
- [16] U. A. Qidwai, "A lamp-lego experience of motivating minority students to study engineering," *SIGCSE Bull.*, vol. 39, no. 4, Dec. 2007.
- [17] R. Koitz and W. Slany, "Empirical comparison of visual to hybrid formula manipulation in educational programming languages for teenagers," in *Proceedings of the 5th Workshop on Evaluation and Usability of Programming Languages and Tools*, ser. PLATEAU '14. ACM, 2014.
- [18] Lejos Community, "Lejos," (21/4/2017). [Online]. Available: <http://www.lejos.org/>
- [19] Pololu, "Pololu 3pi," (21/4/2017). [Online]. Available: <https://www.pololu.com/product/975/>
- [20] J. Jackson, "Microsoft robotics studio: A technical introduction," *Robotics & Automation Magazine, IEEE*, vol. 14, no. 4, pp. 82–87, 2007.
- [21] W. P. Dann, S. Cooper, and R. Pausch, *Learning to Program with Alice (w/CD ROM)*. Prentice Hall Press, 2011.
- [22] Google, "Blockly," (21/4/2017). [Online]. Available: <https://developers.google.com/blockly/>
- [23] "http://mindstorms.lego.com/en-us/software/default.aspx."
- [24] M. Resnick, J. Maloney, A. Monroy-Hernández, N. Rusk, E. Eastmond, K. Brennan, A. Millner, E. Rosenbaum, J. Silver, B. Silverman *et al.*, "Scratch: programming for all," *Communications of the ACM*, vol. 52, no. 11, pp. 60–67, 2009.
- [25] S. Papert, *Mindstorms: Children, computers, and powerful ideas*. Basic Books, Inc., 1980.
- [26] M. McNally, M. Goldweber, B. Fagin, and F. Klassner, "Do lego mindstorms robots have a future in CS education?" in *SIGCSE*, vol. 6, 2006, pp. 61–62.
- [27] E. Rosenbaum, E. Eastmond, and D. Mellis, "Empowering programmability for tangibles," 2010.
- [28] Lego, "Enchanting," (21/4/2017). [Online]. Available: <http://enchanting.robotclub.ab.ca/>
- [29] C. Kelleher and R. Pausch, "Lowering the barriers to programming: A taxonomy of programming environments and languages for novice programmers," *ACM Computing Surveys (CSUR)*, vol. 37, no. 2, pp. 83–137, 2005.
- [30] Arduino, "Arduino," (21/4/2017). [Online]. Available: <http://www.arduino.cc/>
- [31] X.-W. T. Liu, "An intuitive and flexible architecture for intelligent mobile robots," Ph.D. dissertation, The University of Manitoba, 2005.
- [32] BlindReviewAuthor, "VL/HCC VL Study 2017 Companion Site," (21/4/2017). [Online]. Available: <https://sites.google.com/view/vl-empiricalstudy/home>
- [33] R. Van Solingen, V. Basili, G. Caldiera, and H. D. Rombach, "Goal question metric (GQM) approach," *Encyclopedia of Software Engineering*, 2002.
- [34] G. Sim and M. Horton, "Investigating children's opinions of games: Fun toolkit vs. this or that," in *Proc. 11th International Conference on Interaction Design and Children*. ACM, 2012, pp. 70–77.
- [35] B. Kitchenham, L. Madeyski, P. Brereton, S. Charters, S. Gibbs, and A. Pohthong, "Robust statistical methods for empirical software engineering," *Empirical Software Engineering*, vol. in press, p. 54, 2016.