

This provisional PDF corresponds to the article as it appeared upon acceptance.

A copyedited and fully formatted version will be made available soon.

The final version may contain major or minor changes.

## The Performance Improvement-score algorithm applied to Endoscopic Stone Treatment step 1 protocol.

Domenico VENEZIANO, giulio PATRUNO, Michele TALSO, Theodore TOKAS, Silvia PROIETTI, Angelo PORRECA, Guido KAMPHUIS, Shekhar BIYANI, Esteban EMILIANI, Marcos CEPEDA DELGADO, Lopez MARIA DE MAR PEREZ, Roberto MIANO, Stefania FERRETTI, Nicola MACCHIONE, Panagiotis KALLIDONIS, Emanuele MONTANARI, Giovanni TRIPEPI, Achilles PLOUMIDIS, Giovanni CACCIAMANI, Estevão LIMA, Bhaskar K. SOMANI

*Minerva Urologica e Nefrologica* 2020 Aug 04

DOI: 10.23736/S0393-2249.20.03747-9

Article type: Original article

© 2020 EDIZIONI MINERVA MEDICA

Article first published online: August 4, 2020

Manuscript accepted: June 3, 2020

Manuscript revised: May 27, 2020

Manuscript received: December 18, 2019

Subscription: Information about subscribing to Minerva Medica journals is online at:

<http://www.minervamedica.it/en/how-to-order-journals.php>

Reprints and permissions: For information about reprints and permissions send an email to:

[journals.dept@minervamedica.it](mailto:journals.dept@minervamedica.it) - [journals2.dept@minervamedica.it](mailto:journals2.dept@minervamedica.it) - [journals6.dept@minervamedica.it](mailto:journals6.dept@minervamedica.it)

## Treatment step 1 protocol.

Veneziano Domenico<sup>\*1,2,3</sup>, Patruno Giulio<sup>4</sup>, Talso Michele<sup>5</sup>, Tokas Theodore<sup>6</sup>, Proietti Silvia<sup>7</sup>, Porreca Angelo<sup>8</sup>, Kamphuis Guido<sup>9</sup>, Biyani Shekhar<sup>10</sup>, Emiliani Esteban<sup>11</sup>, Cepeda Delgado Marcos<sup>12</sup>, Maria de Mar Perez Lopez<sup>13</sup>, Miano Roberto<sup>14</sup>, Ferretti Stefania<sup>15</sup>, Macchione Nicola<sup>16</sup>, Kallidonis Panagiotis<sup>17</sup>, Montanari Emanuele<sup>18</sup>, Tripepi Giovanni<sup>19</sup>, Ploumidis Achilles<sup>20</sup>, Cacciamani Giovanni<sup>21</sup>, Lima Estevao<sup>1,2</sup>, Somani Bhaskar<sup>22</sup>

1) Life and Health Sciences Research Institute (ICVS), School of Medicine, University of Minho, Braga, Portugal

2) ICVS/3B's, PT Government Associate Laboratory, Braga/Guimarães, Portugal

3) Dept of Urology and Kidney Transplant, Grande Ospedale Metropolitano, Reggio Calabria, Italy

4) Dept of Urology, AO San Giovanni Addolorata, Rome, Italy

5) Dept of Urology, ASST Vimercate, Italy

6) Department of Urology and Andrology, General Hospital, Hall in Tirol, Austria.

7) Dept of Urology, San Raffaele-Turro Hospital, Milan, Italy

8) Dept of Urology, Policlinico Abano Terme, Italy

9) Dept of Urology, AMC University Hospital, Amsterdam, The Netherlands

10) Department of Urology, St. James's University Hospital Leeds Teaching Hospitals NHS, Leeds

11) Dept of Urology, Fundació Puigvert, Barcelona, Spain

12) Dept of Urology, Rio Hortega University, Valladolid, Spain

13) Dept of Urology, Centro de Cirugía de Mínima Invasión Jesús Usón, Cáceres, Spain

14) Dept of Urology, Università Torvergata, Rome, Italy

15) Dept of Urology, AOU di Parma, Italy

16) Dept of Urology, Ospedale San Paolo, Milan, Italy

17) Dept of Urology, University of Patras, Greece

18) Dept of Urology, Policlinico Universitario, Milan, Italy

19) CNR IFC, U.O. of Reggio Calabria, Italy

20) Department of Urology, Athens Medical Centre, Athens, Greece

21) USC Institute of Urology & Catherine and Joseph Aresty Department of Urology, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA;

22) Department of Urology, University of Southampton, Southampton, UK

**Abstract count: 244**

**Word count: 2291**

**Corresponding author:**

Domenico Veneziano MD, FEBU

**Keywords:**

EST s1, Pi-score, Hands-on training, Assessment, Endourology

**Abbreviations:**

HoT: Hands on Training

ART in Flexible: Advanced Residents Training in Flexible

Pi: Performance improvement

Pi-score: Performance Improvement score

E-BLUS: European Basic Laparoscopic Urological Skills

AUA-BLUS: American Urological Association Basic Laparoscopic Urological Skills

EST-S1: Endoscopic Stone Treatment – step 1

URS: UreteroRenoScopy

**Abstract**

## BACKGROUND

Pi-score (Performance Improvement score) has been proven to be reliable to measure performance improvement during E-BLUS hands-on training sessions. Our study is aimed to adapt and test the score to EST s1 (Endoscopic Stone Treatment step 1) protocol, in consideration of its worldwide adoption for practical training.

## METHODS

The Pi-score algorithm considers time measurement and number of errors from two different repetitions (first and fifth) of the same training task and compares them to the relative task goals, to produce an objective score. Data were obtained from the first edition of 'ART in Flexible Course', during 4 courses in Barcelona and Milan. Collected data were independently analysed by the experts for Pi assessment. Their scores were compared for inter-rater reliability. The average scores from all tutors were then compared to the PI-score provided by our algorithm for each participant, in order to verify their statistical correlation. Kappa Statistics was used for comparison analysis.

## RESULTS

16 Hands-on Training expert tutors and 47 3rd year residents in Urology were involved. Concordance found between the 16 proctors' scores was the following: Task1=0.30 ("fair"); Task2=0.18 ("slight"); Task3=0.10 ("slight"); Task4=0.20, ("slight"). Concordance between Pi-score results and proctor average scores per-participant was the following: Task1=0.74 ("substantial"); Task2=0.71 ("substantial"); Task3=0.46 ("moderate"); Task4=0.49 ("moderate").

## CONCLUSION

Our exploratory study demonstrates that Pi-score can be effectively adapted to EST s1. Our algorithm successfully provided an objective score that equals the average performance improvement scores assigned by of a cohort of experts, in relation to a small amount of training attempts.

## BACKGROUND

Performing surgical training outside the operative room (OR) is considered fundamental in order to prepare for clinical practice [1]. The aim of simulated lab and model training is focused on surgical skills' details, in order to improve technical performances and to avoid using the patients as training platforms. Surgical training relies mostly on the curriculum, rather than on the simulator used, [2] which means that the tutor plays a critical role in it.

Hands-on Training (HoT) has become increasingly popular in urological congresses and events and allows trainees and urologists to move their first steps under the guidance of tutors. To analyse their improvement, tutors observe the trainees along their trials and measure their achievements. This process is not objective or standardized and relies on personal expertise and opinion. Measuring performance improvement can be very useful to develop more reliable training methodologies and optimize any educational procedure. Laparoscopy has represented a giant leap in clinical practice, as well as surgical training, allowing clear view to the surgeon, its assistants and observers. Learning basic laparoscopic skills has been codified in renown and standardized curricula, such as E-BLUS[3] and AUA-BLUS[4]. Similarly, an endoscopic stone treatment training curriculum, named Endoscopic Stone Treatment step 1 (EST-s1) has been proposed to teach and assess trainees on endoscopic procedures[5], [6]. This consists in 4 tasks: Task 1, Flexible Cystoscopy; Task 2, Rigid Cystoscopy; Task 3, Semi-rigid Ureteroscopy; Task 4, Flexible Ureterorenoscopy. In every teaching protocol, performance improvement of practical skills is evaluated by expert's opinion, and is potentially biased by tutor's teaching ability, standardization of instruments, tasks or goals. The judgement of a group of experts, even when based on objective measures such as number of errors and time to complete tasks, can greatly vary[7].

In 2018 a new algorithm for performance improvement in laparoscopic skills has been proposed[7]. The Pi-score (performance improvement score) provides an objective measurement of performance improvement over a small number of trials, by considering multiple variables and using established quality standards as a benchmark. It has been tested on the EBLUS tasks and provided an objective score that equals the average Pi assessment of a cohort of experts.

In this exploratory study we evaluated the efficacy of the PI-score algorithm applied to the endoscopic steps of EST-s1, to evaluate its reliability against the average score of a group of expert tutors.

## METHODS

Study of performance improvement was done on the four EST-s1 tasks: rigid cystoscopy, flexible cystoscopy, semi-rigid Ureterorenoscopy (URS), flexible URS. Rules and task goals were applied in accordance with what has already described in literature [5]. Five attempts were used, because this modality was applicable to the majority of HoTs and enough to identify an improvement.

#### *Pi-score algorithm*

The Pi-score [7] has been designed to create an objective and fast correlation between different results within a HoT session. The Pi-score considers and measures the two variables collected during each EST-s1 task: number of errors and task completion time. Ideal timing and errors are defined following the established EST-s1 protocol[5]. Following the EBLUS Pi-score study[7], dedicated cut-offs have been set to highlight different levels of performance improvement on each EST-s1 task. Cut-offs were pre-determined and then trimmed following experts' opinion. A participant could therefore achieve a low (1), good (2), excellent (3) or outstanding (4) improvement. The cut-offs defined for each task are summarized in Table 1.

#### *Data collection*

In 2018, we collected data from 4 hands-on training courses, the ART-in-Flexible Course 1<sup>st</sup> edition, in Milan and Barcelona. Each course lasted 2 days and was delivered by the same tutor (D.V.), using the same series of tips and tricks aimed to optimize performance on the EST-s1 training tasks. The setup adopted during the courses included standard EST s1 simulators: the endoscopic simulator by Cook Medical (Bloomington, IN, USA) and the Coloplast K-Box (Humlebaek, DK). Both simulators were modified to meet the EST-s1 requirements. Two participants were allowed in each training station, with time availability enough to complete 5 runs of each task. The first trial was performed as a baseline check with time and error collection, before receiving any practical hint or guidance by the tutor. Afterwards, the tutor provided the participants standardized practical suggestions, to enhance and improve their performance. Run 1 and 5 were considered for this study. Time and errors from trial number 5 were collected as well, to monitor performance improvement after the learning session.

A score-sheet (Table 2) was used to collect the time taken for completion of all tasks, together with the number of errors along the course. Participants demographics (age, gender) and previous experience with endoscopic procedures or technological devices was also collected with a Likert scale score ranged from 1 to 5 (1- very low, 5-very high). The tutor personally oversaw this phase as well, to minimise biases and ensure complete data collection.



A control arm of 16 experienced tutors was arranged to study the reliability of the PI score. Just those who had at least four years of past experience as official HoT proctor on behalf of the European Association of Urology (EAU) were eligible. Each proctor was provided with an anonymized database, containing information (time + errors) from trials 1 and 5 for each participant enrolled. Based on the provided data, the proctor was asked to provide subjective evaluation of performance-improvement for each participant. Their evaluation had to be formulated on a scale from 1 to 4: 1-low improvement, 2-good improvement, 3-excellent improvement, 4-outstanding improvement. Inter-rater reliability was performed between the scores collected, to understand how much evaluation of performance improvement may change, even within a cohort of proctors is composed by experts. Afterwards, the average between the scores collected from all tutors for each participant was calculated. The resulting average scores were then compared to those provided by our PI-score calculator, in order to verify their statistical correlation and finally analyse the reliability of the algorithm. As a secondary goal, this study provided information about the efficacy of EST s1 as a teaching tool.

### *Statistical Analysis*

Data were summarized as mean and standard deviation (SD) or as absolute number, as appropriate. Cohen's Kappa statistics (for two raters and a binary response – i.e. with no degree of disagreement), Weighted Kappa statistics (which considers the degree of disagreement between two raters) or Kappa Statistics for multiple raters [8] were used to measure reliability among raters, as appropriate. One rater-removed analysis was applied to assess the influence by each single rater. The degree of the agreement analysed by Kappa Statistics was scored according to the following scale: 0 – no agreement; 0-0.2 – slight agreement; 0.2- 0.4 – fair agreement; 0.4-0.6 – moderate agreement; 0.6- 0.8 – substantial agreement; 0.8-1.0 – almost perfect agreement. A P-Value <0.05 was considered statistically significant. All statistical analyses were done by SPSS for Windows (version 22), IBM, Chicago, Illinois, USA or by STATA for Windows (version 13), Lakeway Drive, College Station, USA.

## **RESULTS**

Forty-seven 3<sup>rd</sup> year residents in Urology were enrolled in the study from 4 courses. Mean age of the participants was 29 years ( $\pm 1,2$ ) with a male:female ratio of 32:15. Personal technological expertise was scored with an average of 2,8 ( $\pm 0,9$ ) out of 5, while personal endoscopic expertise before the beginning of the course was on

average 2,5 ( $\pm 0,8$ ) out of 5.

On Flexible Cystoscopy, inter-rater agreement was “fair” (Kappa=0.30) with Kappa ranging from 0.21 to 0.35 with one rater-removed analysis (P=0,00001). Regarding Rigid Cystoscopy, performance improvement agreement among the 16 proctors was “slight” (Kappa=0.18) with Kappa ranging from 0.17 to 0.24 with one rater-removed analysis (P=0,00001). For Semi-rigid URS, expert agreement was “slight” (Kappa=0.10) with Kappa ranging from 0.09 to 0.11 with one rater-removed analysis (P=0,00001). For the final task, Flexible URS, the agreement was “slight” (Kappa=0.19) with Kappa ranging from 0.11 to 0.23 with one rater-removed analysis (P=0,00001).

After comparing the scores provided from the 16 experts, the average score for each participant and for each task was compared to the PI-scores produced by the algorithm, with the following results.

On task 1 the agreement between the expert average scores and the PI-scores was “substantial” (Kappa=0.74, 95% CI: 0.58 - 0.89) (P=0,00001). On task 2 the agreement between the expert average scores and the PI-scores was “substantial” (Kappa=0.71, 95% CI: 0.56 - 0.86) (P=0,00001). On task 3 the agreement between the expert average scores and the PI-scores was “moderate” (Kappa=0.46, 95% CI: 0.27 – 0,65) (P=0,00001). On task 4 the agreement between the expert average scores and the PI-scores was “moderate” (Kappa=0.49, 95% CI: 0,3 – 0,68) (P=0,00001). Agreement data are summarized in Figure 1. Overall average Pi-score for all participants was 44,1 ( $\pm 9,2$ ). No relevant correlation was found between the PI-score and the age of participants. No relevant correlation was found between the reported technological or endoscopic background and the PI-scores obtained.

## DISCUSSION

The Pi-score has been developed to objectively analyse performance improvement, optimize training methodologies and eventually spot those who demonstrate specific surgical talent. In the approach of a new skill, the improvement is higher in the beginning, while it becomes slower as the learning curve approaches a plateau [8]. For this reason, the original PI score was tested on the very first series of training trials [7]. The previous study already highlighted that expert opinion in this early learning phase is not reliable, with low agreement between the different raters. This happens especially because it's not usually feasible for a tutor to personally follow several trainees at the same time. This evidence was confirmed by the present study, which shows even lower agreement between the different tutors, regarding the performance improvement of the participants enrolled. This is

probably due to the novelty of the EST s1 protocol, with new tips and tricks and learning curves, as opposed to the already renowned E-BLUS [3], [11]. Another aspect to be considered is that assessing a single task trial is different than assessing the performance improvement of an individual over a series of trials, as in this case pre-established metrics are missing and scoring is still highly subjective. This justifies why the highest agreement between tutors was described as “fair” ( $Kappa=0.30$ ) on task 1, rigid cystoscopy. This issue was overcome by comparing the PI-score results to the average scores of the tutors instead of their individual ones. The PI score applied to EST-s1 provided results that were matching experts’ average opinion, with concordance measured as “substantial” on task 1 and task 2 and “moderate” on task 3 and 4. For this reason the algorithm was used, even if still under testing, to select the participants who demonstrated more progression and who were more suitable to proceed to advanced courses.

Moreover, analysing the Pi-score data from course to course, allowed the faculty to improve and refine the teaching methodology and the hands-on training session duration time, which is core to the success of a training event. The overall average Pi-score registered between the 47 participants was 44,1 ( $\pm 9,2$ ), which corresponds to an “average performance improvement” (table 1), with several pikes rated as “excellent improvement” (11/47 on task 1, 6/47 on task 2, 4/47 on task 3 and 17/47 on task 4), which added further validity evidence to EST s1 as a teaching tool (secondary goal of our study).

The PI-score has proven its efficacy on different fields of urology and could be tested on different procedures and medical specialities or, in general, on any task that includes measurable practice. Pi-score might be also used to reinforce the assessment of a simulator or a curriculum as a pure teaching tool.

#### *Limitations of the study*

Our study confirms the findings of the previous one, but shares some of its limitations: the size of the sample used, which is intrinsic to the type of events selected. Allowing each participant to perform the same tasks several times requires time and simulation devices. Future studies could be conducted in training hospitals, equipped with the required simulation devices, allowing multi-centric data collection from more participants. Another limitation was related to the poorly standardized teaching methodology of EST s1, which didn’t allow a high concordance between tutor judgement. Increasing data collection and tutor experience could also help to trim the quality cut-offs of the algorithm, making it more reliable for automated performance assessment.

#### *Future applications*



The Pi-score can be used to analyse several educational aspects. While being designed to objectively assess performance improvement of an individual, it can be applied to the analysis of a whole course, by calculating the average scores of all participants. Overall course results can be afterwards used to judge the course itself, in relation to duration, contents or tutor teaching skills, depending on the variable considered. Pi-score could be also applied to individual selection, based on personal learning skills, especially when considering big data volume. Pi-score might be lastly used to be integrated in a machine learning process to provide automatic tailored teaching on a large scale.

To ensure full replicability of our study, Pi-score is available for free use on the webpage [www.domenicoveneziano.it/piscore](http://www.domenicoveneziano.it/piscore)

## CONCLUSIONS

Our exploratory study demonstrated that subjective evaluation of performance improvement is extremely variable and needs therefore to be aided by dedicated tools. The use of algorithms could help to objectively measure performance improvement. The PI-score, initially developed for basic laparoscopic skills, our preliminary evaluation of the algorithm provided an objective score for EST-s1 that can equal the average performance improvement scores assigned by of a cohort of experts, in relation to a small amount of training attempts.

## REFERENCES

- [1] B. A. Parsons, N. S. Blencowe, A. D. Hollowood, and J. R. Grant, "Surgical training: The impact of changes in curriculum and experience," *J. Surg. Educ.*, 2011.
- [2] R. Satava, "The future of surgical simulation and surgical robotics.," *Bull. Am. Coll. Surg.*, vol. 92, no. 3, pp. 13–9, 2007.
- [3] W. M. Brinkman *et al.*, "Results of the European Basic Laparoscopic Urological Skills Examination," *Eur. Urol.*, vol. 65, no. 2, pp. 490–496, Feb. 2014.
- [4] T. M. Kowalewski *et al.*, "Validation of the AUA BLUS Tasks," *J. Urol.*, 2016.
- [5] D. Veneziano *et al.*, "Validation of the endoscopic stone treatment step 1 ( EST - s1 ): a novel EAU training and assessment tool for basic endoscopic stone treatment skills — a collaborative work by ESU , ESUT and EULIS," *World J. Urol.*, vol. 1, no. 0123456789, 2019.
- [6] D. Veneziano *et al.*, "Evolution and Uptake of the Endoscopic Stone Treatment Step 1 (EST-s1) Protocol: Establishment, Validation, and Assessment in a Collaboration by the European School of Urology and the Uro-Technology and Urolithiasis Sections," *Eur. Urol.*, vol. 74, no. 3, 2018.

- [7] D. Veneziano *et al.*, "Performance Improvement (Pi) score: an algorithm to score Pi objectively during E-BLUS hands-on training sessions. A European Association of Urology, Section of Uro-Technology (ESUT) project," *BJU Int.*, 2019.
- [8] L. S. Feldman, J. Cao, A. Andalib, S. Fraser, and G. M. Fried, "A method to characterize the learning curve for performance of a fundamental laparoscopic simulator task: Defining 'learning plateau' and 'learning rate,'" *Surgery*, 2009.
- [9] D. Carrion *et al.*, "Current status of urology surgical training in Europe: an ESRU-ESU-ESUT collaborative study.," *World J Urol*, vol. Apr 13, no. 10.1007/s00345-019-02763-1., 2019.
- [10] A. Volpe *et al.*, "Pilot Validation Study of the European Association of Urology Robotic Training Curriculum," *Eur. Urol.*, 2015.
- [11] B. K. Somani *et al.*, "The European Urology Residents Education Programme Hands-on Training Format: 4 Years of Hands-on Training Improvements from the European School of Urology," *Eur. Urol. Focus*, 2018.

Improvement	Flex cyst	Rigid cyst	Semi-rigid URS	Flex URS
Low	0 – 25	0 – 30	0 – 31	0 – 31
Average	25 – 51	30 – 57	31 – 62	31 – 55
Excellent	51 – 85	57 – 90	62 – 95	55 – 81
Outstanding	>85	>90	>95	>81

Table 1: EST s1 Pi-score cut-offs

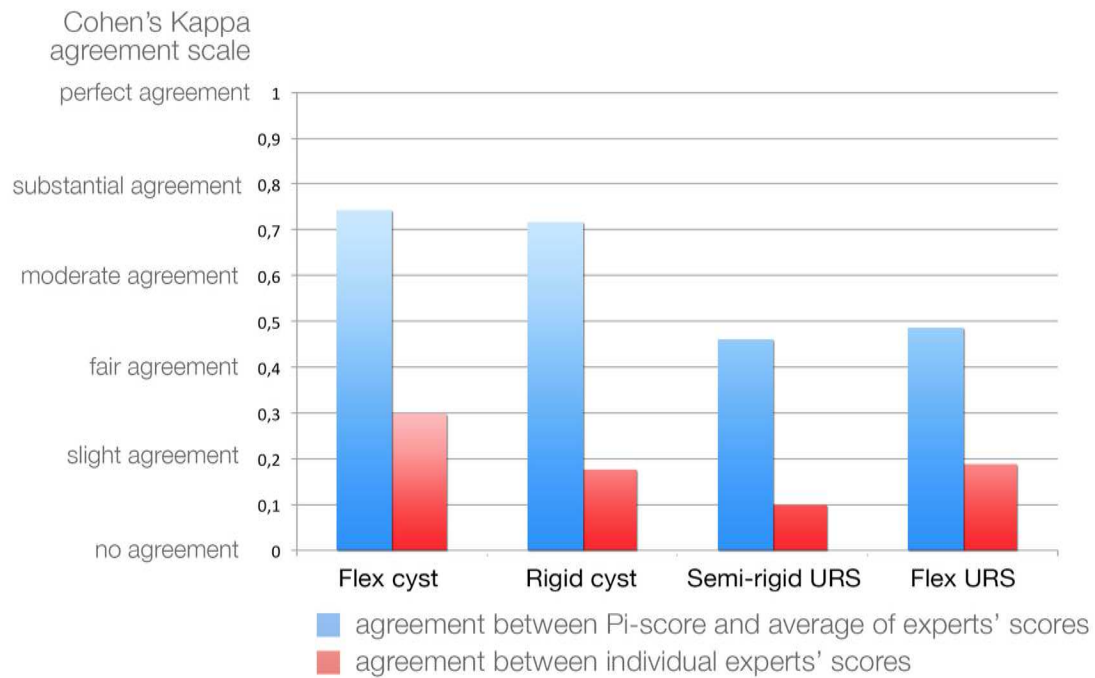
Name Surname: \_\_\_\_\_ Age: \_\_\_\_\_ PGY: \_\_\_\_\_  
 Date: \_\_\_\_\_

Endourology expertise (1 poor – 5 very high): 1 2 3 4 5  
 General technology expertise (1 poor – 5 very high): 1 2 3 4 5  
 Have you ever taken part to an endourology course? YES NO

Please use the template to take note of task time and number of errors during the course. The data will be used for research and quality purposes.

	Flex cystoscopy Ref – 0:36						Rigid Cystoscopy Ref – 0:41						Semi-Rigid URS Ref – 1:07						Flexible URS Ref – 3:43									
1st trial	min		sec				min		sec				min		sec				min		sec							
N° errors	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6
2nd trial	min		sec				min		sec				min		sec				min		sec							
N° errors	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6
3rd trial	min		sec				min		sec				min		sec				min		sec							
N° errors	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6
4th trial	min		sec				min		sec				min		sec				min		sec							
N° errors	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6
5th trial	min		sec				min		sec				min		sec				min		sec							
N° errors	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6	0	1	2	3	4	5	6

Table 2: Scoring sheet used along the courses.



1