

Instrumentation of Learning Situation using Automated Speech Transcription: A Prototyping Approach

Vincent Bettenfeld¹, Salima Mdhaffar², Christophe Choquet¹ and Claudine Piau-Toffolon¹

¹LIUM-IEIAH, Le Mans Université, Avenue Olivier Messiaen, 72085 LE MANS CEDEX 9, France

²

Keywords: Learning Web Environment, Transcription Aided Learning, User Needs Assessment.

Abstract: This paper presents the ongoing conception of a set of tools, based on live transcription of speech during lectures and designed to instrument traditional lectures as well as web conferences or hybrid learning situations. The toolset exploits speech and interactions taking place during courses, keeps track of them and facilitates their reuse both in students' studies and in future iterations of the course delivered by the teacher. Its goal is to help students stay focused on the teacher's explanations and offer them greater possibilities of interactions. The prototype was conceived with an approach based on the analysis of communicational and informational needs of the end users, especially in regard to the instrumentation possibilities offered by the innovative technologies considered in the project. In this paper, we detail the different tools produced in order to offer synchronous and asynchronous support to the learning activity. We describe a real-life test as well as changes brought to the device afterwards, and finally we describe the first experiment conducted with the device.

1 INTRODUCTION

The PASTEL project, standing for Performing Automated Speech Transcription for Enhancing Learning, is a research project driven by LIUM, LS2N, CREN and Orange Labs. The project's goal is to explore the potential of synchronous speech transcription and application in specific teaching situations. This technology allows to generate a textual version of the teacher's or the students' speech, and immediately use it in order to help them in their teaching or learning activities. Other interests of the project include an editorial toolset for teachers to save, edit and reuse the material produced or gathered in class. This edited material can be exploited during future sessions of the course, or become the support of another course in a different format such as an online class.

The project focuses on pedagogical situations such as lectures and group work, mainly in higher education. These situations can take place either in the classroom, on online platforms using videoconferencing systems, or in hybrid classes. These situations are complex and need to be instrumented with flexible and adaptive tools, offering various configurations for diverse actors. Using transcription technologies, we aim to create such tools and adequately

instrument these situations.

The speech transcription allows human actors to access the textual version of a sentence a few seconds after it was pronounced, and browse the whole text as they wish. This tool can help students solving comprehension problems caused by hearing and language barrier (Ryba et al., 2006), or allows them to use down time to read again a more complex section. Having access to the recording of a lecture decreases students' stress: they express trust in these recordings, which leads them to take fewer notes and concentrate on the teacher's explanations (Heilesen, 2010). To our knowledge, usage of synchronous transcriptions is very limited in pedagogical situations, particularly in higher education in which the vocabulary used can be very specific to the taught domain. As part of the project, other technologies, such as real-time material recommendation and thematic segmentation, will be additionally tested.

Our goal is to study the usability of these tools, and overall the usability of synchronous speech-to-text transcription in class context. In the next part, we present a state of the art of research in technology enhanced learning (TEL) and automatic speech transcription (AST) motivating the project and our iterative design methodology. Moodle plugins and their

interfaces were developed in an iterative way, this toolset is experimented in real conditions and reengineered in a new iterative session. The first and second iterations of the project are then described and discussed in the following parts before giving some concluding remarks related to future work to be done.

2 PROJECT MOTIVATION AND DESIGN METHODOLOGY

2.1 Project Motivation

2.1.1 Research in Technology-enhanced Learning

The rise of hybrid classes and video conferences provided institutions with recording and sharing tools both software and hardware (microphones, cameras, clickers). With this project, we want to take advantage of these devices. The goal is to offer tools simple enough to be compatible with the teaching activity, and efficient enough to be adopted. The methodology being design-based research (Wang and Hanafin, 2005), these tools are based on existing needs, tested in authentic context and improved through iterative cycles. Both researchers and teachers are engaged in this procedure.

As Internet usage is widespread, a great number of teachers share their pedagogical material with students online. Massive Open Online Courses (MOOCs) are now common, and offer a great flexibility for learners. Another type of online courses is Small Private Online Courses (SPOCs). They exploit MOOCs' strengths, i.e. their tools and flexibility, with a smaller number of learners. The edition of lecture videos and their segmentation is a tedious task, which existing works tackle automatically using slides (Ngo et al., 2003)(Uke and Thool, 2012). This task can be intimidating for teachers, especially if they must consider other resources than slides, such as the textual version of lessons and relevant external documents. Designing a software tool to collect, store and exploit this material directly during traditional classroom sessions is one of the motivations of the project.

In this project, we intent to explore how the product of real-time speech transcription can help learners and teachers. Existing works have proposed the transcription of lectures in real-time (Iglesias et al., 2016), and the originality of our proposition is providing new documents and interactions in real time. The satisfaction of user needs is researched, particularly in terms of information. These needs could be

satisfied by content derived from the transcription, or by using this transcription as a support of communication. Finally, this project is the opportunity to study how to index and browsing a great quantity of data growing and formatted in real-time.

2.1.2 Research in Speech Transcription

Automatic Speech Recognition (ASR) aims to automatically convert speech from a recorded audio signal into a text. Different researches in the literature have demonstrated the advantages of synchronous speech transcription for online courses (Ho et al., 2005)(Hwang et al., 2012). For example, Ho and al. (Ho et al., 2005) argued that synchronous speech transcription helps non-native English students to better understand lectures that are delivered in English. (Shadiev et al., 2014) mentioned that synchronous speech transcription can help (1) students with cognitive or physical disabilities who need additional support, have difficulties in reading, writing, or spelling by giving them a note-taking for all the teacher's speech (2) online students if the quality of audio is not good during communications, (3) non-native speakers, (4) students in traditional learning environment by allowing the educators to take a proactive, rather than a reactive approach to teach students with different learning styles and (5) students in collaborative academic activities.

Over the past decades of research, the field of automatic speech recognition has made considerable progress and many algorithms and techniques have been proposed and identified. However, a limit in an ASR is its inability to recognize unknown words which do not appear in the training data of the language model. A language model is one of the most important components in an ASR which assigns a probability to a sequence of words. The assigned probability guides the ASR to choose which sequences of words are possible to recognize. For the language model, the training data is in the form of word sequences from real speech, which were manually or automatically transcribed. It can also be from written text (e.g. web documents, news articles, books). This data is very hard to acquire and very expensive.

The aim of this work is to transcribe lecture speeches, which are very different from other types of audio content (e.g. news audio, documentary). Educational domain is characterized by the presence of different courses which are also given by different teachers. It is very hard to build a training data that allows a relevant use for all courses' topics. Since manually generating transcriptions is a challenging task, the classical idea is to use an existent language model

and to adapt this language model with domain's data. Domain's data means data related to the addressed task. So, the first motivation in regard to ASR is to adapt the language model by supplementing the training data with some additional data from the domain. The main difficulties are: (1) which data to collect, (2) what is the source of data and (3) which queries should be used to collect these data.

In order to automatically enrich the course with relevant external documents, a segmentation step must be realized. This segmentation must be a thematic segmentation to enrich each part. The goal of thematic segmentation is to detect the borders of homogeneous zones at the level of the content. It will be necessary to characterize borders in order to link each thematic zone with pedagogical documents available in an external knowledge base. The first difficulty is to define the theme concept's notion in a transcription that is monothematic (the main object of the course). The second difficulty will come from the real-time aspect of the thematic segmentation which, to our knowledge, has never been discussed so far.

2.2 Design

2.2.1 Considered Pedagogical Situations

Though diverse pedagogical situations will be considered in the project, current development has focused on lecture situations, either in presence of students or in hybrid configuration. They have been considered first because the interactions in these contexts are more simple. Indeed, the teacher interacts with his audience as a whole and student-to-student interactions are very limited. Hybrid situations grant the possibility for remote students to attend to the lecture by watching the teacher through their interface. However, the teacher cannot visually evaluate if their learning activity takes place as planned. By considering these situations, a way to convey the information needed by the teacher must be found.

2.2.2 Design Process

The initial phase of the project was a study of practices among students and teachers in higher education. This first study allowed us to determinate which functionalities were pertinent to offer in this context, and guided the development of different tools based on the transcription service.

Our toolset is conceived by an iterative process, of which two cycles are described in this article. Our goal is to experiment this toolset in real conditions with actual end users, in order to explore the most relevant use for each technology in the considered con-

texts. This cyclic methodology lets us quickly assess users' behavior and feedback, as well as operate a fast re-engineering of the interfaces. In addition of lectures, the project will later focus on group work sessions, and post-session reuse of the generated content.

3 FIRST ITERATION

3.1 User Needs Assessment

The first step of the project consisted in studying existing practices before any instrumentation, as well as possible uses for tools planned to be developed in the project. A study based on a twofold approach, both quantitative and qualitative, was conducted (Crétin-Pirolli et al., 2017). First, a questionnaire was built and distributed to 94 students engaged in a computer science master's degree. Sixty-two questionnaires were collected and exploited. Meanwhile, a semi-directive interview was conducted with an associate professor giving lectures as part of the computer science master's program at Le Mans Université.

In regard to lectures, among students encountering difficulties, two thirds (66%) reported experiencing downtime at a point during class, as well as difficulties in understanding the class concepts (reported by 61% of them). Only a minority is using external resources given that 75% of students do not search for additional material such as definitions, texts or graphics on the Internet during class.

Considering the course instrumentation based on transcription, students are mainly interested by the possibility to communicate to the teacher which points were not understood. A majority of students also estimate the synchronous transcription is useful, in addition to the possibility of reading and exploiting the text after class.

For his part, the professor assumes the access to the transcribed speech is a good thing, albeit students should still be taking notes. He considers that note-taking is a part of the learning process. He does not want students to be distracted from their learning activity, so the resource recommendation system use should be occasional and brief, in order to prevent cognitive overload.

Needs and a priori acceptance lead to think that live transcription availability is relevant. However, synchronous availability of the transcription was not a central expectation in the eyes of students experiencing no hearing difficulties, or language barriers. Other tools based on this live transcription were relevant in regards of needs and existing practice, but exploitation of the transcription itself was considered more

pertinent in asynchronous time. After a synthesis of the various needs and constraints, relevant functionalities which could be integrated using the technologies of interest of the project were listed. They are summarized in the following section.

3.2 User-tool Interaction

The toolset is conceived to instrument usual classrooms and hybrid classes. It is designed to support learning all through the module duration. A set of functionalities has been selected, regardless of topic taught, but in regards of the needs assessment. They were chosen to foster learning during classes and during personal review of lessons.

3.2.1 During Classes

Students. Our goal was to allow remote students to access the same material as users present in class, mainly the teacher's explanation and a view of the slideshow projected. Despite their remote location, their possibilities of interactions with the teacher should be similar to the ones offered to the students present in class. All students would profit to be recommended resources which can be consulted during downtime, or saved for later use. This extra material illustrates the teacher's speech and is suggested at the right time to prevent students from having to navigate in a list of resources while they are already engaged in listening activities.

Teachers. The teachers need indicators (e.g. a pedagogical significant variable extracted from, or elaborated by the help of data automatically collected during the session) to compensate for the facts that some or all students are out of sight. Since teaching is also a very engaging activity, they need very synthetic aid to make decisions about the progression of the lecture in real-time. To engage and be able to get feedback from their audience, the tool must let them propose interactions for large auditoriums, in the fashion of clickers. Closed-ended questions are relevant, but teachers will also need the possibility to gather open-ended feedback, either concerning the lesson topic or the practical situation (e.g. "I see some students have trouble connecting to the class, what seems to be the issue?").

3.2.2 Between Classes

Students. After the lesson has ended, students need to access to the transcription, the slides and the resources offered to attendees. If they need to spend more time to study the material, they need the possibility to replay the video recorded during class. In this

case they should be able to associate each moment of the video with the particular slide or resources proposed at this moment.

Teachers. Teachers need access to the interactions history to take into account the difficulty level of their class and the need for examples/details. An overview of frequently asked question would give the opportunity to answer frequent questions beforehand and generally taking remarks into account. They have to be able to retrieve and potentially export the recorded video, either as a whole or partially. Ideally, they would be able to divide the class recording to produce a series of video clips, each one detailing a theme. These clips could be suggested as complement to class, as learning material in a SPOC or in a context of flipped classroom.

3.3 Plugin Description

The toolset was built as a Moodle plugin, and can thus be accessed as a web page. This plugin use technologies of speech-to-text transcription and resource recommendation, each hosted on a separate server managed by the research team who developed it (see Fig. 1).

Speech-to-text Technology. The teacher is equipped with a lapel microphone (as well as filmed). This microphone sends the sound data to our server which organizes the coordinates streams. Sound is transcribed, and sentences are sent to the students as soon as the system evaluates that they were transcribed with enough confidence in regard to the system. This validation and upload takes less than ten seconds, which coincide with the video streaming latency.

Recommendation Technology. The resources recommendation system takes as parameter the product of the transcription. In real time, this system recommends documents dealing with similar concepts. A set of interesting concepts can be extracted from the analysis of a sentence pronounced by the teacher, or a student's question. The system then fetches labelled resources either freely on the Internet or in a limited pool moderated by humans.

PASTEL as a Moodle Plugin. Moodle is an open-source learning management system on which teachers can create online course and enroll their students. The toolset was developed as a Moodle plugin so that it can be integrated to any existing course platform, next to existent material. Students enrolled in a class can access our tool as they would access other activities or resources, and teachers can create a virtual classroom in the same way they would add a page to their course.

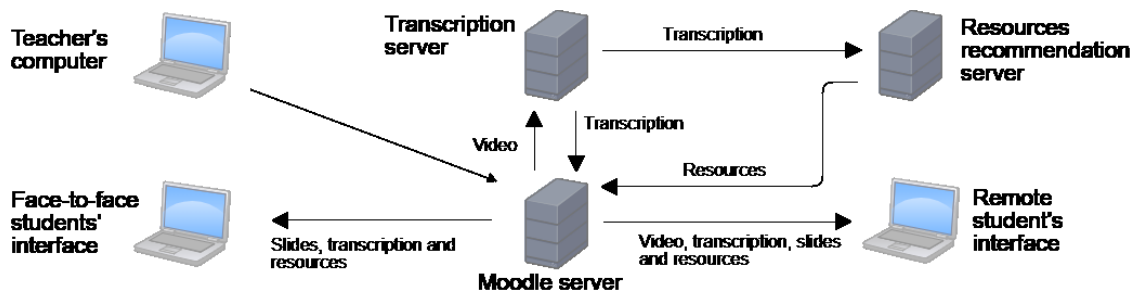


Figure 1: Data streams exchanged between the different modules.

3.4 First Interface Version Design

In this section we present the different components implemented in the initial stage of the prototype. This first interface was elaborated by the research team in regards of detected needs, and implemented in order to gather user feedback (Bettenfeld et al., 2018). Given the varying needs and scenarios, students and teachers have access to their own version of the toolset.

Components of the Student Window. The first version of the plugin provided two modes of content browsing: synchronous and asynchronous. In synchronous mode, slide display was coordinated in real-time with the projection in the classroom. In asynchronous mode, even though data was received and stored, the display kept focused on a given slide or point of the transcription, allowing the student to spend a longer time reading. On the left part of the screen, students could watch the teacher's video stream in real-time. On the bottom left, the lecture slides were displayed and could be zoomed on. Underneath, they could submit their questions or their needs in a text field. The text was then sent to the teacher, and also analyzed by the resource recommendation system. The transcription display area was located in the center of the screen and was updated synchronously. A note-taking panel was associated to each paragraph of the transcription. Besides writing down their notes, students could notify a need for further information to the teacher. The material recommendation system also took in consideration this alert, and analyzed the concepts being explained or cited in the paragraph to provide relevant information. This system offered a set of links on the right part of the screen, evolving in real time, redirecting to external resources. At the top of the screen, a wide timeline of the lecture was displayed. It was a navigation tool designed to facilitate browsing slides and the transcription window quickly.

Components of the Teacher Window. This window is offered to be used while giving a lecture. As this task requires time and focus, functionalities provided are simple and quick to use, and the display components are kept concise. A portion of the screen gave visual feedback of the camera, displaying the same video stream as the students received (see A on Fig. 2). In addition to this page designed for the teacher, a second page displayed the slides destined to be projected to the audience (opened clicking F on Fig. 2). The teacher could navigate the slideshow using the mouse or the left and right arrow keys. The right part of the screen was dedicated to a text feed showing open-ended questions asked by the audience in real time (see E on Fig. 2). At the bottom of the screen was presented a set of indicators, two of them shown as bar graphs (C on Fig. 2). The first one displayed the proportion of alerts sent by students estimating the lecture's pace was too quick. The second bar displayed the proportion of students looking at material corresponding to a past moment. The last indicator was a table (D on Fig. 2) detailing the three slides on which the greater number of student expressed a difficulty. The resources recommended to students were displayed in real time, in a table (B on Fig. 2). The display was more concise as the teacher did not need to rate, or get a preview of the resources. Only the titles appeared, allowing the teacher to recommend a one particular resource to students, or to showcase it on the projected view.

3.5 User Test

The prototype was tested with students and a teacher in actual class context in order to evaluate usability, to list the benefits of using the system, to be able to analyze usage, and consequently to improve tools supporting the least satisfying tasks. One of the main objectives is to check that the amount of information provided to testers was not a source of cognitive overload. Providing text in real time during a

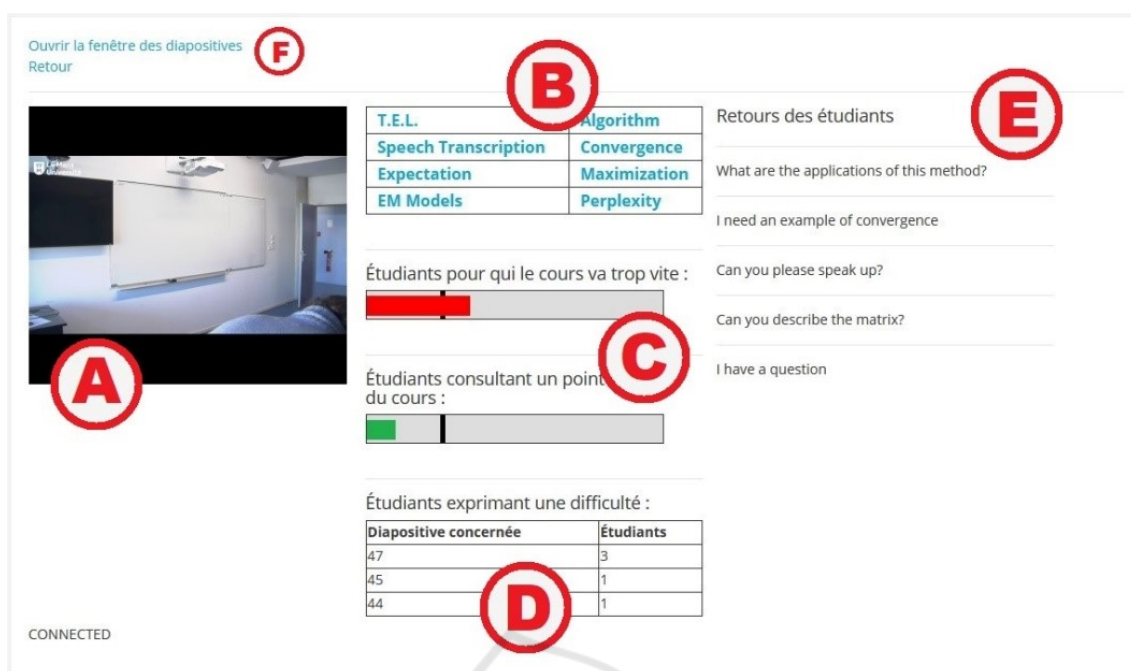


Figure 2: Example of the teacher’s interface: (a) Webcam feedback (b) Recommended resources (c) Indicators (d) Detailed difficulty indicator (e) Open-ended question feed (f) Link opening the slideshow and link closing this window.

learning activity can trigger such overload, but conversely saving lecture content for later use can relieve students’ memory, as it was studied with podcasts (Traphagan et al., 2010). The user test took place December 12, 2017 in hybrid configuration, both on Le Mans campus and on Nantes campus. The teacher was situated on the former location accompanied by three students while the remaining 7 students were located on the other campus. In both places, participants were filmed, their screen’s activity was recorded and they agreed to take part in a post-experiment interview. Students’ behavior during the lecture was mainly concentrated on the teacher and slides, and their use of the transcription and resources was more occasional. To navigate between different slides, students mostly used the previous and next slide display buttons. Their use of the whole timeline to browse content was limited. During the interview, they expressed that manipulating the timeline requires additional thinking, albeit short, they could not invest time and concentration in. The students also made explicit their frustration about the substance of recommended material. The content was not adapted to their master’s degree level: a lot of suggested information was considered trivial at this point of the cursus, such as definitions of basic concepts used during the lecture. No major problem was identified regarding the teacher interface.

4 SECOND ITERATION

4.1 Modifications of the Interface

Slideshow as Core Material. In the first version of the prototype, an important section of the screen was dedicated to the transcription, and slides were displayed in a smaller section. The relative size did not coincide with the importance of each component’s use. Students primarily zoomed on the slides and occasionally downsized them to peek at the transcription. The size of each of these component has been switched in the interface offered by default to better coincide with the students’ use (see Fig. 3).

Interface Flexibility. To prevent the students from hiding every tool whenever they need to zoom on the slide displayed, a system of panels was implemented. Using this system, students are able to select the tools they wish to keep on screen. The slideshow is displayed as big as possible depending on the remaining space. Ultimately it can be displayed at full size if every tool is hidden, for maximum readability. Students can hide and show every tool panel at any moment according to their needs. To support the use case of students reading the transcription comfortably, it is possible to expand the transcription area.

Navigation in Asynchronously Generated Content. As the previous way of browsing the previous

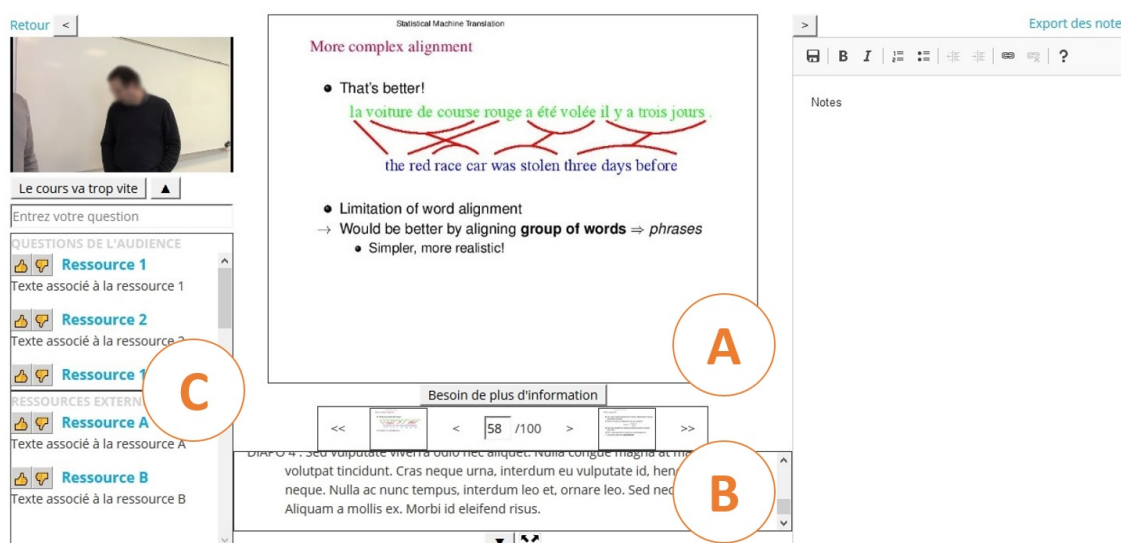


Figure 3: New student interface, showing the slideshow (A), transcription (B), and resources (C).

parts of the translation and slideshow necessitated too much concentration, the system has been simplified. Instead of a timeline representing the entirety of the class time, users have access to the previous and the next slide of the slideshow, with preview of both of them to facilitate recall. The possibility of entering the number of the desired slide has also been implemented, allowing users to quickly access a slide referenced in the teacher's speech or their notes. If the transcription window is extended (as shown on Fig. 5), clicking on a paragraph updates the slide view to display the slide projected during this part of the speech.

4.2 Real Conditions Experiment

4.2.1 Protocol

The second version of the prototype was tested on March 22nd 2018, at Laval during an information and communication lecture. Beforehand, the teacher and students received a presentation of the toolset in order to facilitate later use. During the experiment, 13 students were located in the same room as the teacher and 16 others were located in a remote classroom. During the one-hour lecture, students were given access to the platform and their activity was monitored. Afterwards, the teacher and students were interviewed.

4.2.2 Results

Students were satisfied with the quality of the translation. They reported a frustration in regards of real-time communications as their open-ended questions

were not correctly transmitted to the teacher due to network limitations, and the teacher's pedagogical scenario did not include time dedicated to review the audience's questions.

Students were confused by the number of tools available simultaneously. They are themselves aware that a comprehensive interface could be useful to users familiar with the system. Yet, as users still discovering the toolset, they are not comfortable enough to both concentrate on the lecture and use every functionality offered. They still expressed the need for more flexibility, particularly the possibility of changing the position of tools on screen. An analysis of 4 students activity recordings showed that changing the disposition of the interface, either for exploration purpose or comfort, was the action most performed by students (constituting 31,4% of all actions). Different students had different opinions concerning which element should be displayed at a large size at the core of the interface.

Students expressed very few questions or alerts, with a maximum of two questions per student during the experiment session. The note-taking activity is much more contrasted : 28.4% of actions belong to the note-taking category, but most of these were undertaken by two students. Other observed students only took notes twice each during the session. Consultation of the resources panel represents 19.6% of actions, almost equally divided between browsing the list and the actual consultation of links provided.

The teacher was not accustomed to monitoring the indicators on his computer screen in real time, and quickly abandoned this behavior to adopt a more usual posture. His evaluation of the students' activity

was hence limited to the audience physically present in the classroom.

5 ON-GOING WORKS AND FUTURE PERSPECTIVES

Given the high number of class configurations and personal preferences towards practices, we plan to further develop the tool by adding new tools, and a more flexible interface.

The pedagogical situation considered next is the group work situation. The constraints in this situation are different from those considered in this paper. The communications between actors is more complex than during a lecture as all students and the teacher must be able to communicate. The activity of the students varies greatly depending on both the work expected and personal behaviour of each student. Additionally, the condition under which the transcription software operates are less controlled. Microphones, student voices and vocabulary used are very diverse, and a challenge for the system. The explored solution is a lesser important use of the transcription and an emphasis on student's activity indicators.

Currently, post-session access to course content is limited to an export of the text posted, generated and typed. The video recorded can also be watched separately. We plan to develop a browsing system allowing students to retrieve all this content on a single interface, synchronized together. Setting the replay to a particular point in time would display the corresponding pieces of transcription and resources. Due to the large amount of content, readability and browsing can become a tedious task; a thematic segmentation system will be implemented in order to divide the session into smaller chapters, easier to label and to browse individually.

As for the teacher, post-session use of traces and generated content is for now limited to the display of the contents of the Moodle database. Future work involves the conception of an editorial toolset allowing to retrieve and export more easily the content produced and the feedback collected during a particular lecture. This would allow the teacher to easily produce and share online pedagogical material relevant to a particular theme, greatly reducing the time investment required to set up a SPOC or a flipped classroom configuration.

As the prototype comes closer to a stable version, we plan to release it to the Moodle plugins directory and make it available to a greater public.

REFERENCES

- Bettenfeld, V., Choquet, C., and Piau-Toffolon, C. (2018). Lecture instrumentation based on synchronous speech transcription. In *18th International Conference on Advanced Learning Technologies*, pages pp. 11–15, Bombay, India.
- Crétin-Pirolli, R., Pirolli, F., and Bettenfeld, V. (2017). Projet pastel, performing automated speech transcription for enhancing learning, analyse des besoins. Unpublished deliverable document.
- Heilesen, S. (2010). Podcasts for efficient learning. In *Proceedings of the 7th International Conference on Networked Learning 2010*, pages 971–972. L. Dirckinck-Holmfeld, et. al.
- Ho, I., Kiyohara, H., Sugimoto, A., and Yana, K. (2005). Enhancing global and synchronous distance learning and teaching by using instant transcript and translation. *Cyberworlds*.
- Hwang, W. Y., Shadiev, R., Kuo, T. C., and Chen, N. S. (2012). Effects of speech-to-text recognition application on learning performance in synchronous cyber classrooms. *Journal of Educational Technology & Society*.
- Iglesias, A., Jiménez, J., Revuelta, P., and Moreno, L. (2016). Avoiding communication barriers in the classroom: the apeinta project. *Interactive Learning Environments*, 24(4):829–843.
- Ngo, C.-W., Wang, F., and Pong, T.-C. (2003). Structuring lecture videos for distance learning applications. In *Multimedia Software Engineering, 2003. Proceedings. Fifth International Symposium on*, pages 215–222. IEEE.
- Ryba, K., McIvor, T., Shakir, M., and Paez, D. (2006). Liberated learning: Analysis of university students' perceptions and experiences with continuous automated speech recognition. *E-Journal of Instructional Science and Technology*, 9(1):n1.
- Shadiev, R., Hwang, W. Y., Chen, N. S., and Yueh-Min, H. (2014). Review of speech-to-text recognition technology for enhancing learning. *Journal of Educational Technology & Society*.
- Traphagan, T., Kucsera, J. V., and Kishi, K. (2010). Impact of class lecture webcasting on attendance and learning. *Educational Technology Research and Development*, 58(1).
- Uke, N. J. and Thool, R. C. (2012). Segmentation and organization of lecture video based on visual contents. *International Journal of e-Education, e-Business, e-Management and e-Learning*, 2(2):132.
- Wang, F. and Hannafin, M. (2005). Design-based research and technology-enhanced learning environments. *educational technology research and development*, 53(4), 5-23. *Educational Technology Research and Development*, 53:5–23.