

Cohen, Cathy and Guichon, Nicolas (2016). "Analysing multimodal resources in pedagogical online exchanges: Methodological issues and challenges". In Caws, Catherine, and Hamel, Marie-Josée (eds), *Language-Learner Computer Interactions: Theory, methodology and CALL applications*. John Benjamins. Pp. 187-213.

Analysing multimodal resources in pedagogical online exchanges: Methodological issues and challenges

Cathy Cohen, ESPE Université Lyon 1
Nicolas Guichon, Université Lyon 2

Abstract

This chapter focuses on the contribution to webconferencing based pedagogical synchronous interactions of meaning-making multimodal resources (spoken language as well as gesture, gaze, body posture and movement). The first part of the chapter explores different methodological approaches to the analysis of multimodal semiotic resources in online pedagogical interactions. Having presented an overview of what research into synchronous web-mediated online interaction can bring to the field of CALL, we discuss the importance of determining the relevant units of analysis which will impact on the granularity of transcription and orient the ensuing analyses. With reference to three of our own studies, we then explore different methods for studying multimodal online exchanges depending on the research questions and units of analysis under investigation. To illustrate the various ethical, epistemological and methodological issues at play in the qualitative examination of multimodal corpora, the second part of the chapter presents a case study that identifies the different steps involved when studying online pedagogical exchanges, from the initial data collection phase to the transcription of extracts of the corpus for publication.

Key words

multimodal resources; web-mediated pedagogical interaction; units of analysis; webcam; transcription; multimodal corpora

Introduction

As a result of globalization and easy Internet access, opportunities for exposure to foreign languages have greatly increased over the past two decades (Kern 2014). Not only can language learners access all types of documents (e.g., films, audio and video documents, written texts, images) quickly and simply, they can also exchange synchronously or asynchronously with speakers of the target language, opening up seemingly unlimited possibilities for foreign language contact and potential learning. These might be informal social interactions as learners seek out opportunities to use the target language with their peers, but they may also be specifically designed as pedagogical exchanges between a language teacher and learner, or between two learners under the coordination of a language teacher. Indeed, more and more language learning courses take place online, often between language teachers in one country, and language learners in another. Such courses may involve both asynchronous (e.g., email or blogging) and synchronous (e.g., text chat or

videoconferencing) tools. As a result, new interaction patterns and norms are constantly developing and these combine a broad range of semiotic modes (Sindoni 2013), which potentially offer new and diverse opportunities for learning.

The current chapter focuses on pedagogical synchronous interactions which use desktop videoconferencing (henceforth DVC), described by Kern as “a quintessential technological support for providing communicative practice with speakers at a distance, since it is the closest approximation to face-to-face conversation” (2014: 344). This powerful instrument to learn languages is an Internet-based system enabling two or more people located in different places to communicate online with simultaneous two-way audio and video transmission (Sindoni 2013). The video transmission, made possible thanks to a webcam on each participant’s computer, gives access to several meaning-making modes, including spoken language, but also other multimodal elements such as gesture, gaze, body posture and movement. With the growing number of online language courses and telecollaboration projects, it is clearly important for CALL practitioners to gain a better understanding of how these multimodal resources contribute to the pedagogical setting and to learning contexts, and also how the different semiotic resources are orchestrated in interactive technology mediated situations (Stockwell 2010).

This chapter will analyze the contribution of multimodal resources to pedagogical online exchanges. The first part explores, the different methodological approaches to the analysis of multimodal semiotic resources in online pedagogical interactions. We begin by briefly reviewing recent literature in order to take an overview of what research into synchronous web-mediated online interaction can bring to the field of CALL. The issues of determining the relevant units of analysis will be discussed as the latter have a clear impact on the granularity (i.e., the amount of detail provided by researchers) of transcription and orient the ensuing analyses (Ellis and Barkhuizen 2005). Then, with reference to three of our own studies, we explore different methods that can be employed to study multimodal pedagogical exchanges depending on the research questions and the units of analysis under investigation. In the three studies, our focus is on the role played by technological mediation in online pedagogical exchanges and in particular, on the affordances provided by the webcam (see Chapter 3 in this volume).

To illustrate the different ethical, epistemological and methodological issues at play in the qualitative examination of multimodal corpora, the second part of the chapter will present a case study that identifies the different steps involved in the study of online pedagogical exchanges, from the initial data collection phase to the transcription of extracts of the corpus for publication. The case study is an extract from Study 2 which is presented in the first part of this chapter.

Methodological approaches to the study of multimodal pedagogical interactions

In this section, we focus on different methodological approaches that can be employed to analyze how multimodal semiotic resources function in online pedagogical interactions. Studies exploring how these interactions are mediated and organized by the webcam are still quite limited and different units of analysis have been the focus of recent research. Determining the relevant units of analysis is important as they have a clear impact on the type of data collected (quantitative or qualitative, see Table 1), on the granularity of transcription and will orient the ensuing analyses (Ellis and Barkhuizen 2005).

We use the term *unit of analysis* to refer to the general phenomenon under investigation. Once the unit of analysis has been identified, it has to be operationalized by researchers who must then select the variable(s) which they are going to investigate. These are the features which the researchers believe constitute the unit of analysis (see Table 1 for examples). Several examples taken from the field of pedagogical DVC interactions are provided here to illustrate this. The unit of analysis for Wang's (2007) study was design principles for videoconferencing tasks. One of the components she explored was the role played by the webcam image in task completion. Using personal observation and post-session interviews with the small group of learners who participated in the study, she concluded that facial expressions and gestures visible via the webcam were key features that facilitated task completion. Satar (2013) focused on how social presence was established in online pedagogical DVC interactions. She explored how the trainee teachers interacting with one another used gaze, and how they compensated for the impossibility of direct eye contact. She identified a range of different uses of the webcam and highlighted the importance of eye contact for the establishment of social presence in online multimodal interactions. Guichon and Wigham (in review) explored the potential of the webcam for language teaching, focusing particularly on the unit of analysis of framing, in other words how trainee teachers framed themselves in front of the webcam and, as a result, what information was made visible to their learners within the frame of the video shot. So, they investigated how trainee teachers made use of the affordances of the webcam to produce non-verbal cues that could be beneficial for mutual comprehension (see Study 3 below for more details). Their results emphasized the need for trainee teachers to enhance their critical semiotic awareness, including paying closer attention to framing, thus enabling them to gain a finer perception of the image they projected of themselves. In so doing, it is hypothesized that they should be able to take greater advantage of the potential of the webcam and, as a consequence, increase their online teacher presence.

Different methods can be employed to study pedagogical online exchanges and researchers' choice of method will depend on the research questions they wish to investigate and the objectives of their study. We will take three examples from our own work to illustrate different approaches. In all three, we are interested in the role played by technological mediation in online pedagogical exchanges and our particular focus is on the affordances provided by the webcam. Table 1 provides an overview of these studies which will be discussed in turn below.

	Duration	Type of data	Task	Design	Number of participants	Unit of analysis	Features/Variables studied	
Study 1	1 interaction lasting around 10 minutes per student	Quantitative	Describe pictures	4	Experimental	40	Learner perceptions of online interaction	Feeling of psychological and physical presence; understanding of and by teacher; quality, naturalness and enjoyment of interaction
							Rhythm of interaction	Silences; overlaps; turn duration; number of words
							Word search episodes	Frequency; duration
Study 2	1 interaction lasting around 10 minutes per student	Qualitative	Describe pictures	4	Experimental	3	Word search episodes	Multimodal orchestration of speech and non-verbal features (e.g., gaze, nods, gestures, facial expressions)
Study 3	1 weekly interaction lasting around 40 minutes over a 6-week period	Quantitative and qualitative	Range of different tasks and open-ended conversation		Ecological	12	Framing choices	Teachers' semiotic self-awareness
						3	Visibility of gestures in and out of the webcam	

Table 1: Overview of studies on affordances of the webcam

Study 1: Quantitative approach on experimental data

The first study, reported fully in Guichon and Cohen (2014), adopted a quantitative methodology and had an experimental design. In this study, we explored the impact of the webcam on an online interaction by comparing several dependent variables between an audioconferencing and a videoconferencing condition, using Skype (<http://www.skype.com>). In the audioconferencing condition, the webcam was switched off, whereas it was on in the videoconferencing condition. Our objective was to assess the webcam's contribution to the interaction. There were three research questions, each of which explored different units of analysis which we felt might operate differentially in the two experimental conditions. The first was learner perceptions, which were probed using a short post-task Likert scale questionnaire to gauge learners' feelings of: (1) the teacher's psychological and physical presence, (2) understanding of and by the teacher, and (3) the quality, naturalness and enjoyment of the conversation. The second explored the rhythm of the interactions by measuring silences, overlaps, turn duration and number of words. The third focused on word search episodes by measuring their frequency and duration. Before the experiment began, we had clear hypotheses which stated that being able to see one's interlocutor would make a difference to the online pedagogical interaction. In other words, we stated that we expected to find a statistically significant difference between all the dependent measures under investigation in the audioconferencing and videoconferencing conditions. Furthermore, for the dependent measures relating to learner perceptions, we predicted that the videoconferencing condition would be received more favorably than the audioconferencing condition.

The independent variables were strictly controlled before the experiment began. Forty French students who had a similar level of English, the foreign language they were learning at university, took part in the experiment. Twenty of them were put in the videoconferencing condition and 20 in the audioconferencing condition. Indeed, in order to be able to carry out certain statistical tests, it was necessary to have at least 20 participants in each condition. Statistical tests were used to verify that there were no significant differences between the two groups in terms of sex, age, English level, familiarity with online communication tools and attitudes towards speaking English. Had there been differences between the two groups at this stage, we could not have been sure whether our results were due to initial group differences or rather to differences resulting from the testing conditions. In the experiment, each student interacted individually with the same unknown native English-speaking teacher who was always in the same setting. Furthermore, they all did exactly the same task, which consisted in describing four previously unseen photographs. The duration of the interaction for all participants was set at around ten minutes.

In order to compare the different dependent variables between the two experimental conditions and assess the contribution of the webcam, it was necessary to carry out a quantitative study. In other words, we had to be able to measure the different variables in the two experimental conditions to see how they compared. So, for example, number of silences and word search episodes were counted and turn durations were measured (see *Annotation* below for more details as to how this was achieved). All the data were then imported into SPSS (<http://www-01.ibm.com/software/fr/analytics/spss/>), a computer-based statistical package for analysis, allowing statistical comparisons to be made between participants in the two conditions. Our results showed that, contrary to our predictions, there were fewer differences than we had anticipated between the videoconferencing and audioconferencing conditions on the dependent measures, with few comparisons reaching statistical significance. The main difference was the greater number of student silences in the audioconferencing condition.

This first study was clearly time consuming in terms of data collection and analysis. It also involved many people: 40 students, a teacher, an assistant who helped organize the data collection sessions, four research assistants to transcribe and annotate the data (see *Annotation* below) and two researchers who analyzed the data and wrote up the research for publication. Although the differences between the results obtained from the two experimental conditions were far less clear-cut than we had expected, the results were nevertheless thought provoking. Indeed, we considered that although from a quantitative point of view the presence of the webcam did not seem to have a great impact on the pedagogical interactions with regard to the units of analysis which were investigated, the webcam image could nevertheless be facilitative, could modify the quality of the mediated interaction and that the reality was in fact considerably more complex than our findings seemed to show. Hence these results also highlighted the limitations of using quantitative data to grasp the more subtle interactional aspects in a multimodal learner corpus. Furthermore, our results provide a good example of the iterative process of research, with the first more generic experiment being a necessary step to reveal the need to explore particular parts of our corpus using a much finer grained analysis. This led us to conduct our second study.

Study 2: Qualitative approach on experimental data

In this study (Cohen and Guichon 2014), we carried out a qualitative and descriptive analysis on small sections of the videoconferencing data taken from the first experimental study. In

other words, we used part of the same corpus used in Study 1, but this time, to conduct a microanalysis. The analysis focused on short sections of just three of the 20 videoconferencing interactions, in order to examine how the learners and the teacher used the webcam strategically at different times during their exchanges.

Since we were particularly interested in training language teachers to utilize the affordances of the webcam during pedagogical online interactions and to develop their critical semiotic awareness, we considered that only a fine-grained analysis of non-verbal behavior in the videoconferencing condition would enable us to identify when and how the interaction was facilitated by the appropriate use of the webcam by participants.

The methodology employed in Study 2 was quite different from the first. This time, we worked within the Conversation analysis (CA) paradigm, as articulated in work initially conducted by gesture specialists (e.g., McNeill 1992) and more recently pursued by researchers working on gesture in the field of Second Language Acquisition such as McCafferty and Stam (2008) and Tellier and Stam (2010). We adapted the methodology of these authors who focus on face-to-face pedagogical interactions in order to investigate pedagogical computer-mediated interactions. We also integrated an approach from the broader domain of multimodal discourse analysis, as applied by Norris (2004) and Baldry and Thibault (2006) whose work is not conducted in the pedagogical field. Finally, our approach was influenced by recent work carried out by Sindoni (2013) who has explored non-pedagogical online interactions using a multimodal approach. In other words, the methodological approach we adopted was influenced by work conducted in several domains of scientific research. By combining and adapting elements from these different areas, we created a method suitable for analyses in our own field of investigation, i.e., the study of multimodal resources in pedagogical online exchanges.

In this second study, we explored the contribution to meaning making of several nonverbal semiotic resources other than speech and investigated how they helped the teacher to manage the online pedagogical interaction and how they were orchestrated. Each of these semiotic modes will now be presented briefly, with specific reference as to how they function in an online DVC interaction.

We considered *proxemics*, that is to say the physical distance individuals take up in relation to one another and to objects in their environment. Proxemics functions quite differently when interacting online using DVC, since interactants are not in the same location. Sindoni observes that “distance is not established by those who interact, but between one participant and one machine. This distance foregrounds the *representation* of distance among users.” (2013: 56). Added to this, whatever position the user chooses, because he has constant access to his own image in the smaller frame on his computer screen, he is able to monitor and manipulate the image he wishes to project to his interlocutor (Sindoni 2013). This affordance provided by web-mediated communication also gives the user greater control over the construction and negotiation of social space.

We examined different types of *gesture*, defined as the use of the hands and other body parts for communicative purposes (MODE 2012). We focused in particular on those gestures which were visible in the webcam: iconic gestures representing an action or an object; metaphoric gestures illustrating an abstract concept or idea; and deictic gestures used to point towards concrete or abstract spaces. Our objective here was to assess what type of information was communicated by these gestures and to what extent they appeared to facilitate (or not) the

online exchange. For instance, were they transmitting some information to the interlocutor to complement or accompany what was said in the verbal channel (co-verbal gestures)? Or were they self-regulatory gestures, produced unintentionally to help speakers to think, thereby allowing them to maintain a sense of coherence for themselves (McCafferty 2008)? To what extent were they visible in the webcam?

Head movements, which may convey meaning between interlocutors (e.g., nodding in agreement; shaking one's head from side to side to convey disagreement; holding one's head quite still while fixing one's gaze on someone to indicate concentration and focus), were also considered.

Finally, we were interested in *eye contact, gaze* and *facial expressions*. Compared to face-to-face conversation, gaze management is very different in online video interactions. With the current state of technology used in DVC systems, it is impossible for speakers to make direct eye contact with one another (see De Chanay 2011). When speakers direct their eyes to their interlocutor's image on their computer screen, their eyes are slightly lowered, so not aimed directly at their interlocutor's eyes. They can choose to look directly at the webcam which gives the interlocutor the impression that he is being looked at straight in the eyes, but in so doing, paradoxically the speaker can no longer focus on the interlocutor's image on the screen (De Chanay 2011). So, not only are there fewer visible gestures to facilitate communication and intercomprehension in videoconferencing interactions, but there is also the impossibility of mutual gaze. Cosnier and Develotte (2011) hypothesize that speakers compensate for this through facial expressions which become more important and seem to be more numerous and perhaps over-exaggerated in videoconferencing interactions compared to face-to-face conversations, precisely to compensate for the lack of visible hand and arm gestures.

The different non-verbal semiotic modes have been discussed separately here, but of course during any chosen communicative event, they are operating simultaneously and, as Sindoni (2013: 69) argues "Ensembles of semiotic resources [...] produce effects that differ from those produced by a single semiotic resource *and* from the mere *sum* of semiotic resources". A transcript and microanalysis taken from this study corpus is provided below (see *Transcript and analysis presentation*) as an illustration of our approach. Since the study was exploratory, our hypotheses emerged progressively as the data were explored. Three angles of analysis became apparent: (1) self-regulatory versus co-verbal gestures; (2) gestures which contribute something to the construction of the message versus gestures which potentially cause interference and are distracting and (3) redundant gestures which duplicate what is said in the verbal channel versus to complementary gestures which add some new information.

This qualitative study provided us with rich and complex data, enabling us to gain insights into the multimodal orchestration of the different semiotic resources in an online pedagogical interaction. However, we were using data collected for a study carried out in experimental conditions – the interaction duration was fixed; it was the first time that both the teacher and the learners had met and taken part in an online pedagogical interaction. So, the findings may have been attributable, to some degree at least, either to the novelty of the learning situation and/or the task learners were asked to carry out. In other words, the conditions of this second study, and indeed the first, lacked ecological validity. Thus in our third study, we tried to address this methodological shortcoming.

Study 3: Quantitative and qualitative approach on ecological data

As shown in Table 1, the corpus for the third study was collected in ecological conditions. The context was a telecollaborative project¹ in which 12 trainee teachers of French as a foreign language met for online sessions in French with undergraduate Business students at an Irish university. Each trainee teacher met with the same learner (or pair of learners) once a week for approximately 40 minutes over a six-week period. Over this period, the trainee teachers proposed a range of different interactional tasks to their learners. So, unlike Study 2, which was conducted in experimental conditions, i.e., it was set up with the sole purpose of conducting an experiment to test our different hypotheses, Study 3 used data collected from an online course that was set up between two universities with learner training in mind: helping Irish learners to develop their interactional skills in French, and helping students training to be French teachers to develop their online teaching skills. Thus this teaching and learning situation was not set up initially for research purposes but the data collected from the online sessions were used subsequently to conduct research.

The research carried out in this study (Guichon and Wigham, in review) focused on very specific elements taken from the sizeable corpus that was collected. As in the previous two studies, we were interested in how participants used the affordances of the webcam, but this time the particular focus was on framing, i.e., how the trainee teachers framed themselves in front of the webcam and, as a result, what information was made visible to their learners within the frame of the video shot. For the qualitative part of the study, the same method of analysis was used as in Study 2. Two questions were explored here. Firstly, in order to study teachers' framing choices, screenshot images were taken of the 12 trainees each week over six weeks, at around minute 17 of their online interaction. A quantitative approach was adopted to provide an indication of the frequency of the trainees' different framing choices along a continuum, from extreme close-up shot, to close-up, to head and shoulder shot, to head and torso shot. In parallel, a qualitative approach was used to conduct a fine-grained analysis on the same data and, in particular, how the trainees positioned their gestures in relation to the webcam over the six-week course.

The findings revealed that, head and shoulder shots, followed by close-up shots of themselves were those most favored by the trainee teachers. Furthermore, qualitative analysis of the data showed that certain trainee teachers adjusted the position of some of their gestures, in particular highly communicative iconic and deictic gestures, so that they were framed and therefore more likely to be visible to learners and, therefore, potentially helpful for learner comprehension. Furthermore, quantitative analyses revealed that these gestures were held for longer in front of the webcam. So such teaching gestures, which clearly had a communicative purpose, appeared to be produced by these trainee teachers quite intentionally, and consequently were aimed at the webcam and remained visible to the language learners for some time.

The second question investigated in this study explored the communicative functions of gestures that were visible or invisible in the frame. For technical and practical reasons explained fully in the study, data were collected for just three participants for just one session each. The teacher trainees were filmed using DVC with their learners with two distinct recordings. One captured the on-screen activity, so what was visible and audible through the webcam, and an external camera was used to film what lay outside the webcam's view (the

¹ ISMAEL projet: <http://nicolas.guichon.pagesperso-orange.fr/projets.html>

hors champ). When the two sets of recordings were compared, it became clear that the trainee teachers continued to perform many potentially co-verbal gestures which were either invisible or only partially visible in the webcam recordings which only captured a close-up of the head and upper torso area. In contrast, extra-communicative gestures, such as touching their hair or scratching their ear, become much more visible because of the magnifying effect provided by the very restricted view offered through the webcam. Such gestures, which may have gone unnoticed in a face-to-face interaction because of the presence of other broader contextual elements, were more difficult to miss when communicating using DVC. Indeed, if numerous, they could become rather distracting and interfere with communication.

So, the findings of this study highlighted the need to train teachers “to become critically aware of the semiotic effect each type of framing could have on the pedagogical interaction so that they made informed choices to monitor the image they transmitted to their distant learners according to an array of professional preoccupations” (Guichon and Wigham, in review). This ecological study provided valuable information which could be reinvested in future teacher training courses.

Synthesis

We have explored three different studies, each of which explores the role of HCI in online pedagogical exchanges, with a particular focus on the affordances provided by the webcam. From the first generic study which was experimental and quantitative, through to the third study which had a specific focus, was ecological and combined both quantitative and qualitative analyses, we have shown that the method adopted will depend on the research questions under investigation. Both quantitative and qualitative analyses are valid means to explore the data collected, as long as the method is sound and the objective clearly stated. The qualitative microanalysis of a much broader range of units of analysis investigated within the field of webconferencing-supported teaching is certainly to be encouraged in order to further enhance our knowledge of HCI in a pedagogical setting. By putting certain elements of the interaction into the spotlight, we may progressively untangle the complexity of these online pedagogical exchanges.

In the first part of this chapter, we have explored different methods for studying multimodal resources in pedagogical online exchanges. However, in order to be able to conduct the type of analyses presented above, researchers have to ensure that their data are collected and stored in such a way that they can be later transcribed and annotated. Whether the study is quantitative and experimental or qualitative and ecological, numerous transformations are required to progress from the initial data collection stage to the creation of a corpus that can be presented in academic publications or at conferences, and also perhaps be shared among researchers.

In the next section of this chapter, we examine these different stages and investigate the opportunities and challenges concerning the study of data relating to synchronous mediated language learning and teaching.

Reflections on a multimodal approach to synchronous pedagogical online interactions

From the traces of mediated activity to a corpus that can be studied from different perspectives

Any mediated learning activity produces traces: digital traces, currently much used in the field of Learning Analytics, can be computer logs that provide quantitative information (frequency of access, time spent on a task, number of times a given functionality is used, etc.). The aim of these digital traces is to understand and optimize learning and learning environments (Siemens and Baker 2012). Digital traces can also be comprised of “rich histories of interaction” (Bétrancourt, Guichon and Prié 2011: 479) that provide multimodal data and time stamps that can be gathered from digital environments in order to gain an insight into certain teaching and learning phenomena. This second form of traces has been studied by researchers in the field of computer-mediated communication (CMC) for the last 20 years (see for instance Kern 1995; Kost 2008; Pelletieri 2000). Thus, traces collected in forums, blogs, emails, audiographic platforms and DVC have been built into corpora to study the specificities of mediated language learning usually by using conversation and/or interaction analytic tools.

The present section focuses on mediated learning interactions to illustrate how technology helps fashion methodological and scientific research agendas in the field of mediated interactions. Several operations are at play when researchers deal with a data-driven study of multimodal learning and teaching, when they strive to create a corpus that can offer different types of analyses as was illustrated in the first part of this chapter.

If we take the example of a corpus composed of recordings of online learning interactions mediated by a DVC facility, four main operations can be identified: corpus building, annotation, data transcript and presentation. Each of these operations will be explained and illustrated by a case study using data that were initially collected for a larger research project (Guichon and Cohen 2014, discussed in Study 1 above). However, before we do this, it is important to underline the ethical aspects that researchers must respect when dealing with data which include participants’ images.

Ethical considerations

Ethical issues are relevant to all research involving humans. In the case of the type of studies we have described above, which may involve the publication of participants’ images, certain issues should be considered very carefully.

Before recording begins, researchers must obtain written informed consent from participants: first, that they agree to be recorded; second, that they agree to be recorded for research purposes; and third, that they agree that recordings (or screenshots) may be displayed publicly or published (ten Have 1999). If participants consent to all three, they must understand fully what is at stake. For example, will they be recognizable from the recordings (visual, auditory)? Will their faces be blurred/pixelated to avoid recognition? Where will the recordings be shown and where will they be published? Will they be available freely online to anyone (for an (un)limited period of time)? Will participants have access to the recordings before they are used, in order to confirm or cancel their informed consent? (see Yakura 2004 for an excellent discussion of the issues at stake here).

The above questions present real challenges for researchers. First and foremost, if recordings or screenshots are to be used publicly, anonymity cannot be ensured at every stage (Yakura 2004). Secondly, depending on what participants have consented to, researchers may be more restricted in what they can present and/or publish. If, for instance, researchers wish to provide a finely grained analysis of the different non-verbal semiotic modes employed by participants,

but they are only authorized to publish faces which have been blurred, displaying eye contact, gaze and facial expressions becomes impossible, thus “rendering the data unusable for certain lines of linguistic inquiry” (Adolphs and Varter 2013: 149).

How can researchers circumvent this problem in order to preserve and communicate to others some of the richness of their data? To compensate to some extent for the loss of visual information, researchers could provide very detailed written descriptions (Lamy and Hampel 2007). In a recent study by Sindoni (2013), because of reservations expressed by certain participants about the publication of screenshots, she opted to use drawings instead. However, she recognises the drawbacks of this, stating that “they are time-consuming and require specific expertise, so that they can be used selectively, only for very brief and fine-grained analyses. Furthermore, drawings incorporate the researcher and artist’s bias that represent participants in their interactions.” (Sindoni 2013: 71).

Multimodal data collection

Several applications, for instance *Camtasia* (<http://camtasia-for-mac.en.softonic.com>) or *Screen video recorder* (<http://www.dvdvideosoftware.com/fr/products/dvd/Free-Screen-Video-Recorder.htm#.VHBk8Eve5g0>) can be used to capture on-screen activity in an online interaction and this can be converted into a video (for more details, see Chapter 7 in this volume). The advantage of such applications is that they can be installed beforehand on each participant’s computer and once switched on, they capture everything that is visible on the screen and audible around the screen, thus providing researchers with access to all the actions and utterances produced by the participants during the online interaction. Hence, whether the study is experimental or ecological (see above), traces of the mediated activity can be collected with little or no interference on the ecology of the learning situation. This is quite different from classroom-based research that requires more intrusive devices (i.e., video cameras) to collect traces of the observable teaching and learning activities.

While the traces of the mediated learning activity constitute the main material of the study, complementary data have to be collected via consent forms, researchers’ field notes, pre- and post- interviews or questionnaires with the participants to gather crucial information about:

- Ethical dimensions (as discussed above);
- Socio-demographics and learner profiles: age of the participants, gender, relations to one another (in case of an interaction), familiarity with the given program or application, level in target language and motivations, experience in learning or teaching online;
- Pedagogical dimensions: nature of the interaction, tasks, themes, documents used, instructions, place within the curriculum;
- Temporal dimensions: length of each interaction, frequency of interactions (e.g., once a week), duration of module (e.g., a semester);
- Methodological dimensions: how participants were recruited for the study, how their level was assessed, how they are divided (in case of an experimental study that compares two or several groups), what they were told of the aim(s) of the study, precisely how the data collection was organized, how ethical considerations were taken into account (see above);
- Technological dimensions: type of software and hardware used (e.g., desktop or laptop, devices used for recording, etc.).

The conjunction of field notes, questionnaires, interviews, consent forms with the main data thus help create “a dynamic constellation of resources, where meanings are produced through inter-relationships between and within the data sets, permitting the researcher literally to “zoom in” on fine-grained detail and pan out to gain a broader, socially and culturally, situated perspective” (Flewitt et al. 2009: 44).

The data that serve as the illustration for this chapter come from Study 2, discussed above. The reader will need to know a number of elements about the two participants who took part in a larger study (Study 1 discussed above, see Guichon and Cohen 2014). The learner was 20 at the time of the study. His level had been assessed as B2 (according to the Common European Framework of Reference for Languages) and he described himself as a keen language learner. He used Skype for social purposes but had never used it for language learning. It was the first time he had interacted with the 28-year-old female native teacher and this interaction was not part of his usual class. The teacher had several years of experience teaching non-specialist university students in a classroom setting and was a regular user of Skype, mainly for personal communication. However, this was the first time that she had taken part in an online pedagogical interaction. Neither of the two participants was informed of the study’s purpose or hypotheses before the experiment. The task consisted in getting the student to describe four previously unseen photographs. These photographs were chosen because each one contained what were considered to be problematic lexical items likely to trigger word search episodes, chosen as the unit of analysis for this research. The interaction via Skype lasted for about 10 minutes and participants were asked not to use the keyboard.

All the secondary data (field notes, questionnaires and interviews) had to be digitalized and grouped together with the data comprising the traces of the mediated interaction “to reconstitute for researchers, in as many ways as desired, information about the original experience” (Lamy and Hampel 2007: 184) and to enrich subsequent analyses.

Annotation

There are several computer software tools that researchers can use to code audio and video data. Among these, ELAN (<http://tla.mpi.nl/tools/tla-tools/elan/>) is a linguistic annotation tool devised by researchers at the Max Planck Institute (Sloetjes and Wittenburg 2008). Figure 1 shows a sample of the data that were annotated with ELAN with which the researchers can:

1. Access the video stream of one or up to four participants;
2. Play the film of the interaction at will with the usual functionalities to navigate it;
3. View a time line aligned with media time;
4. Transcribe, on the horizontal axis, the utterances of the participants (one layer per participant);
5. Add a new layer for each element they wish to investigate (indicating for instance the onset and the end of a gesture and its description);
6. View annotations of one layer in a tabular form to facilitate reading.

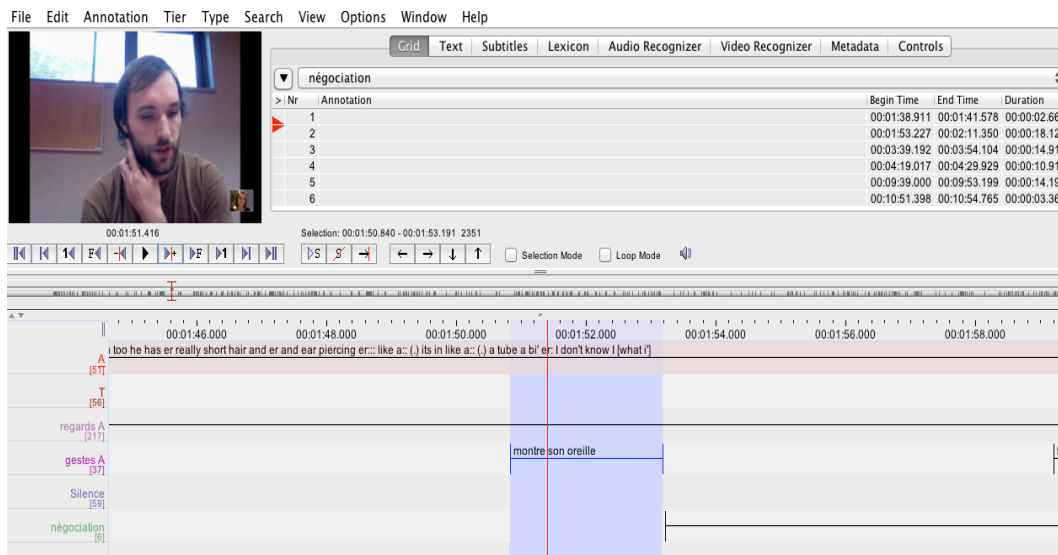


Figure 1: Example of a sample of data annotated with ELAN

With ELAN, there can be as many layers (called tiers) as is deemed useful for a given study (i.e., words, descriptions, events, translations, etc.). As the case study presented here focuses on the verbal and co-verbal behavior of the learner who has to describe four photographs to a distant teacher via Skype, the direction of the eyes, the gestures that were produced (e.g., *points to his ear*), the silences between and within turns were annotated because they all were crucial elements when a learner is speaking in his L2 and is confronted with word search episodes. Researchers working on multimodal data can thus align different features of the interaction, accurately transcribe data across modes and then obtain a variety of views of the annotations that can be connected and synchronized.

The data from the three studies described in the first part of this article were all transcribed using ELAN. Hence, although the first study was quantitative and the second qualitative, the same annotation tool was used for both although the tiers which were the focus of each study were different.

Annotation corresponds to a necessary transformation of the data in view of further analysis. It is a time-consuming and demanding task that requires devising a coding scheme so that all annotations are consistent across different annotators. As noted by Adolphs and Varter (2013:155), coding schemes have to be carefully explained and recorded so that “they can be shared across different research communities and with different community cultures and different representational and analytical needs”. It is methodologically sound to get two different researchers to annotate a sample of the same data in order to ensure that the coding scheme is sound. This can be verified by calculating the inter-rater reliability to determine for instance whether two researchers interpret and code gestures consistently and reach a satisfactory level of agreement. If they fail to do so, the annotation scheme needs to be refined and retested in the same way until satisfactory inter-rater reliability is achieved (Allwood et al. 2007). Yet, as noted by Calbris (2011: 102), “achieving the ideal of scientific objectivity when coding a corpus is a delusion, because coding depends on perception, an essentially pre-interpretative and therefore subjective activity”.

Furthermore, priorities and research questions have to be carefully defined beforehand so that the granularity of the annotations does not evolve. Researchers such as Flewitt et al. (2009) underline that annotation already corresponds to a first level of analysis since it entails

selecting certain features of the mediated interaction and leaving others out according to both a research rationale and agenda.

Once the data have been annotated, they can then be organized into a coherent and structured corpus (see Chapter 10 in this volume for a full account of corpus building and sharing). They may also be put on a server allowing them be shared with other researchers. In order to do this, close attention has to be paid to the formats of the data so that they are compatible with different computer operating systems. Providing researchers with clear information as to how to access the data, specifying all the contextual information (see above) and ethical dimensions (e.g., what can be used for analysis and what cannot be used for conferences or publications because participants have withdrawn their permission) are important steps to make the corpus *usable*, *searchable* and *sharable*. The field of CMC would greatly profit from having more researchers working on the same corpora: not only would it reduce the costs associated with corpus building, transcription and annotation, but it would provide researchers with the opportunity of examining the same data using different tools, methods and research questions and would therefore produce results that can present more significance and reliability to the community at large.

Transcript and analysis presentation

Once the data have been organized into a coherent corpus, analyses can be made starting by the making of the transcript. Bezemer (2014) allocates two functions to the making of a multimodal transcription. The first function of transcription is *epistemological* and consists of a detailed analysis of a sample of an interaction in order to “gain a wealth of insights into the situated construction of social reality, including insights in the collaborative achievements of people, their formation of identities and power relations, and the socially and culturally shaped categories through which they see the world” (Bezemer 2014: 155). The second function is *rhetorical* in that the transcript is designed to provide a visual transformation of the trace of the interaction that can be shared with readers in a scientific publication. Transcripts chosen and prepared for an article are not illustrations of a given approach or theory but are both the starting-point of the analysis and the empirical evidence that supports an interpretation and can be shown as such to readers. The researcher must therefore find an appropriate time-scale (e.g., a few turns, an episode, a task, a series of tasks, a whole interaction) to study a phenomenon (for instance negotiation of meaning in a mediated pedagogical interaction) and then define the boundaries of the focal episode. Making the transcript may also involve refining the initial research questions and determining what precise features will be attended to.

For our study on videoconference-based language teaching, it seemed crucial to understand how the distant teacher helped the learner during word search episodes and used the semiotic resources (such as gestures, facial expressions and speech) at her disposal. It was equally important to examine how the learner used different resources to signal a lack of lexical knowledge and how meaning was negotiated with the native teacher. The interplay of gestures, head and body movements, gaze and facial expressions produced by both participants while the learner was trying to describe a photograph became features that were selected as especially important for the transcript (see Figure 2). Although conventions used for Conversation analysis can be adjusted to multimodal transcription, new questions arise concerning the representation of co-verbal resources (gesture, gaze) with text, drawings or video stills and the alignment of these different representations so that the reader can capture how verbal and nonverbal resources interact (see Figure 2). Ochs (1979 as cited in Flewitt et

al. 2009: 45) underlined the theoretical importance of transcription, arguing that “the mode of data presentation not only reflects subjectively established research aims, but also inevitably directs research findings”. For instance in Figure 2, the choice of presenting, when relevant, the images of the two interlocutors side by side (e.g., images 5 and 6) was made because we felt that the detail of their facial expressions, smiles and micro-gestures within the same turn was necessary to understand minutely the adjustments that occurred during such an interaction. Such a transcription allows a vertical linear representation of turns and makes it possible to unpack the different modes at play “via a zigzagged reading” (Sindoni 2013: 82). Working iteratively on the transcript and on the accompanying text (see Table 2) helps refine both because they oblige researchers to give saliency to certain features in the transcript (simultaneousness of different phenomena, interaction between different semiotic modes, etc.), while the text that they write has to deploy textual resources to recount them. Neither the transcript nor the text can stand alone; rather they function as two faces of the proposed analysis.

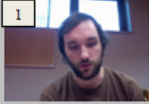
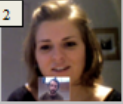
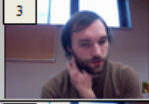

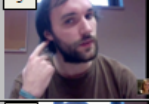
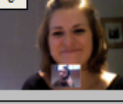
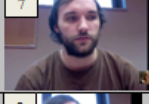
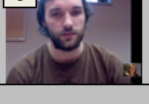
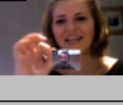

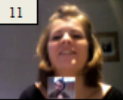
		LEARNER IMAGE	TEACHER IMAGE
1. Learner	The third young people is er a man too he has er really short hair and er <i>(he looks down at the photographs)</i> <i>(she is focused on the screen and produces a slight smile)</i>	1 	2 
	an ear piercing er:: like a:: (.) <i>(touches his ear while looking down)</i>	3 	
	<i>(looks up and looks at screen)</i> it's in like a:: (.) a tube <i>(makes a gesture to represent a round hole)</i>	4 	
	<i>(looks at the screen)</i> a bi:: er: <i>(points to his ear)</i>	5 	6 
	<i>(turns his face from the screen)</i> I don't know I [what I']	7 	
	<i>(looks at screen)</i>	8 	
			9 
	2. Teacher	[xx] (.) it's big/ <i>(mirrors learner's gesture (see 4) and looks at the screen with a smile)</i>	
3. Learner	Yeah it's big (.) it makes a hole in his ear:: <i>(touches his ear again and looks down)</i>	10 	
4. Teacher	OK <i>(nods and smiles)</i>		11 

Figure 2: Multimodal transcript of a word search episode

In turn 1, certain marks of hesitation, long pauses and self-admonishments ("I don't know") signal a communication breakdown while the learner is trying to find a way to describe the unknown lexical item. By touching his own ear repeatedly and miming a hole with his fingers, the learner is not only making his search visible to the teacher but is negotiating the meaning with her and looking for signs of her understanding. Her smile in image 6 suggests that she seems to understand what he is trying to describe, although he pursues his description in an attempt to be even more precise. As is visible in image 7, the student has what Goodwin and Goodwin (1986) would describe as a "thinking face", indicating to the teacher that he is still searching for the exact term, before he looks directly at the screen in image 8 – suggesting he wants confirmation from the teacher that she understands precisely what he is trying to describe. This search triggers a smile from the teacher and the mirror gesture (image 9) of that of the learner, which indicates that the teacher ratifies the description to a certain extent and that the interaction can continue while she is giving him her full attention by looking directly at the screen. Once the association of the verbal and nonverbal messages seem to have reached their objective, the learner verbally adds an element ("a hole") and gives redundant information by prodding his index finger at his ear again making sure that the teacher has understood the lexical item (she nods in image 11) even if the precise word has not been found.

Table 2: Textual analysis of the episode

There is no stabilized way of making multimodal transcripts although more and more researchers (see for instance Bezemer 2014; Flewitt et al. 2009; Norris 2004; Sindoni 2013) have devised astute ways of approaching this. Reading these authors, several considerations arise in relation to the units of analysis that can be selected, ethical dimensions that have to be attended to, the readability and the presentation of multimodal transcription.

First, turns of speech that constitute the conventional unit of analysis in Conversation analysis may not be as pertinent for multimodal analysis because, as noted by Flewitt et al. (2009: 45), "as soon as multiple modes are included, the notion of speech turns becomes problematic as other modes contribute meanings to exchanges during the silences between spoken turns". New units of analysis have thus to be devised to capture the specificity of multimodal interactions. For example, what is a speech turn when an individual uses written chat and speech simultaneously? Second, multimodal transcription makes participants identifiable, which makes it even more crucial to be vigilant about ethical considerations (as discussed above). Finally, researchers must establish a careful balance between the representation of all the features that are to be considered in a truly multimodal interaction and what a reader – even a seasoned one – is able to capture when confronted with a thick rendering of multimodality. As noted by Flewitt et al., "the perceptual difficulties for the audience of 'reading' genuinely multimodal transcription might outweigh the advantage of its descriptive 'purity'" (2009: 47). Eventually, there will be new ways of presenting multimodal data along with more traditional paper-based publication that will truly render the multimodal nature of such data. How to transform multimodal data in order to make them accessible with various degrees of complexity or presentational choices constitutes one direction for future research.

Drawing conclusions

Once transcriptions are completed, researchers can proceed to analyses such as the one proposed in Table 2. If their approach is *quantitative*, all the annotations can be exported to "applications that are able to perform statistics on the results" (Wittenburg et al. 2006).

Quantitative studies can thus give insight into a certain number of phenomena that can be relevant to understanding online learning and teaching. For instance, the number of pauses, the frequency of overlaps and the length of turns can shed light on the rhythm of a given interaction. The number of gestures and facial expressions produced by the participants could also give indications as to the communication potential of videoconferencing. The main outcome of quantitative studies concerns the identification of interactional patterns.

Although some examples of quantitative studies can be found, studies usually rely on qualitative approaches to data and focus on short episodes. As Bezemer says:

Making a transcript is an invaluable analytical exercise: by forcing yourself to attend to the details of a strip of interaction you gain a wealth of insights into the situated construction of social reality, including insights in the collaborative achievements of people, their formation of identities and power relations, and the socially and culturally shaped categories through which they see the world. (2014: X)

Yet, Adolphs and Varter point out that the community of researchers interested in multimodal analysis might profit from adopting a mixed approach and combining, when possible and pertinent, the conversation analysis of small samples of data with a corpus linguistics-based methodological approach. Thus, with the inclusion of large-scale data sets such an approach could extend “the potential for research into behavioral, gestural and linguistic features” (2013: 145).

Conclusion

In this chapter, we have shown the importance of taking into account the array of technologies that accompany the fabrication, analysis and transformation of interactional data. With ever refined software and transcription techniques, interactional linguistics has come to integrate into its agenda the intrinsically multimodal nature of interactions (Détienne and Traverso 2009). This is even more apparent when the interactions under study are themselves mediated by technologies, as is the case with videoconferencing-based exchanges. Technologies thus facilitate the gathering of interactional data and allow researchers to search them, replay them at will, annotate them with different degrees of granularity, visualize them from different perspectives, and structure them according to different scientific agendas (Erickson 1999). Not only do these technologies change the way researchers approach data, they also require them to develop new technical and methodological skills. As we have seen with the various steps involved in the collection, transcription and analysis of multimodal data, the different techniques at play mostly concern the representation of data. Each transformation of the data results in a new object that can be subject to yet another transformation, until the refinement is complete enough to yield a satisfactory comprehension of the phenomena under study. This points to the essential work of representations that “serve as resources for communicating and meaning-making” to the scientific community and beyond (Ivarsson, Linderöth and Saljö 2009: 201) and are “achieved by combining symbolic tools and physical resources” (ibid: 202).

The kinds of studies we have conducted not only help us to uncover the interplay of the different multimodal semiotic resources in online teaching environments but, ultimately serve to improve the design of teacher training programs (e.g., how to use the affordances of the

webcam in online interactions, how to pay attention to learner needs thanks to visual cues) so as to enhance learner computer interaction in online webcam-mediated exchanges.

References

- Adolphs, S. & Varter, R. (2013). *Spoken corpus linguistics. From monomodal to multimodal*. New York, NY: Routledge.
- Allwood, J., Ahlsén, E., Lund, J. & Sundqvist, J. (2007). Multimodality in own communication management. In J. Toivanen & P. Juel Henriksen (Eds.) *Current Trends in Research on Spoken Language in the Nordic Countries, Vol. II* (pp. 10-19). Oulu: Oulu University Press,.
- Baldry, A. & Thibault, P.J. (2006). *Multimodal Transcription and Text Analysis*. Equinox: London.
- Bétrancourt, M., Guichon, N., & Prié, Y (2011). Assessing the use of a Trace-Based Synchronous Tool for distant language tutoring. In H. Spada, G. Stahl, N. Miyake & N. Law (Eds.) *Connecting Computer-Supported Collaborative Learning to Policy and Practice: CSCL2011 Conference Proceedings*. Volume I — Long Papers (pp. 478-485).
- Bezemer, J. (2014). How to transcribe multimodal interaction?. In C.D. Maier & S. Norris (Eds.). *Texts, Images and Interaction: A Reader in Multimodality* (pp. 155-169). Boston: Mouton de Gruyter.
- Calibris, G. (2011). *Elements of Meaning in Gesture* (translated by Mary M. Copple). Amsterdam/Philadelphia: John Benjamins.
- Cohen, C. & Guichon, N. (2014). Researching nonverbal dimensions in synchronous videoconferenced-based interactions. Presentation at *CALICO Conference*, University of Athens, OH, USA.
- Cosnier, J. & Develotte, C. (2011). Le face à face en ligne, approche éthologique. In C. Develotte, R. Kern & M.-N. Lamy (Eds). *Décrire la conversation en ligne: Le face à face distanciel* (pp. 27-50). Lyon: ENS Éditions.
- De Chanay, H. (2011). La construction de l'éthos dans les conversations en ligne. In C. Détienne, F. & Traverso, V. (Eds.). (2009). *Méthodologies d'analyse de situations coopératives de conception*. Nancy: Presses Universitaires de Nancy.
- Develotte, R. Kern & M.-N. Lamy (Eds). *Décrire la conversation en ligne: Le face à face distanciel* (pp. 27-50). Lyon: ENS Éditions.
- Ellis, R. & Barkhuizen, G.P. (2005). *Analysing learner language*. Oxford: Oxford University Press.
- Erickson, T. (1999). Persistent conversation: an introduction. *Journal of computer-mediated communication*, 4(4), article 1. <http://jcmc.indiana.edu/vol4/issue4/ericksonintro.html>
- Flewitt, R., Hampel, R., Hauck, M. & Lancaster L. (2009). What are multimodal transcription and data? In C. Jewitt (Ed.). *The Routledge handbook of multimodal analysis* (pp. 40-53). London: Routledge.
- Guichon, N. & Cohen, C. (2014). The impact of the webcam on an online L2 interaction. *Canadian Modern Language Journal*, 70(3), 331–354.
- Guichon, N. & Wigham, C.R. (in review) Examining synchronous online language teaching from a semiotic perspective.
- Ivarsson, J., Linderöth, J. & Saljö, R. (2009). Representation in practices. A socio-cultural approach to multimodality in reasoning. In C. Jewitt (Ed.). *The Routledge handbook of multimodal analysis* (pp. 201-212). London: Routledge.
- Kern, R. G. (1995). Restructuring classroom interaction with networked computers: Effects on quantity and quality of language production. *The Modern Language Journal*, 79, 457-476.

- Kern, R.G. (2014). Technology as Pharmakon: The Promise and Perils of the Internet for Foreign Language Education. *The Modern Language Journal*, 98(1), 340-358.
- Kost, C. R. (2008). Use of communication strategies in a synchronous CMC environment. In S. Sieloff Magnan (Ed.). *Mediating discourse online* (pp. 153-189). Amsterdam: John Benjamins.
- Lamy, M.-N. & Hampel, R. (2007). *Online communication in language learning and teaching*. London: Palgrave Macmillan.
- McCafferty, S.G. (2008). Material foundations for second language acquisition: Gesture, metaphor and internalization. In S.G. McCafferty & G. Stam (Eds.). *Gesture: Second Language Acquisition and Classroom Research* (pp. 47-65). New York: Routledge.
- McCafferty, S.G. & Stam, G. (Eds.). (2008). *Gesture: Second Language Acquisition and Classroom Research*. New York: Routledge.
- McNeill, D. (1992). *Hand and mind: What gestures reveal about thought*. Chicago: The University of Chicago Press.
- MODE (2012). *Glossary of multimodal terms*. <http://multimodalityglossary.wordpress.com/>. 16 November 2014.
- Norris, S. (2004). *Analyzing multimodal interaction – A methodological framework*. New York: Routledge.
- Pellettieri, J., (2000). Negotiation in cyberspace: The role of chatting in the development of grammatical competence. In M. Warschauer & R. Kern. (Eds.). *Network-based Language Teaching: Concepts and Practices* (pp. 59–86). New York: Cambridge University Press.
- Satar, H.M. (2013). Multimodal language learner interactions via desktop videoconferencing within a framework of social presence: Gaze. *ReCall*, 25(1), 122-142.
- Siemens, G. & Baker, R.S.J.d. (2012). Learning Analytics and Educational Data Mining: Towards Communication and Collaboration. *Proceedings of the 2nd International Conference on Learning Analytics and Knowledge*.
- Sindoni, M.G. (2013). *Spoken and written discourse in online interactions: a multimodal approach*. New York: Routledge.
- Sloetjes, H. & Wittenburg, P. (2008). Annotation by category – ELAN and ISO DCR. In *Proceedings of the 6th International Conference on Language Resources and Evaluation (LREC 2008)*.
- Stockwell, G. (2010). Effects of multimodality in computer-mediated communication tasks, in M. Thomas & H. Reinders (Eds.). *Task-based language learning and teaching with technology* (pp. 83-104). London: Continuum.
- Tellier, M. & Stam, G. (2010). Stratégies verbales et gestuelles dans l'explication lexicale d'un verbe d'action. In V. Rivière (Ed.). *Spécificités et diversité des interactions didactiques* (pp. 357-374). Paris: Rivecourt Éditions.
- ten Have, P. (1999). *Doing conversation analysis: a practical guide*. London: Sage Publications
- Wang, Y. (2007). Task Design in Videoconferencing-Supported Distance Language Learning. *CALICO Journal*, 24(3), 591-630.
- Wittenburg, P., Brugman, H., Russel, A., Klassmann, A. & Sloetjes, H. (2006). ELAN: a Professional Framework for Multimodality Research. *Proceedings of Language Resources and Evaluation Conference*, 1556-1559.
- Yakura, E. K. (2004.) Informed Consent” and Other Ethical Conundrums in Videotaping Interactions. In P. Levine & R. Scollon (Eds.). *Discourse and technology – multimodal discourse analysis* (pp. 184-195). Washington DC: Georgetown University Press.

Software

Camtasia: <http://camtasia-for-mac.en.softonic.com>

ELAN: <http://tla.mpi.nl/tools/tla-tools/elan/>

Screen Video Recorder: <http://www.dvdvideosoft.com/fr/products/dvd/Free-Screen-Video-Recorder.htm#.VHBk8Eve5g0>

Skype: <http://www.skype.com>

SPSS: <http://www-01.ibm.com/software/fr/analytics/spss/>