

The CoMeRe French CMC corpora and their modeling in TEI

Thierry Chanier, Céline Poudat, Ciara Wigham

ird-cmc-rennes :

International Research Days: Social Media and CMC Corpora for the eHumanities
23-24h October 2015

CoMeRe (*Communication Médiée par les Réseaux*): a reference corpus of French CMC (2013-14)

<http://comere.org>

<http://hdl.handle.net/11403/comere>



Project supported by the national consortium *Corpus-écrits*, sub-part of *Huma-Num*, and *Ortolang*

- ❖ **People:** 14 research. from 8 research units. Coord: Chanier, T (Clermont), Poudat, C. & Sagot, B (Paris), Longhi, J. (Cergy), Antoniadis, G. (Grenoble)

Objective: Kernel corpus assembling existing corpora of different CMC genres and new corpora build on data extracted from the Internet. These heterogeneous corpora will be structured and processed in a uniform way, complemented with metadata. CoMeRe will be released as OpenData through the national infrastructure Ortolang, following constraints which will be reused for the forthcoming “*Corpus de Référence du Français*”.

Variety + Standards + Open Access

Variety + Standards + Open Access

SMS - cmr-smslareunion - cmr-smsalpes	Tweets - cmr-politweets	Email - cmr-simuligne	Text chat - cmr-getalp.org - cmr-favi - cmr-simuligne	Multimodal - cmr-copeas - cmr-tridem06
Wiki discussions - cmr-wikiconflits	Weblog - cmr-infral	Discussion forum - cmr-simuligne		Multimodal + 3D - cmr-archi21

Variety + Standards + Open Access

- ❖ People often wonder: "what did you choose the *Text Encoding Initiative* to encode multimodal interactions?
- ❖ These interactions can be viewed as text
 - BALDRY & THIBAULT (2006) consider "texts to be meaning-making events whose functions are defined in particular social contexts," following HALLIDAY (1989:10) "any instance of living language that is playing a role some part in a context of situation, we shall call it a text. It may be either spoken or written, or indeed in any other medium of expression that we like to think of."
- ❖ Mainstream of oral corpora are encoded into TEI
- ❖ TEI offers a very rich way to describe the project corpus (on top of the interactions set)
- ❖ Opportunity to work at a European level

Variety + Standards + Open Access

- ◊ **“Availability and Access:** the data must be available as a whole and at no more than a reasonable reproduction cost, preferably by downloading over the internet. The data must also be available in a convenient and modifiable form.
- ◊ **Reuse and Redistribution:** the data must be provided under terms that permit reuse and redistribution including the intermixing with other datasets. The data must be machine-readable.
- ◊ **Universal Participation:** everyone must be able to use, reuse and redistribute – there should be no discrimination against fields of endeavor or against persons or groups. For example, ‘non-commercial’ restrictions that would prevent ‘commercial’ use, or restrictions of use for certain purposes (e.g. only in education), are not allowed. “OpenDefinition.org



Example of CoMeRe licences

- ❖ Falaise, A. (2014). *Corpus de français tchaté getalp.org*. [cmr-getalp_org].



- ❖ Antoniadis, G (2014). *Corpus de SMS réels dans les Alpes, smsalpes* . [cmr-smsalpes].



- ❖ Longhi, J., Marinica, C., Borzic, B., Alkhouli, A. Polititweets. (2014). *Corpus de tweets provenant de comptes politiques influents*. [cmr-polititweets]



- ❖ Ledegen, G. (2014). *Grand corpus de sms smslareunion* . [cmr-smslareunion]



- ❖ Yun, H. & Chanier, T. (2014). *Corpus d'apprentissage FAVI (Français académique virtuel international)*. [cmr-favi].



- ❖ Abendroth-Timmer, D., Bechtel, M., Chanier T. & Ciekanski, M. (2014). *Corpus d'apprentissage INFRAL (Interculturel Franco-Allemand en Ligne)*. [cmr-infral]



- ❖ Reffay, C. Chanier, T. Lamy, M.-N. & Betbeder, M.-L. (2014) *Corpus d'apprentissage Interactions Simuligne (Simulation en ligne en apprentissage des langues)*. [cmr-simuligne]



Open Resources and TOols for LANGuage

ORTOLANG is an EQUIPEX project accepted in February 2012 in the framework *d'avenir*. Its aim is to construct a network infrastructure including a repository (corpora, lexicons, dictionaries etc.) and readily available, well-documented Expected outcomes comprise:

- promoting research on analysis, modelling and automatic processing at the highest international level thanks to effective resource pooling;
- facilitating the use and transfer of resources and tools set up with industrial partners, notably SMEs which often cannot develop such language processing given the cost of investment;
- promoting French language and the regional languages of France, acquired by public laboratories.



ORTOLANG is a service for the language, which is complemented by the services offered by [Huma-Num](#) (very large research infrastructure).

The screenshot shows the ORTOLANG website interface. At the top, there is a navigation bar with the title "ORTOLANG (beta)" and a search bar labeled "Search for corpora". Below the navigation bar is a sidebar with links: Home, Corpora (which is highlighted in blue), Integrated Projects, Tools, Lexicons, and Information. The main content area displays a grid of corpora entries. One entry, "CoMeRe (Communication médiée par les réseaux)", is circled in red. The grid includes other entries like MC4, LIDILEM, and L'EST. A red arrow points from the top left towards the sidebar.

Corpus	Description	Published on
MC4	Modélisation Contrastive et Computationnelle des Chaînes de Coréférence	Published on 2015-06-15
LIDILEM	Littéracie Avancée	Published on 2015-05-01
CoMeRe (Communication médiée par les réseaux)	Salut s que NOM_4> c dc d pr sa	Published on 2015-05-01
L'EST	TGOF	Published on
Corpus 14	S	Published on

Corpora repository in ORTOLANG

<http://hdl.handle.net/11403/comere>

Cuurent list of corpora

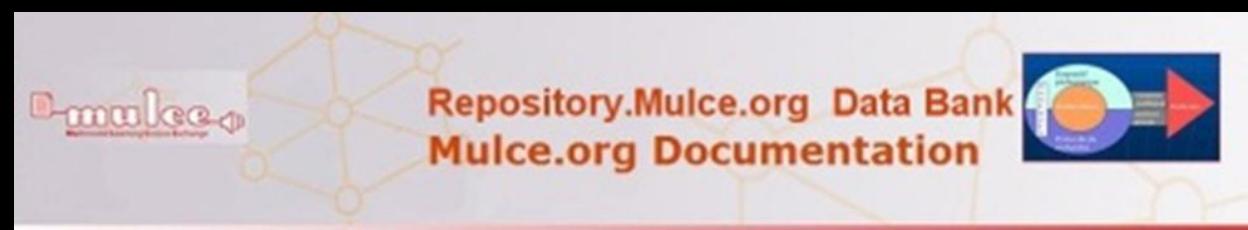
- ❖ 1) Antoniadis, G (2014). *Corpus de SMS réels dans les Alpes, smsalpes* [corpus]. In Chanier T. (ed.) Banque de corpus CoMeRe. Ortolang.fr : Nancy. [<http://hdl.handle.net/11403/comere/cmr-smsalpes>]
- ❖ 2) Falaise, A. (2014). *Corpus de français tchaté getalp_org* [corpus] . In Chanier T. (ed) Banque de corpus CoMeRe Banque de corpus CoMeRe. Ortolang.fr : Nancy. [http://hdl.handle.net/11403/comere/cmr-getalp_org]
- ❖ 3) Ledegen, G. (2014). *Grand corpus de sms SMS La Réunion* [corpus]
- ❖ 4) Reffay, C. Chanier, T. Lamy, M.-N. & Betbeder, M.-L. (2014). *Corpus Interactions Simuligne (Simulation en ligne en apprentissage des langues)* [corpus]...
- ❖ 5) Yun, H. & Chanier, T. (2014). *Corpus d'apprentissage FAVI (Français académique virtuel international)* [corpus...]
- ❖ 6) Abendroth-Timmer, D., Bechtel, M., Chanier T. & Ciekanski, M. (2014). *Corpus d'apprentissage INFRAL (Interculturel Franco-Allemand en Ligne)*. [corpus]...
- ❖ 7) Longhi, J., Marinica, C., Borzic, B. & Alkhouli, A. (2014) *Corpus de tweets provenant de comptes politiques influents*. [corpus]...
- ❖ 8) Chanier, T. & Audras, I. (2015). Tridem06 corpus: intercultural competence in online exolingual group exchanges [...]
- ❖ 9) Chanier, T. & Wigham, C.R. (2015). Archi21 corpus: collaborative language and architectural learning in Second Life [...]
- ❖ 10) Chanier, T., Reffay, C., Betbeder, M-L., Ciekanski, M. & Lamy, M-N. (2015). Copéas corpus: online language learning within an audiographic environment [...]
- ❖ 11) Poudat,C., Grabar , N. Kun, J. & Paloque-Berges, C. (2015). Corpus wikiconflits, conflits dans le Wikipédia francophone [...]

Corpora composed of verbal acts

Ref	Tokens	Partici.	Posts	Envir.	
(Antoniadis, 2014)	449 313	359	22 052	SMS	→ Informal business
(Falaise, 2014)	35 M	25 000	3 M	textchat	→ Informal
(Ledegen, 2014)	357 000	850	22 000	SMS	→ Informal
(Reffay et al., 2014)	600 000	67 + 4 groups	- textchat: 6 790 - emails: 2 030 - forums: 2 686	LMS	→ education
(Yun, Chanier, 2014)	77 605	31 + 2 courses	7 750	textchat	→ education
(Abendroth-Timmer et al., 2014)	273 546	26 + 4 groups	1 200	Blog	→ education
(Longhi, Marinica, 2014)	567 851	205	34273	Tweet	→ politics
(Poudat et al., 2015)	489 000	3971	4456	Wiki discussions	→ science

verbal & non-verbal acts (LETEC corpora)

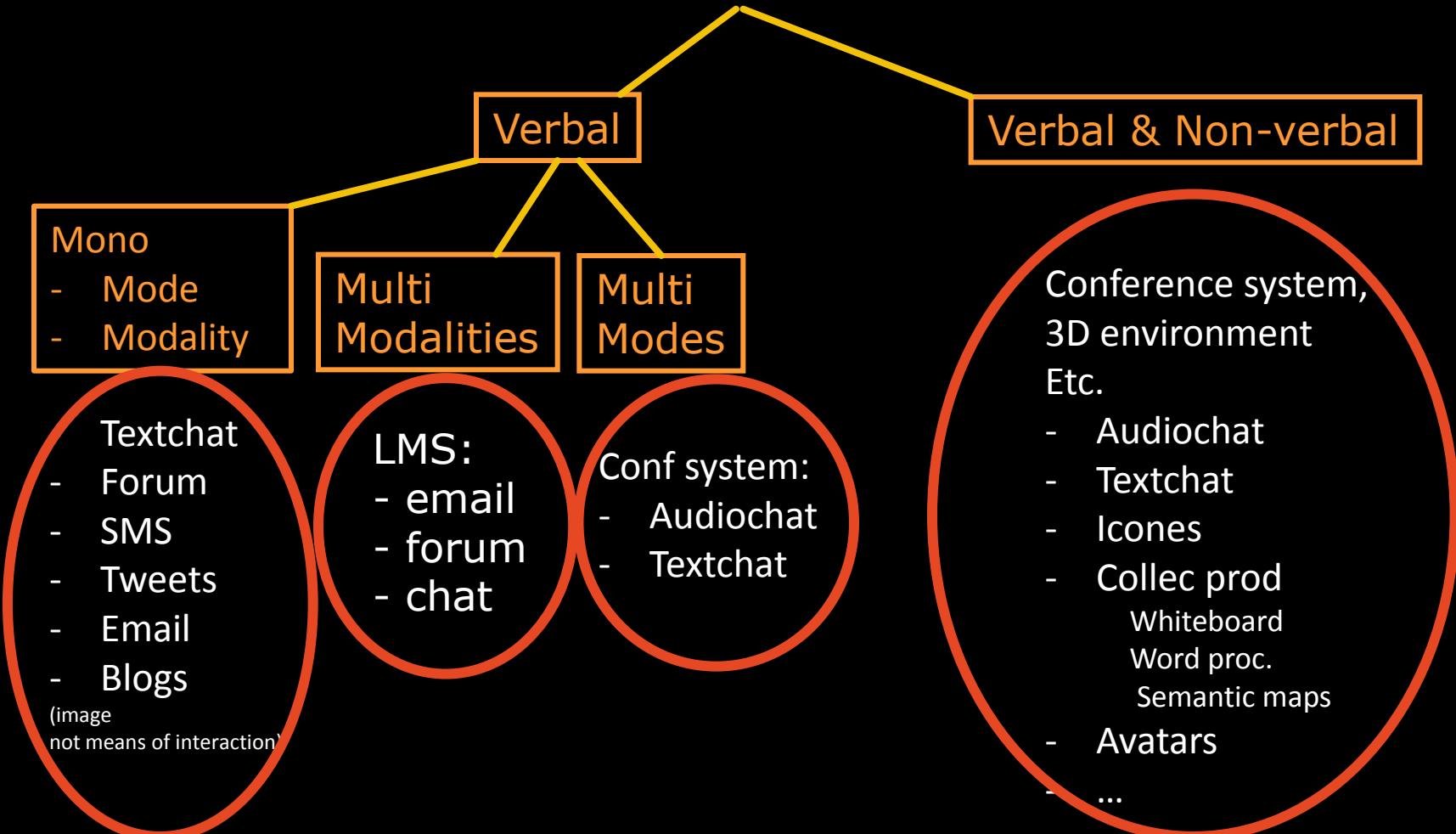
Ref	Tokens	Partici.	Posts, U, Prod	Envir.
(Chanier & Audras, 2015)	184 594	62 + 12 groups	- audio: 2 809 - chat: 248 - non-verbal: 1 058 - blog: 779	Conference system
(Chanier & Wigham, 2015)	27 912	18 + 4 groups	- audio: 1 690 - chat: 669 - non-verbal: 2 452	3D environment
(Chanier , Reffay et al., 2015)	127 228	16 + 2 groups	- audio: 7 718 - chat: 1 566 - non-verbal: 5 790	Conference system



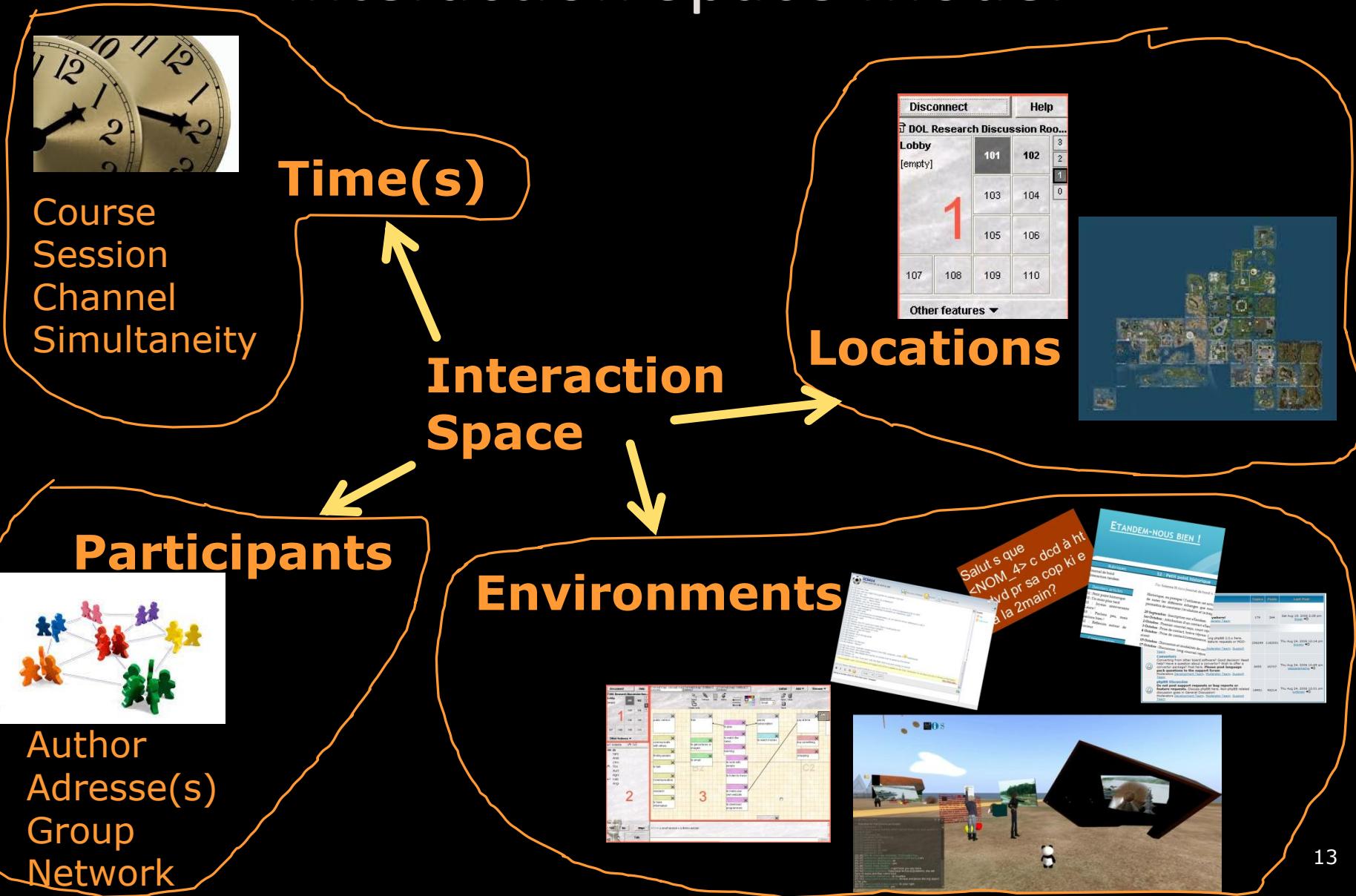
Interaction Space Model

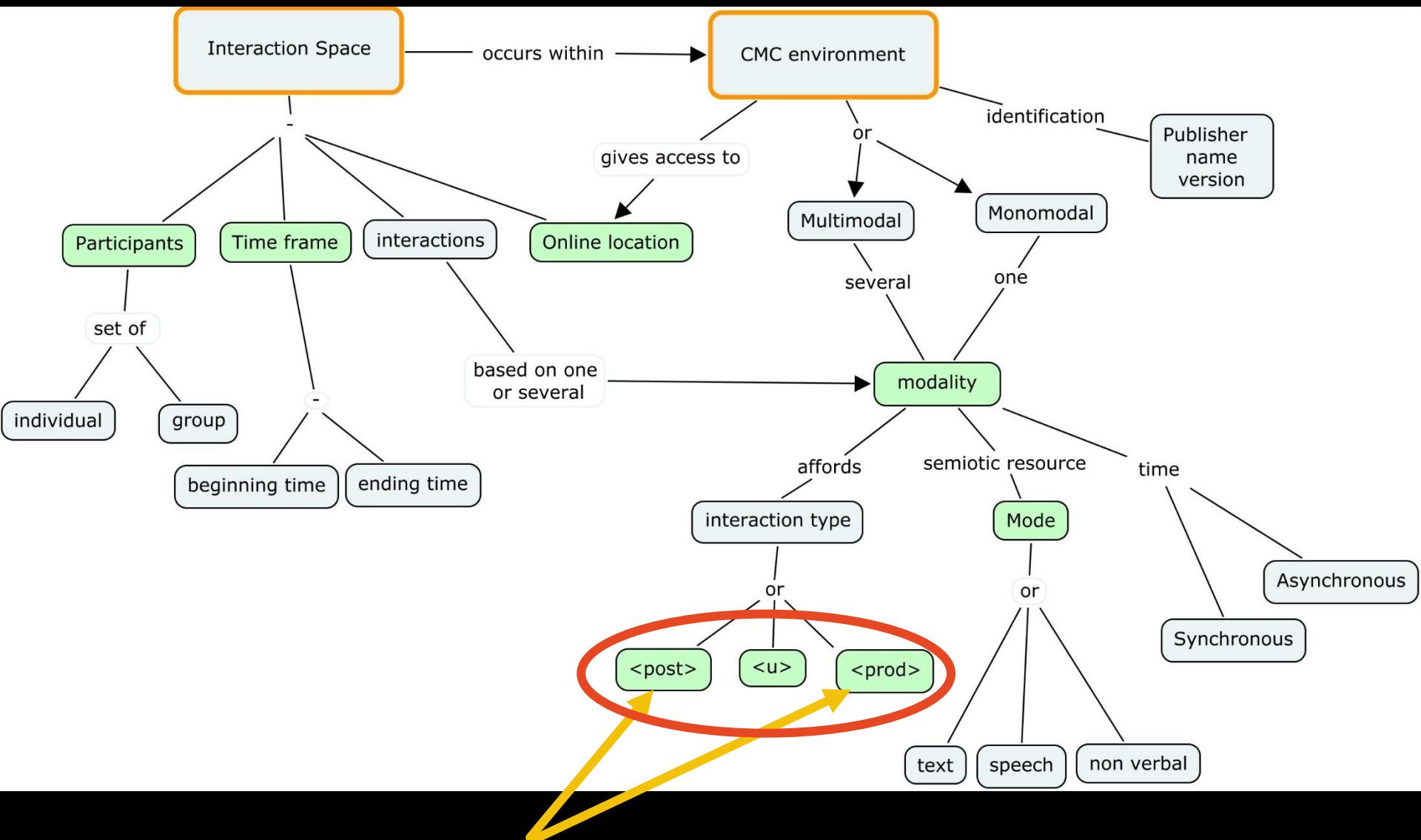
Implementation in CoMeRe corpora

Environments



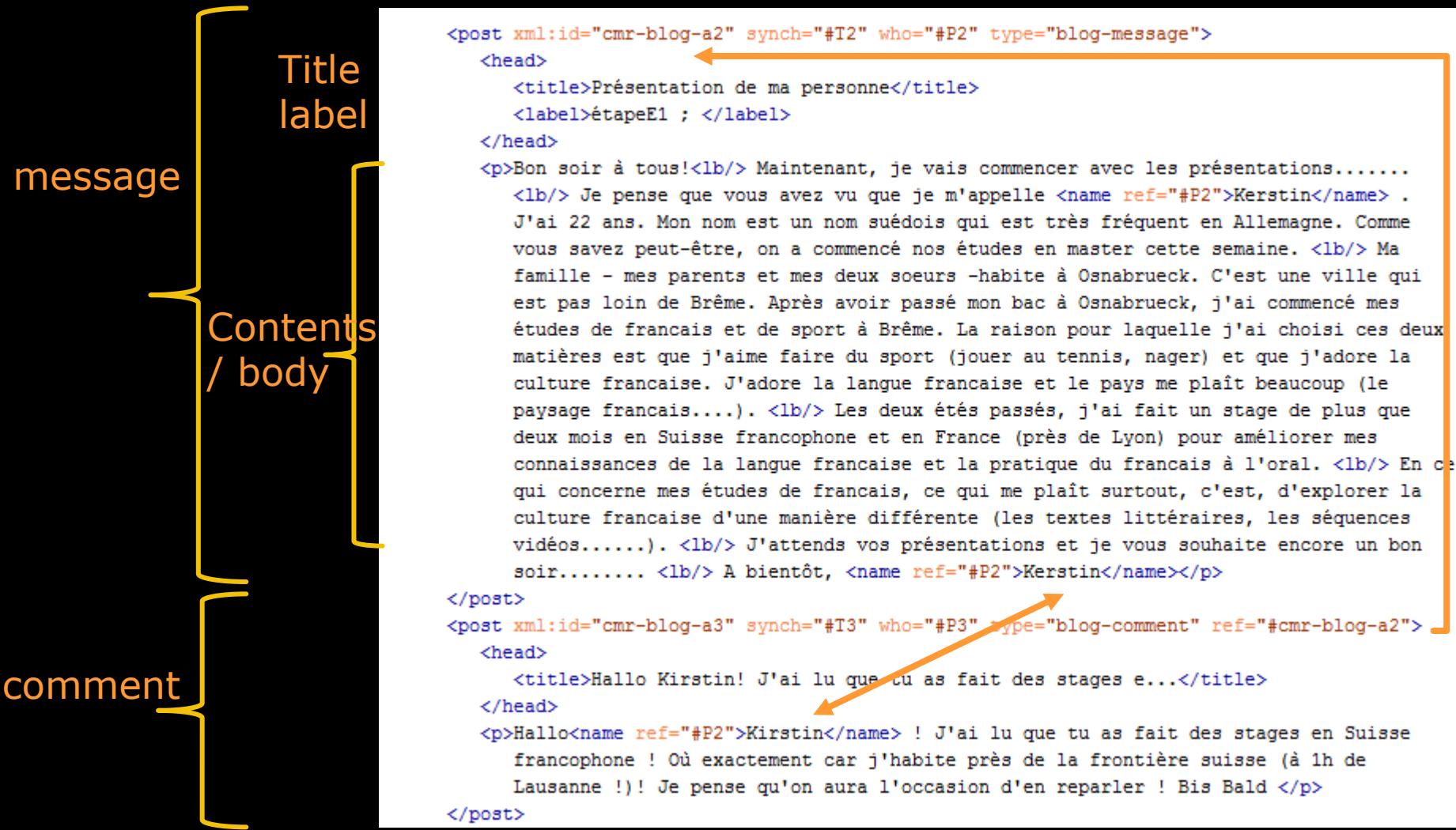
Interaction Space Model





New macro-level elements

Blog: message and comment



More complex (cumbersome?): email (here), forum

Response
to what?

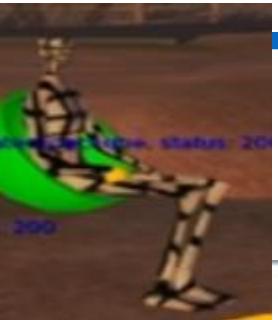
Sent to
whom?
Read by
whom?

May
contain
HTML,
Table,etc.

Attached
doc

```
<post xml:id="#cmr-Simu-Aq-At-Outbox-0004" when="2001-05-01T10:26:00" who="#cmr-Simu-At" type="email-message" ref="#cmr-Simu-Aq-At-Inbox-0003">
    <head>
        <title>mon bio</title>
        <listPerson>
            <person corresp="#cmr-Simu-A12">
                <event type="SendTo" />
                <label>SendTo</label>
            </person>
            <person corresp="#cmr-Simu-A12">
                <event type="Read" when="2001-05-01T10:26:00">
                    <label>Read</label>
                </event>
            </person>
        </listPerson>
    </head>
    <p>Bonjour<name ref="#cmr-Simu-A12"
type="person"><forename>Bruce</forename></name>, .et merci de m'avoir éclairé
sur la nature de votre travail. Il ne me semble pas moins passionnant pour
autant, et cela doit être une tâche infinie d'anticiper sur les causes des
pannes possibles... il doit y avoir tellement de possibilités. A propos des
accents, il me semble que les personnes qui utilisent les claviers anglais ont
du mal à en mettre dans le courriel. Cela dit, il n'est pas interdit d'ouvrir un
fil de discussion sur les accents dans le forum, je suis sûre que vous aurez du
succès, car il me semble que c'est une question qui préoccupe tout le monde !
Amicalement, <name ref="#cmr-Simu-At"
type="person"><forename>Anna</forename></name>
    </p>
    <trailer>
        <ref type="attached_file">symbols.xls</ref>
    </trailer>
</post>
```

Modality interplay



1.5 mn video

tingrabu

ok for me this presentation was become too fast because it's always the same in our architectural school we have not time and too quickly sorry and we can't do good images because it's less time euh I don't know [...] and it's a big matter because we always talk about teleportation [...] an everyday lack of time ok thank you quentinrez and this is very difficult [...]

tfrez2

it went too quickly
or it was too early in the week?
ok
too quickly means you didn't have enough time

romeorez

i think it was to early too

quentinrez

yes, it's an everyday lack of time

* Paper: (Wigham & Chanier, 2013) CALL journal

* Data: (Wigham, 2013) LETEC corpus



Multimodalité : Verbal et non verbal

Table 3. Classification of communication acts in *Second Life*.

Communication mode	Communication modality	Act type and transcription code	Explanation
Verbal	Audio (voice chat)	Audio act (aud)	Verbal act in the full duplex public audio channel
	Text chat	Silence (sil) Text chat act (tc)	Interval between two audio acts greater than three seconds message entered in the public text chat window
Non-verbal	Proxemics	Movement (mvt)	Avatar movement in the environment, e.g. avatar sits down, flies, walks backwards
		Entrance into/exit from the environment (eex)	Avatar enters or exits the synthetic world
	Kinesics	Kinesic (kin)	Avatar gestures and movements made by an avatar's body part, e.g. nod, point, clap
	Production	Production (prod)	Production or display of an object in the <i>Second Life</i> environment

(Wigham & Chanier, 2013)

Audio

kinesics
chat
chat
chat
chat
chat
chat
chat

Modality interplay

```
<u xml:id="a191" who="#tingrabu" start="#ts373" end="#ts430">ok hm for me this presentation was
hm <pause dur="PT1S"/> become too fast because it's always the same in our
architecture school euh we have not time and hm <pause dur="PT1S"/> too
quickly sorry and hm <pause dur="PT1S"/> we can't do good images because euh
[...] may be I don't know <vocal> <desc>chuckles</desc></vocal></u>
<prod xml:id="a192" who="#romeorez" start="#ts376" end="#ts377" type="body" subtype="kinesics">
<code>eat(popcorn)</code></prod>
<post xml:id="a195" who="#tfrez2" start="#ts380" end="#ts381" type="chat-message">
<p>it went too quickly?</p></post>
<post xml:id="a197" who="#tfrez2" start="#ts384" end="#ts385" type="chat-message">
<p>or it was too early in the week?</p></post>
<post xml:id="a200" who="#romeorez" start="#ts392" end="#ts393" type="chat-message">
<p>i think it was to early</p></post>
<post xml:id="a203" who="#tfrez2" start="#ts396" end="#ts398" type="chat-message">
<p>too early ok</p></post>
<post xml:id="a204" who="#tfrez2" start="#ts399" end="#ts401" type="chat-message">
<p>too quickly means that you didn't have enough time to speak</p></post>
<post xml:id="a207" who="#quentinrez" start="#ts405" end="#ts406" type="chat-message">
<p>yes, it's an everyday lack of time</p></post>
```

Detailing the corpus project in TEI

To support reuse by other researchers

Using TEI header



Archi21 corpus: collaborative and architectural learning in Life

This page: <http://hdl.handle.net/11403/comere/cmr-archi21/cmr-archi21-tei-v1>

Back to corpus main page: <http://hdl.handle.net/11403/comere/cmr-archi21>

Download the TEI file: <http://hdl.handle.net/11403/comere/cmr-archi21/cmr-archi21-tei-v1.xml>

- [Overview](#)
- [Rationale for this corpus](#)
- [Description of the Interaction Space](#)
- [Extracts of Interactions](#)
- [Publication Statement and Rights](#)

How to cite this resource

Chanier, T. & Wigham, C.R. (2015). Archi21 corpus: collaborative language and architectural learning in S [cmr-archi21-tei-v1 ; <http://hdl.handle.net/11403/comere/cmr-archi21/cmr-archi21-tei-v1>]

Overview of the corpus

The first version of this corpus, under the LETEC standard - corpus for learning -, (Chanier, T. & Wigham,

Multimodality environment: general features and affordances (LMS- Learning Management System)

LMS

textchat

email

forum

```
<classDecl>
  <taxonomy>
    <category xml:id="WebCT">
      <catDesc>WebCT Online Learning Management System, version 10.1.1. It is a communication tool (all textual modalities are described).
      <category xml:id="ActivityStructure">
        <catDesc>The Interaction Spaces of one learning group or one learning activity (hence <gi>div</gi> of ActivityStructure).
        <category xml:id="chat">
          <catDesc>Interactions (sets of elements <gi>post</gi> of textchat) are organized in chatrooms (<gi>div</gi> level), which are of different types of posts.
            <textDesc xml:lang="en-GB">
              <channel mode="w" xml:lang="en-GB"><term>textchat</term>
              <constitution>Messages typed by participants in a chatroom.</constitution>
              <derivation type="original"/>
              <domain type="education"/>
              <factuality type="fact"/>
              <interaction type="complete" active="plural" passive="many"><note>Synchronous discussion tool. All members of a chatroom. Chatrooms can only be opened and closed by tutors and teachers.</note></interaction>
              <preparedness type="spontaneous"/>
              <purpose degree="high"><note>related to the group activity</note></purpose>
            </textDesc>
          </catDesc>
        <category xml:id="chat-message">
        <category xml:id="chat-event">
          <category xml:id="connexion">
            <catDesc>participant enters into the textchat room</catDesc>
          <category xml:id="deconnection">
            <catDesc>participant leaves the textchat room</catDesc>
          </category>
        <category xml:id="email">
          <catDesc>The email is a communication tool integrated into the LMS. Email messages are in a (<gi>div</gi>).</catDesc>
          [...]
        </catDesc>
      </category>
      <category xml:id="email-message">
        <catDesc>if a message (<gi>post</gi>) has a <att>parent</att>, it is a response to a previous one. More details, see <gi>tagsDecl</gi>.</catDesc>
      </category>
    <category xml:id="forum">
      <catDesc>A discussion forum has a title and is organized in threads (set of posts).</catDesc>
    </category>
  </taxonomy>
</classDecl>
```

CoMeRE model (ODD)

The screenshot shows two pages from the TEI wiki, both highlighted with a yellow border.

SIG:Computer-Mediated Communication

- navigation
 - Main Page
 - TEI website
 - idleTalk
 - Current events
 - Recent changes
 - Random page
 - Help
- search
 - Go
 - Search
- toolbox
 - What links here
 - Related changes
 - Upload file
 - Special pages

Motivation

In the past three decades, computer-mediated communication (CMC) has become a major part of our daily lives. Even though there's been a lot of research and development in this field, there is still no common standards for the representation of CMC in digital corpora. This page aims to provide a basic schema for representing CMC in TEI, based on the existing work in the field.

SIG:CMC/Draft: A basic schema for representing CMC in TEI

- navigation
 - Main Page
 - TEI website
 - idleTalk
 - Current events
 - Recent changes
 - Random page
 - Help
- search
 - Go
 - Search
- toolbox
 - What links here
 - Related changes
 - Upload file
 - Special pages
 - Printable version
 - Permanent link

Contents [hide]

- Status of this draft
- Interaction types
 - Interaction
- Examples of macrostructures
 - Multimodal example
 - Textchat
 - SMS
 - Discussion forum
 - Wikipedia discussion
 - Blog
 - Email
 - Tweets
- Macrostructure: the `<post>` element
 - The element
 - Attributes for `<post>`
- Macrostructure & multimodality: `u` and `prod` elements
 - Types of divisions for the interaction space
 - `<prod>` element
 - Elements than cannot be used as an act of type `prod`
 - Attributes of `<u>` and `<prod>`

Many more examples here

http://wiki.tei-c.org/index.php/SIG:CMC/Draft:_A_metadata_schema_for_CMC

http://wiki.tei-c.org/index.php/SIG:CMC/CoMeRe_schema_draft_forRepresenting_CMC_in_TEI_2014_29



CoMeRe team

Documentation and events : <http://comere.org>
Repository: <http://hdl.handle.net/11403/comere>