

# Investigations in Collaborative Multi-Party Discourse

Andrew J. Cowell, Michelle L. Gregory  
Pacific Northwest National Laboratory

901 Battelle Blvd.

Richland WA 99354

509-375-4548

{andrew.cowell; michelle}@pnl.gov

## Abstract

In this paper, we discuss the efforts underway at the Pacific Northwest National Laboratory in understanding the dynamics of multi-party discourse across a number of communication modalities, such as email, instant messaging traffic and chat data. Only by understanding how individuals communicate through these new media technologies can we hope to successfully design and implement the social media applications of the future. Two prototype systems are discussed: (1) the Conversation Analysis Toolkit (ChAT) is a set of experimental computational linguistic components that enables users to easily identify topics or persons of interest within multi-party conversations, including who talked to whom, when, the entities that were discussed, etc., and (2) the Retrospective Analysis of Communication Events (RACE) application, leveraging many of the ChAT components, is an application built specifically for knowledge workers and focuses on merging different types of communications data so that the underlying message can be discovered in an efficient, timely fashion.

## Keywords

Anthroposemiotics, Discourse Analysis, Sentiment Analysis, Group Dynamics, Conversation Topic Segmentation.

## 1. Introduction

Computationally supported modalities such as email and instant messaging have had immeasurable effect on the way we work, play and interact with those in our lives. Being able to understand how individuals communicate, the methods they use, their personal preferences, etc. are all part of a field called anthroposemiotics which looks to uncover the mystery behind how we communicate with ourselves (intrapersonal communication), with others (interpersonal communication), within groups (group dynamics) and across cultures (cross-cultural communication). While a great deal of literature exists in each of these fields, there are few operational prototypes that allow for true hands-on investigation. Here, we discuss two projects underway within the Rich Interaction Environments group of the Pacific Northwest National Laboratory aimed at building analytic devices that not only uncover content in messages, but various aspects of anthroposemiotics as well.

## 2. The Conversation Analysis Toolkit (ChAT)

The ability to extract and summarize content from data is a fundamental goal of computational linguistics. As such, a number

of tools exist to automatically categorize, cluster and extract information from documents. However, these tools do not transfer well to data sources that are more conversational in nature, such as multi-party meetings, telephone conversations, email, chat rooms, etc. Given the plethora of these data sources, there is a need to be able to quickly and accurately extract and process pertinent information without having to cull them manually.

In this paper, we present a Conversational Analysis Toolkit (ChAT) that consists of several language processing tools (topic segmentation, affect scoring, named entity extraction) that can be used to automatically annotate conversational data. The processing components have been specially adapted to deal with conversational data. The versatile language processing components in ChAT were developed in modular, open designs so that they can be used independently or be integrated into other analytic tools (such as RACE, discussed later). In addition, all the processing components are domain and source independent, e.g., the topic segmentation does not rely on features specific to a dataset, such as acoustic information from transcripts. Also, all processing components have been built as independent plug-ins to the external processing engine: the input of one does not rely on the output of the others. This lets users choose to include or exclude various processes to suit their needs or even exchange the components with new tools.

### 2.1 ChAT components

The processing components of ChAT include topic segmentation, named entity extraction, sentiment analysis, and a participant role identifier. The named entity extraction is a module based on LinPipe with no special modifications made to accommodate conversational data.

For topic segmentation we employ a windowless method (WLM) for calculating a suitable cohesion signal that does not rely on a sliding window to achieve the requisite smoothing for an effective segmentation, as lexical cohesion methods do. Instead, WLM employs a constrained minimal-spanning tree (MST) algorithm to find and join pairs of elements in a sequence. Of particular interest for our research is the success of WLM on sparse, conversational data. We evaluated WLM's performance on the ICSI meeting corpus by comparing our segmentation results to the results obtained by implementing a lexical cohesion method, LCSeg. Using the 25 hand-segmented meetings, our algorithm achieved a significantly better segmentation for 20 out of 25 documents.

In addition to topic and entity extraction, conversations can also be analyzed by who participated in them, their relationship to one another and their attitude toward topics they discuss. In an initial

attempt to capture participant attitude, we have included a lexically based sentiment analysis component. We employed the General Inquirer (GI). Every utterance is scored for the number of and type affect words it contains. We make use of this data through tracking affect over topic as well as affect by participant. In addition to affect, we also use basic conversational statistics, such as number of words, number of utterances, proportion of questions to statements, etc. to help provide insight to participant roles. For example, in the meeting corpus, we find only a couple of participants ask questions and participate in every topic discussed. We can infer from this that they are the meeting leaders, whereas those participants who only contribute utterances during a specific topic might be experts on that topic.

Through these components, ChAT enables detailed analysis of any kind of data that has multiple participants and sparse content. It allows users to easily abstract patterns from the data that might otherwise be difficult to identify.

### 3. RACE

In RACE, the goal is for a user analyst to effectively look for both episodic and social information across the topics of a multitude of conversations in a variety of modalities. RACE integrates email, instant messaging, text messaging, in-person meetings, phone conversations and teleconferences in addition to chat and newsgroup participation. The goal is to get a more holistic sense of individuals throughout their discrete conversations and communication methods. As a post-hoc analysis tool, RACE aids users by adding system interpretations of affect and social dynamics to the other thread or conversation visualization systems. Content-driven interpretations of group dynamics, affect and social role complement full-text transcripts of the conversations, providing shortcuts to insight.

#### 3.1 RACE components

RACE was developed according to functional requirements outlined by four information analysts. Based on these and the data types RACE supports, we designed an environment can run on three screens simultaneously, be split across three panes (useful for performing analysis on large displays like wall-mounted plasma displays) or on a single screen with the use of a window manager seen in the top right of each view.

The first pane, the corpus view, shown on the far left in Figure 1, allows a user to see all of the conversations in their information space clustered by topic. Each dot in the pane represents a different conversation, email, blog entry, or documents. The

various modalities are represented by different colors or icons, and users can select documents by participant names, modalities, or topics in this pane.

The selected documents appear in the middle pane, the sequence view. Here is where a user will review in detail, a small subset of conversations that they found of interest in the corpus space. Each conversation has an independent time line and can be zoomed out to show the entire conversation or zoomed in to see the individual utterances (these may also be accessed using tool-tips). The conversation titles on the left side of the screen can be expanded to show all the participants involved. Clicking a participant opens a dialog box containing known information about that individual (including any known aliases and other names he or she may use online). A global timeline at the bottom of the screen shows where each conversation falls in sequence.

Once a conversation of interest is identified through the triage process, it can be selected for deeper investigation in the details view in the third pane. This view can enable the analyst to see beyond the individual utterances. Using ChAT components, the details view lets the analyst gain insight into an individual's opinion on the topics discussed. The transcript is color-coded to show the seven dimensions of affect (expression, power, ethics, attainment, skill, accomplishment and transactions), while a graph representation allows the analyst to compare individuals' affect against each other. To ingest the text in different ways, a 'text-to-speech' has been integrated for aural ingest so that one can listen to a conversation while visualizing other aspects of the conversation at the same time. As it steps through the utterances, a group dynamics graphic (based on Erickson's Social Proxy) shows how the individuals relate to each other, highlighting those involved in the conversation and those that are idle. This view also provides a hierarchical view of the topics discussed with the ability to trigger a multi-dimensional visualization that maps participants to topics.

### 4. Conclusion

In this paper, we have presented two prototype systems designed to specifically to deal with multiple forms of conversational data. These prototypes allow users to investigate the content of conversational data in novel ways for quick insight into the information contained in them. What makes these systems unique is that not only can content of the messages be revealed for timely analyses, these systems integrate aspects of anthroposemiotics which uncovers a new layer of information about the participants, their roles, and how they interact.

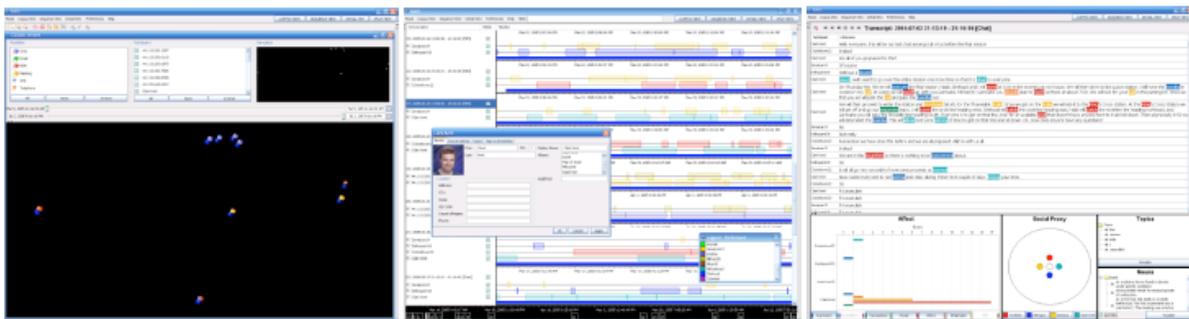


Fig. 1: Three pane RACE display