

# GCA Analyzer: A Python Package for Group Conversation Analysis

Jianjun Xiao<sup>1</sup> and Cixiao Wang<sup>1</sup>

<sup>1</sup> Research Centre of Distance Education, Beijing Normal University, Beijing, People's Republic of China

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Open Journals](#)

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

Group conversation analysis is crucial for understanding social dynamics, learning behaviors, and communication patterns in various settings, particularly in educational contexts. The GCA Analyzer is a Python package that implements a comprehensive set of metrics and methods for analyzing group conversations, with special emphasis on Chinese text processing capabilities. This tool provides quantitative measures for participation patterns, interaction dynamics, content novelty, and communication density, making it especially valuable for researchers in education, social psychology, and communication studies.

The GCA Analyzer builds upon foundational work in group cognition analysis and computational approaches to group communication (Dowell et al., 2019) (Wang & Xiao, 2025). These works provide essential frameworks for understanding how participants interact and contribute in group discussions.

## Statement of Need

The analysis of group conversations is crucial in various fields, particularly in educational research, organizational behavior studies, and online learning environments. While several tools exist for conversation analysis, there is a significant gap in tools that can effectively handle multilingual text and provide comprehensive interaction metrics. Existing tools like ...

The GCA Analyzer addresses these gaps by providing:

1. Robust participation analysis through participation matrices
2. Temporal interaction analysis using sliding windows
3. Content similarity and novelty metrics
4. Social impact and responsivity measurements
5. Visualization capabilities for interaction patterns

These features enable researchers to conduct detailed analyses of group conversations, MOOC interactions, and cross-cultural communication patterns, supporting both research and practical applications in various educational and social contexts. By providing a comprehensive toolkit for quantitative analysis of group dynamics, the GCA Analyzer facilitates deeper insights into collaborative learning processes, team communication effectiveness, and the evolution of ideas within group discussions.

## Installation

Install GCA Analyzer using pip:

```
pip install gca-analyzer
```

## Quick Start

Here's a simple example to analyze a group conversation:

```
from gca_analyzer import GCAAnalyzer

# Initialize the analyzer
analyzer = GCAAnalyzer()

# Load and analyze data
metrics = analyzer.analyze_conversation('conversation_1', data)
print(metrics)
```

## Command Line Usage

```
python -m gca_analyzer --data your_data.csv
```

## Input Data Format

The input data should be a CSV file with the following columns:

- conversation\_id: Identifier for the conversation
- person\_id: Identifier for each participant
- text: The actual message content
- time: Timestamp of the message

## Configuration Options

Command line arguments:

- data: Path to input data file (required)
- output: Output directory for results (default: gca\_results)
- best-window-indices: Window size optimization threshold (default: 0.3)
  - Range: 0.0-1.0
  - Sparse conversations may benefit from smaller thresholds
- console-level: Logging level (default: INFO)
  - Options: DEBUG, INFO, WARNING, ERROR, CRITICAL
- model-name: LLM model for text processing
  - Default: sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2

## Example Results

After running the analysis, you'll get metrics including:

- Participation patterns
- Internal cohesion
- Overall responsivity
- Social impact
- Content newness
- Communication density

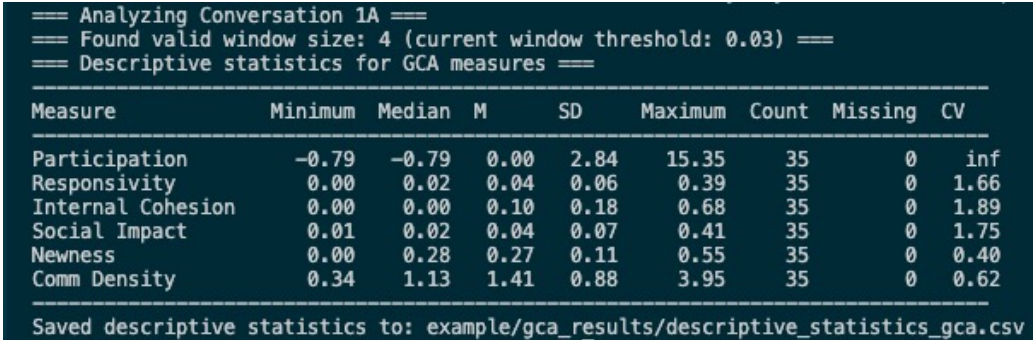


Figure 1: Descriptive Statistics for GCA Measures

64 You'll get interactive and informative visualizations for key GCA measures:

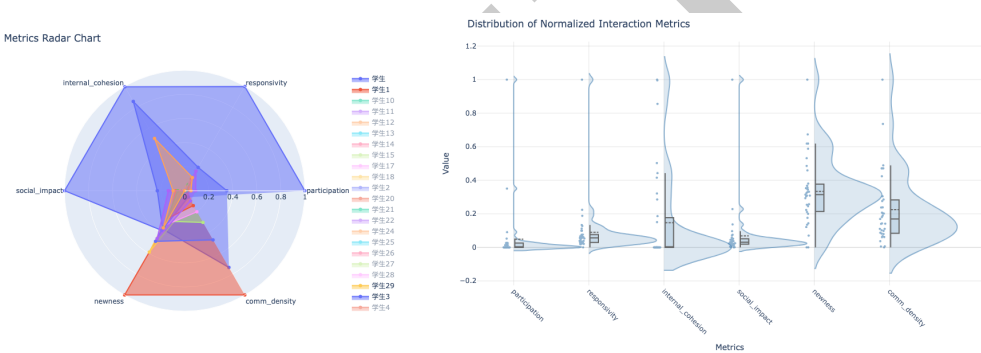


Figure 2: Visualizations for GCA Measures

- 65 ■ **Radar Plots:** Compare multiple measures across participants
- 66 ■ **Distribution Plots:** Visualize the distribution of measures

67 Results are saved as interactive HTML files in the specified output directory, allowing for easy  
68 exploration and sharing of analysis outcomes.

## 69 Mathematics

70 The GCA Analyzer implements several key mathematical formulas for analyzing group conver-  
71 sations:

### 72 Participation Rate

73 For a participant  $a$ , the participation count  $\|P_a\|$  and average participation rate  $\bar{p}_a$  are calculated  
74 as:

75 
$$\|P_a\| = \sum_{t=1}^n M_{a,t}$$

76 
$$\bar{p}_a = \frac{1}{n} \|P_a\|$$

77 where  $M_{a,t}$  is 1 if person  $a$  contributes at time  $t$ , and 0 otherwise, and  $n$  is the total number  
78 of contributions.

### 79 Participation Standard Deviation

80 The participation standard deviation  $\sigma_a$  for participant  $a$  is:

$$\sigma_a = \sqrt{\frac{1}{n-1} \sum_{t=1}^n (M_{a,t} - \bar{p}_a)^2}$$

## Normalized Participation Rate

The normalized participation rate ( $\hat{P}_a$ ) is computed relative to equal participation:

$$\hat{P}_a = \frac{\bar{p}_a - \frac{1}{k}}{\frac{1}{k}}$$

where  $k$  is the number of participants.

## Cross-Cohesion Matrix

The cross-cohesion matrix  $\Xi$  for analyzing temporal interactions is computed as:

$$\Xi_{ab} = \frac{1}{w} \sum_{\tau=1}^w \frac{\sum_{t \geq \tau} M_{a,t-\tau} M_{b,t} S_{t-\tau,t}}{\sum_{t \geq \tau} M_{a,t-\tau} M_{b,t}}$$

where:

- $w$  is the optimal window length
- $S_{t-\tau,t}$  is the cosine similarity between messages at times  $t - \tau$  and  $t$
- $M_{a,t}$  and  $M_{b,t}$  are participation indicators for persons  $a$  and  $b$  at time  $t$

## Internal Cohesion

For each participant  $a$ , internal cohesion is their self-interaction:

$$C_a = \Xi_{aa}$$

## Overall Responsivity

The overall responsivity  $R_a$  for participant  $a$  is:

$$R_a = \frac{1}{k-1} \sum_{b \neq a} \Xi_{ab}$$

## Social Impact

The social impact  $I_a$  for participant  $a$  is:

$$I_a = \frac{1}{k-1} \sum_{b \neq a} \Xi_{ba}$$

## Message Newness

Message newness measures how much new semantic content a message contributes compared to all previous messages in the conversation. For a message vector  $\vec{d}_i$ , we decompose it into two orthogonal components: new content and given (previously discussed) content.

The message vector  $\vec{d}_i$  is first projected onto two complementary subspaces:

1. The subspace  $G_i$  spanned by all previous message vectors  $\{\vec{d}_1, \vec{d}_2, \dots, \vec{d}_{i-1}\}$ , representing the semantic space of previous discussion
2. Its orthogonal complement  $G_i^\perp$ , representing potential new semantic content

The projection results in: - Given content:  $\vec{g}_i = \text{Proj}_{G_i}(\vec{d}_i)$  - New content:  $\vec{n}_i = \text{Proj}_{G_i^\perp}(\vec{d}_i)$

The newness score  $n(c_i)$  is then calculated as the ratio of new content magnitude to total content magnitude:

$$n(c_i) = \frac{\|\vec{n}_i\|}{\|\vec{n}_i\| + \|\vec{g}_i\|}$$

This score ranges from 0 (completely redundant with previous messages) to 1 (entirely new content). A higher score indicates the message contributes more novel semantic content to the discussion.

The overall newness  $N_a$  for participant  $a$  is the average newness across all their messages:

$$N_a = \frac{1}{|P_a|} \sum_{i \in P_a} n(c_i)$$

where  $P_a$  is the set of all contributions by participant  $a$ . This metric reflects how consistently a participant contributes new ideas to the conversation.

## Communication Density

For a message  $c_t$  at time  $t$ , its density  $D_i$  is:

$$D_i = \frac{\|c_t\|}{L_t}$$

where  $L_t$  is the word length of the message.

The average communication density  $\bar{D}_a$  for participant  $a$  is:

$$\bar{D}_a = \frac{1}{|P_a|} \sum_{t \in T_a} D_i$$

## References

- Dowell, N. M. M., Nixon, T. M., & Graesser, A. C. (2019). Group communication analysis: A computational linguistics approach for detecting sociocognitive roles in multiparty interactions. *Behavior Research Methods*, 51(3), 1007–1041. <https://doi.org/10.3758/s13428-018-1102-z>
- Wang, C., & Xiao, J. (2025). A role recognition model based on students' social-behavioural-cognitive-emotional features during collaborative learning. *Interactive Learning Environments*, 0(0), 1–20. <https://doi.org/10.1080/10494820.2024.2442706>