

# Package ‘ConversationAlign’

July 22, 2025

**Type** Package

**Title** Process Text and Compute Linguistic Alignment in Conversation Transcripts

**Version** 0.2.0

**Maintainer** Jamie Reilly <jamie\_reilly@temple.edu>

**Description** Imports conversation transcripts into R, concatenates them into a single dataframe appending event identifiers, cleans and formats the text, then yokes user-specified psycholinguistic database values to each word. 'ConversationAlign' then computes alignment indices between two interlocutors across each transcript for >40 possible semantic, lexical, and affective dimensions. In addition to alignment, 'ConversationAlign' also produces a table of analytics (e.g., token count, type-token-ratio) in a summary table describing your particular text corpus.

**License** GPL (>= 3)

**Encoding** UTF-8

**Depends** R (>= 3.5)

**Imports** DescTools, dplyr (>= 0.4.3), httr, magrittr, purrr, rlang, stringi, stringr, textstem, tibble, tidyr, tidyselect, stats, utils, YRmisc, zoo

**Suggests** devtools, knitr, rmarkdown, testthat (>= 3.0.0)

**URL** <https://github.com/Reilly-ConceptsCognitionLab/ConversationAlign>

**RoxygenNote** 7.3.2

**LazyData** true

**VignetteBuilder** knitr

**Collate** 'ConversationAlign-package.R' 'compute\_auc.R'  
'compute\_lagcorr.R' 'corpus\_analytics.R' 'data.R' 'globals.R'  
'prep\_dyads.R' 'read\_1file.R' 'read\_dyads.R'  
'replacements\_25.R' 'summarize\_dyads.R' 'utils.R' 'zzz.R'

**Config/testthat/edition** 3

**NeedsCompilation** no

**Author** Jamie Reilly [aut, cre] (ORCID:  
<<https://orcid.org/0000-0002-0891-438X>>),

Virginia Ulichney [aut],  
Ben Sacks [aut],  
Sarah Weinstein [ctb],  
Chelsea Helion [ctb],  
Gus Cooney [ctb]

**Repository** CRAN  
**Date/Publication** 2025-07-22 10:10:16 UTC

**Contents**

corpus_analytics . . . . .	2
load_github_data . . . . .	3
MaronGross_2013 . . . . .	3
NurseryRhymes . . . . .	4
NurseryRhymes_Prepmed . . . . .	4
prep_dyads . . . . .	5
read_1file . . . . .	6
read_dyads . . . . .	6
summarize_dyads . . . . .	7
<b>Index</b>	<b>8</b>

---

corpus_analytics	<i>corpus_analytics</i>
------------------	-------------------------

---

**Description**

Produces a table of corpus analytics including numbers of complete observations at each step, word counts, lexical diversity (e.g., TTR), stopword ratios, etc. Granularity of the summary statistics are guided by the user (e.g., by conversation, by conversation and speaker, collapsed all)

**Usage**

```
corpus_analytics(dat_prep)
```

**Arguments**

dat\_prep            takes dataframe produced from the df\_prep() function

**Value**

dataframe with summary statistics (mean, SD, range) for numerous corpus analytics (e.g., token count, type-token-ratio, word-count-per-turn) for the target conversation corpus. Summary data structured in table format for easy export to a journal method section.

---

load_github_data	<i>Load all .rda files from a GitHub data folder into the package environment</i>
------------------	-----------------------------------------------------------------------------------

---

### Description

Load all .rda files from a GitHub data folder into the package environment

### Usage

```
load_github_data(
  repo = "Reilly-ConceptsCognitionLab/ConversationAlign_Data",
  branch = "main",
  data_folder = "data",
  envir = parent.frame()
)
```

### Arguments

repo	GitHub repository (e.g., "username/repo")
branch	Branch name (default: "main")
data_folder	Remote folder containing .rda files (default: "data/")
envir	Environment to load into (default: package namespace)

### Value

nothing, loads data (as rda files) from github repository needed for other package functions

---

MaronGross_2013	<i>Sample Dyadic Interview Transcript: Marc Maron and Terry Gross Radio Interview 2013</i>
-----------------	--------------------------------------------------------------------------------------------

---

### Description

Text and talker information delineated, raw transcript, multiple lines per talker

### Usage

```
MaronGross_2013
```

### Format

## "MaronGross\_2013" A data.frame with 546 obs, 2 vars:

**text** text from interview

**speaker** speaker identity ...

---

NurseryRhymes	<i>Text and talker information delineated, 3 separate nursery rhymes, good for computing analytics and word counts</i>
---------------	------------------------------------------------------------------------------------------------------------------------

---

### Description

Text and talker information delineated, 3 separate nursery rhymes, good for computing analytics and word counts

### Usage

NurseryRhymes

### Format

## "NurseryRhymes" A data.frame with 100 observations, 2 vars:

**Event\_ID** factor 3 different simulated conversations

**Participant\_ID** fictional speaker names, 2 each conversation

**Text\_Raw** simulated language production, actually looped phrases from nursery rhymes ...

---

NurseryRhymes_Prepped	<i>Text and talker information delineated, 3 separate nursery rhymes, good for computing analytics and word counts</i>
-----------------------	------------------------------------------------------------------------------------------------------------------------

---

### Description

Text and talker information delineated, 3 separate nursery rhymes, good for computing analytics and word counts

### Usage

NurseryRhymes\_Prepped

### Format

## "NurseryRhymes\_Prepped" A data.frame with 1507 x 7 observations, 5 vars:

**Event\_ID** factor 3 different simulated conversations

**Participant\_ID** fictional speaker names, 2 each conversation

**Exchange\_Count** sequential numbering of exchanges by conversation, 1 exchange = 2 turns

**Turn\_Count** sequential numbering of turns by conversation

**Text\_Clean** content words

**emo\_anger** raw value of anger salience yoked to each word ...

---

prep\_dyads

---

*prep\_dyads*


---

## Description

Cleans, vectorizes and appends lexical norms to all content words in a language corpus. User guides options for stopword removal and lemmatization. User selects up to three psycholinguistic dimensions to yoke norms on each content word in the original conversation transcript.

## Usage

```
prep_dyads(
  dat_read,
  lemmatize = TRUE,
  omit_stops = TRUE,
  which_stoplist = "Temple_stops25",
  verbose = TRUE
)
```

## Arguments

<code>dat_read</code>	dataframe produced from <code>read_dyads()</code> function
<code>lemmatize</code>	logical, should words be lemmatized (switched to base morphological form), default is TRUE
<code>omit_stops</code>	option to remove stopwords, default TRUE
<code>which_stoplist</code>	user-specified stopword removal method with options including "none", "SMART", "MIT_stops", "CA_OriginalStops", or "Temple_Stopwords25". "Temple_Stopwords25" is the default list
<code>verbose</code>	display detailed output such as error messages and progress (default is TRUE)

## Value

dataframe with text cleaned and vectorized to a one word per-row format. Lexical norms and metadata are appended to each content word. Cleaned text appears under a new column called 'Text\_Clean'. Any selected dimensions (e.g., word length) and metadata are also appended to each word along with speaker identity, turn, and Event\_ID (conversation identifier).

---

read_1file	<i>read_1file</i>
------------	-------------------

---

### Description

Reads pre-formatted dyadic (2 interlocutor) conversation transcript already imported into your R environment.

### Usage

```
read_1file(my_dat)
```

### Arguments

my_dat	one conversation transcript already in the R environment
--------	----------------------------------------------------------

### Value

a dataframe formatted with 'Event\_ID', "Participant\_ID", "Text\_Raw" fields – ready for clean\_dyads()

---

read_dyads	<i>read_dyads</i>
------------	-------------------

---

### Description

Reads pre-formatted dyadic (2 interlocutor) conversation transcripts from your machine. Transcripts must be either csv or txt format. IF you are supplying a txt file, your transcript must be formatted as an otter.ai txt file export. Your options for using csv files are more flexible. ConversationAlign minimally requires a csv file with two columns, denoting interlocutor and text. Each separate conversation transcript should be saved as a separate file. ConversationAlign will use the file names as a document ID. Within the read dyads function, set the my\_path argument as the directory path to the local folder containing your transcripts on your machine (e.g., "my\_transcripts"). Please see our github page for examples of properly formatted transcripts: <https://github.com/Reilly-ConceptsCognitionLab/ConversationAlign>

### Usage

```
read_dyads(my_path = "my_transcripts")
```

### Arguments

my_path	folder of conversation transcripts in csv or txt format
---------	---------------------------------------------------------

### Value

a dataframe where each individual conversation transcript in a user's directory has been concatenated. read\_dyads appends a unique document identifier to each conversation transcript appending its unique filename as a factor level to 'Event\_ID'.

---

summarize_dyads	<i>summarize_dyads</i>
-----------------	------------------------

---

### Description

Calculates and appends 3 measures for quantifying alignment. Appends the averaged value for each selected dimension by turn and speaker. Calculates and Spearman's rank correlation between interlocutor time series and appends by transcript. Calculates the area under the curve of the absolute difference time series between interlocutor time series. The length of the difference time series can be standardized the shortest number of exchanges present in the group using an internally defined resampling function, called with `resample = TRUE`. Spearman's rank correlation and area under the curve become less reliable for dyads under 30 exchanges.

### Usage

```
summarize_dyads(
  df_prep,
  custom_lags = NULL,
  sumdat_only = TRUE,
  corr_type = "Pearson"
)
```

### Arguments

<code>df_prep</code>	produced in the <code>align_dyads</code> function
<code>custom_lags</code>	integer vector, should any lags be added in addition to -2, 0, 2
<code>sumdat_only</code>	default=TRUE, group and summarize data, two rows per conversation, one row for each participant, false will fill down summary statistics across all exchanges
<code>corr_type</code>	option for computing lagged correlations turn-by-turn covariance (default='Pearson')

### Value

either: - a grouped dataframe with summary data aggregated by conversation (`Event_ID`) and participant if `sumdat_only=T`. - the original dataframe 'filled down' with summary data (e.g., AUC, turn-by-turn correlations) for each conversation is `sumdat_only=F`.

# Index

## \* **datasets**

MaronGross\_2013, [3](#)

NurseryRhymes, [4](#)

NurseryRhymes\_Prepped, [4](#)

corpus\_analytics, [2](#)

load\_github\_data, [3](#)

MaronGross\_2013, [3](#)

NurseryRhymes, [4](#)

NurseryRhymes\_Prepped, [4](#)

prep\_dyads, [5](#)

read\_1file, [6](#)

read\_dyads, [6](#)

summarize\_dyads, [7](#)