

# Multimodal Data Analysis

Lecture for ISEA Training Program

Wei Ai <[aiwei@umd.edu](mailto:aiwei@umd.edu)>

University of Maryland

# Let's take a step back...

- What do data mining looks like

# Data Mining: Knowledge Discovery from Data

Data to  
be mined

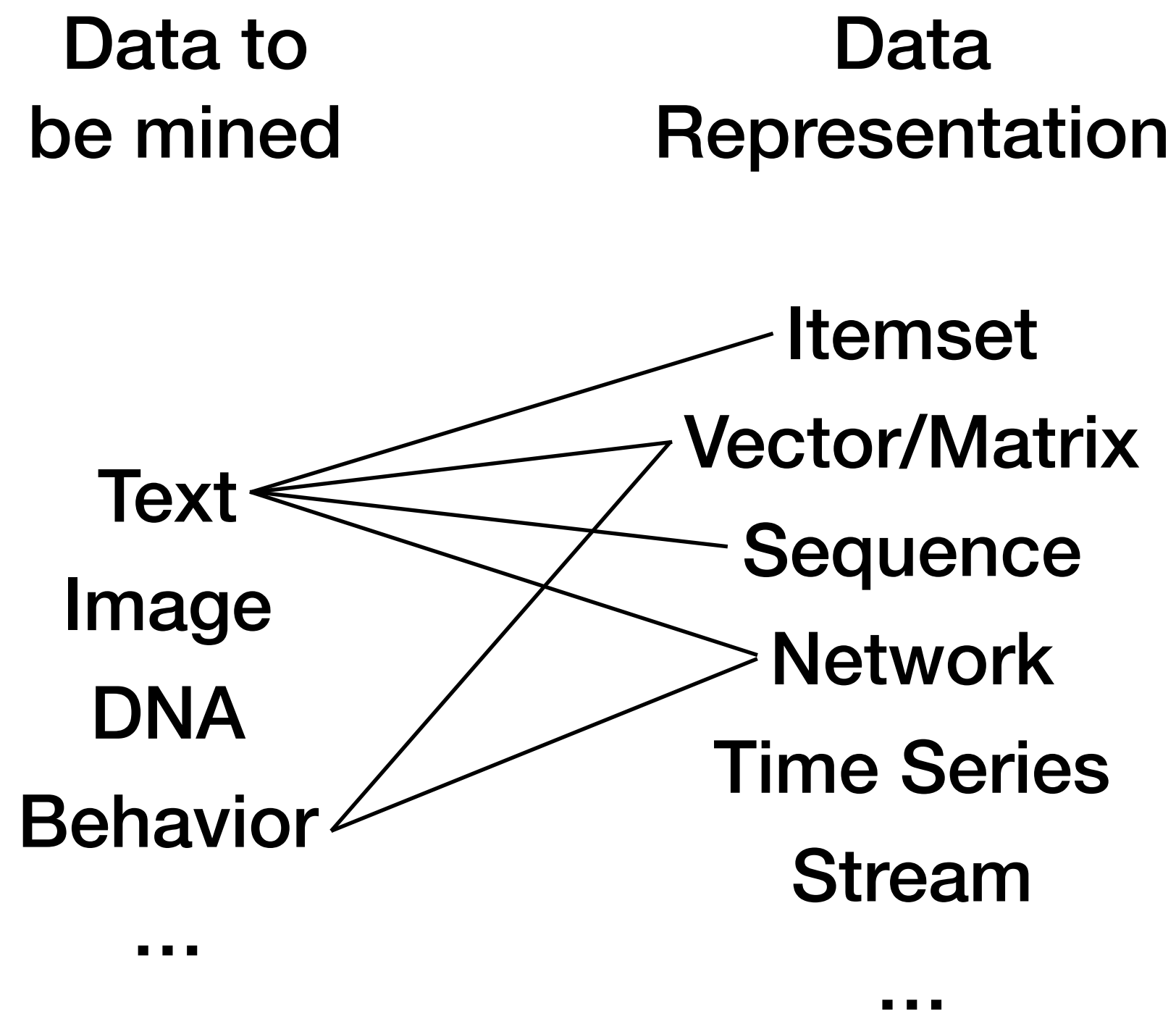
Data  
Representation

Basic  
Functionalities

Knowledges  
(Outcome)

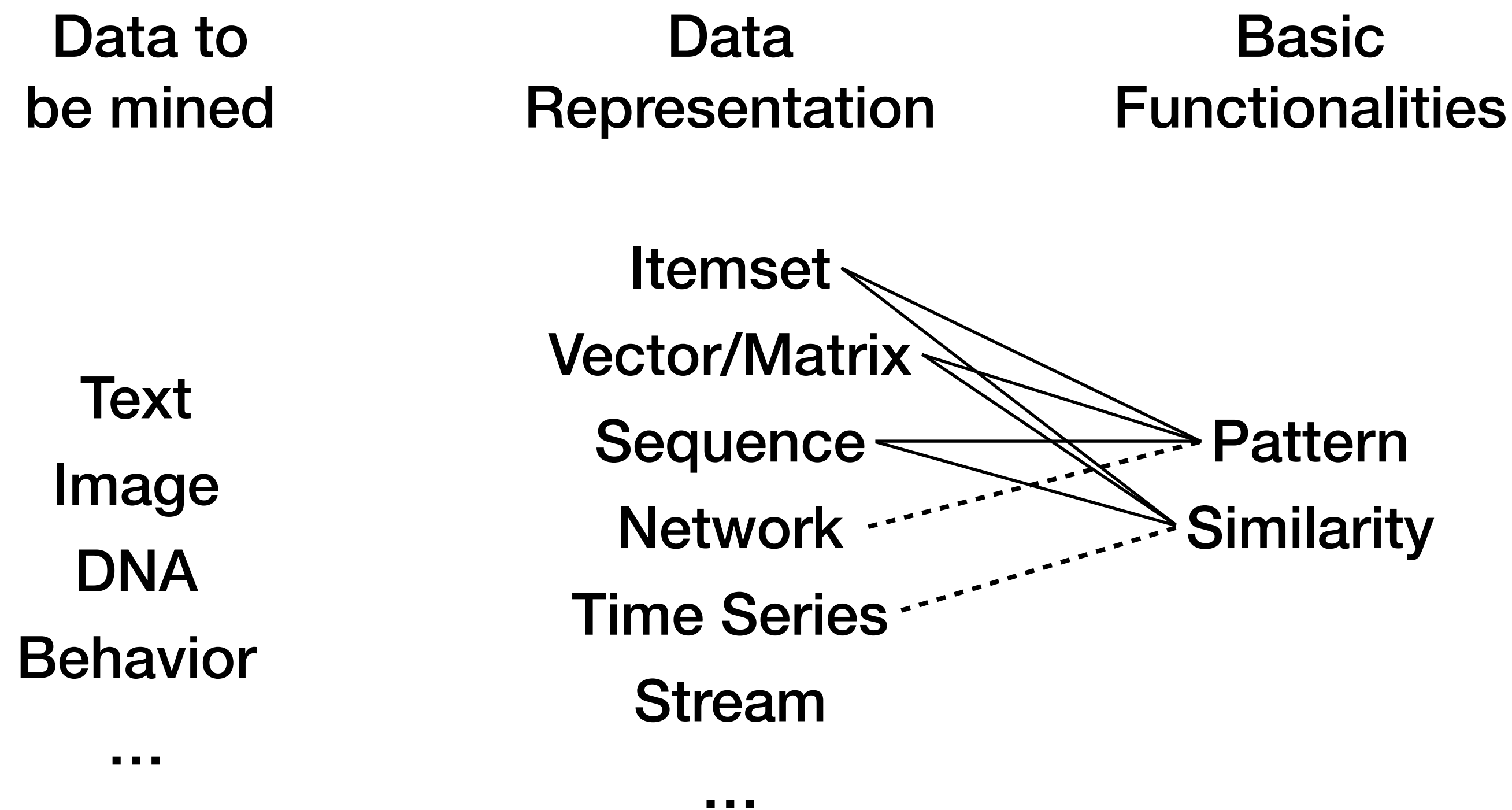
Applications

# Data Mining: Knowledge Discovery from Data



- Text  $\rightarrow$  Itemset:
  - {"to", "be", "or", "not", ...}
- Text  $\rightarrow$  Vector:
  - $\langle 1, 0, 0, 1 \rangle$
- Behavior  $\rightarrow$  Vector:
  - User-Product Rating Matrix
- Behavior  $\rightarrow$  Network:
  - Twitter "Following" network
- Questions to keep in mind:
  - What is the granularity of the analysis
  - What counts as an observation
  - What does each "row" represent in your DataFrame

# Data Mining: Knowledge Discovery from Data



- Itemset + Patterns:
  - Frequent Pattern Mining
- Itemset + Similarity:
  - Jaccard Similarity
- Vector + Similarity:
  - Dot Product
  - Manhattan/Euclidean distance
  - Cosine Similarity
- Sequence + Similarity:
  - Edit distance
  - Shingling
- Network + Pattern:
  - ...

# Data Mining: Knowledge Discovery from Data

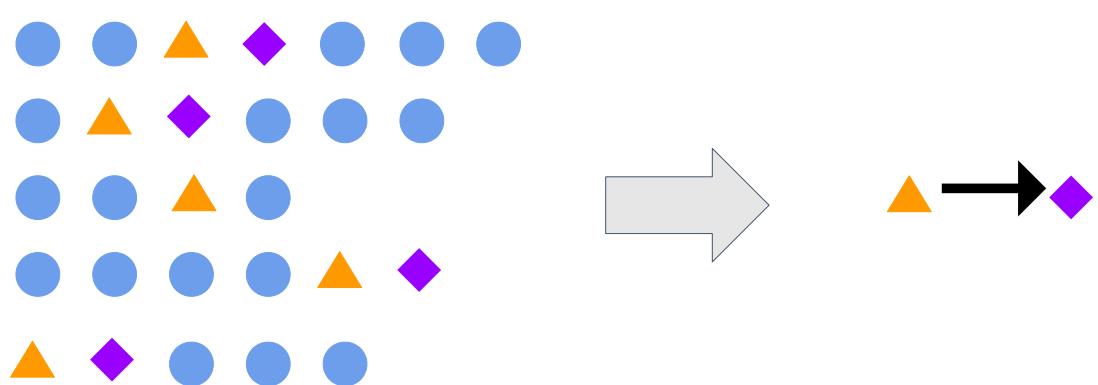
Data to  
be mined

Data  
Representation

Basic  
Functionalities

Knowledges  
(Outcome)

Association



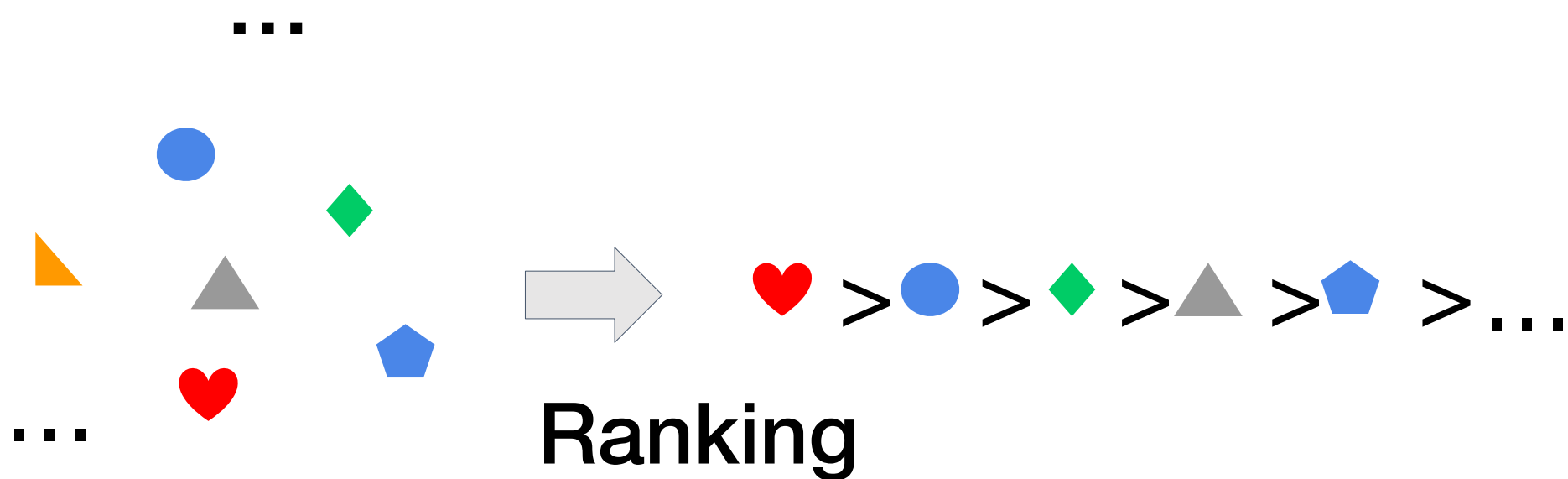
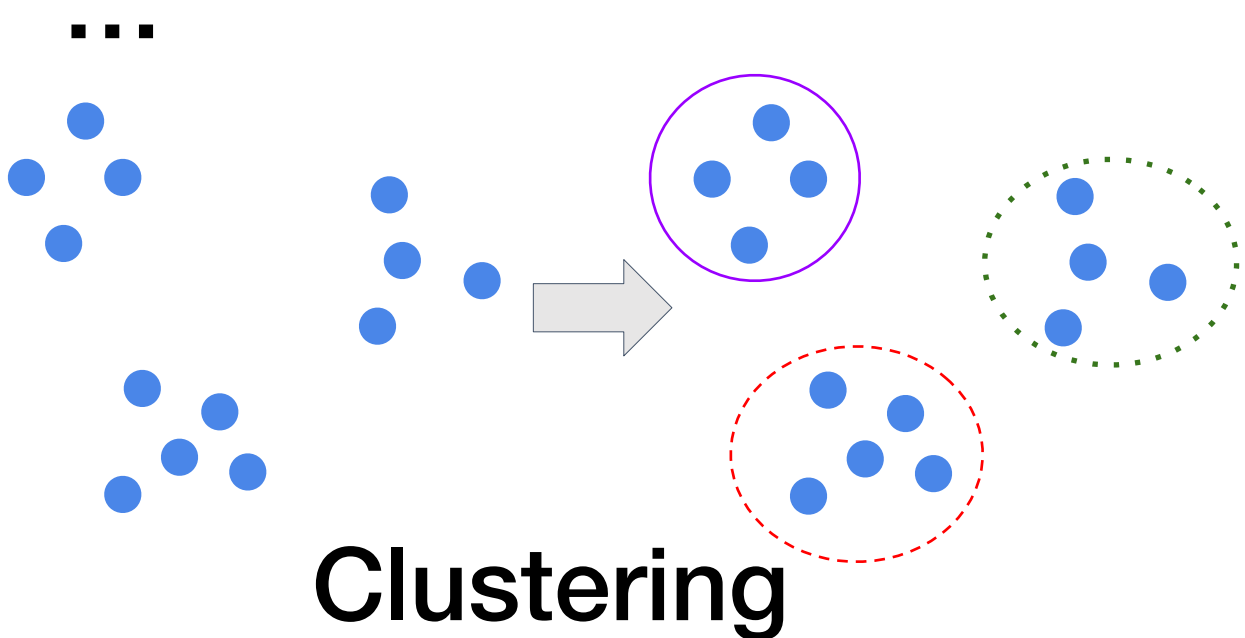
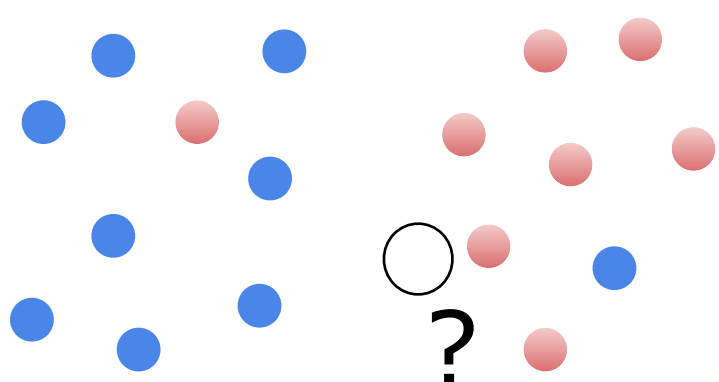
Text  
Image  
DNA  
Behavior  
...

Itemset  
Vector/Matrix  
Sequence  
Network  
Time Series  
Stream  
...

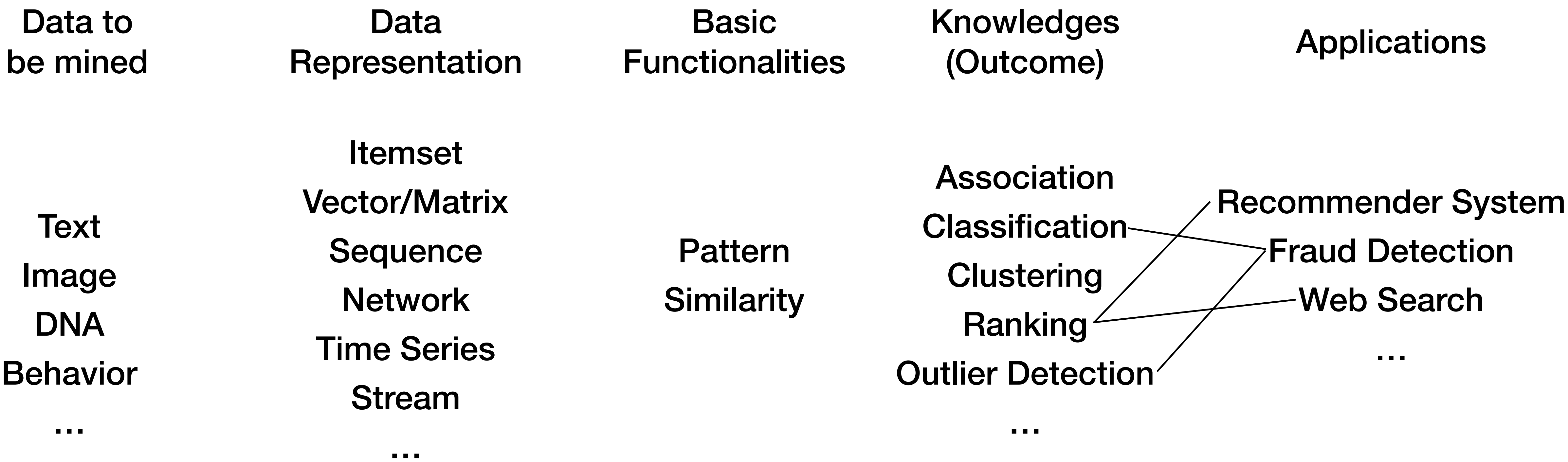
Pattern  
Similarity

Association  
Classification  
Clustering  
Ranking  
Outlier Detection  
...

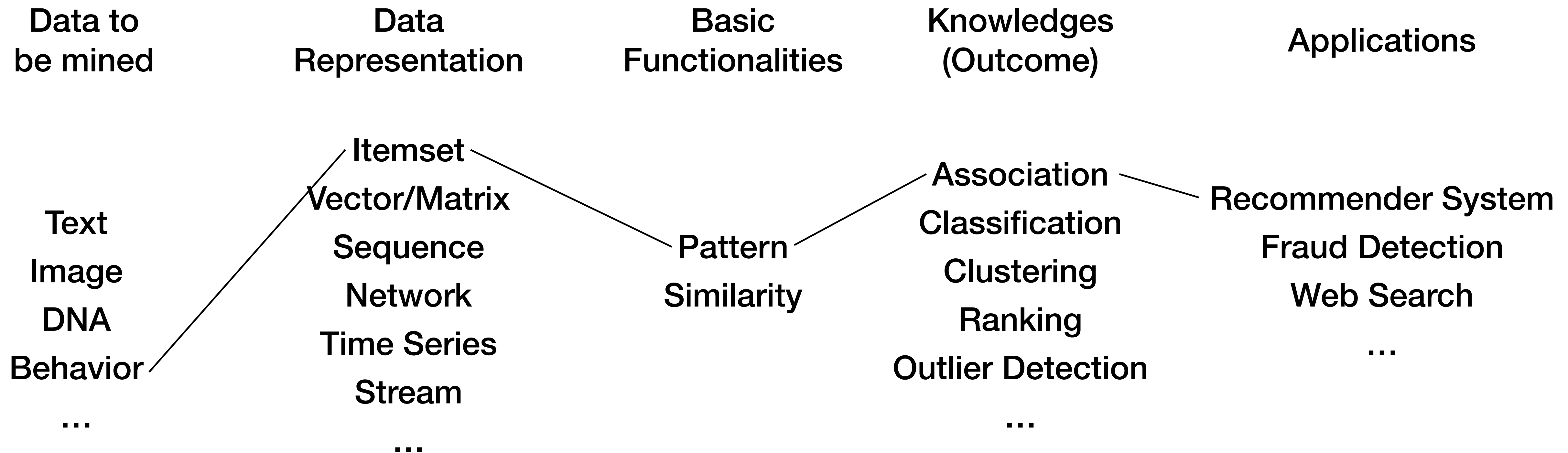
Classification



# Data Mining: Knowledge Discovery from Data

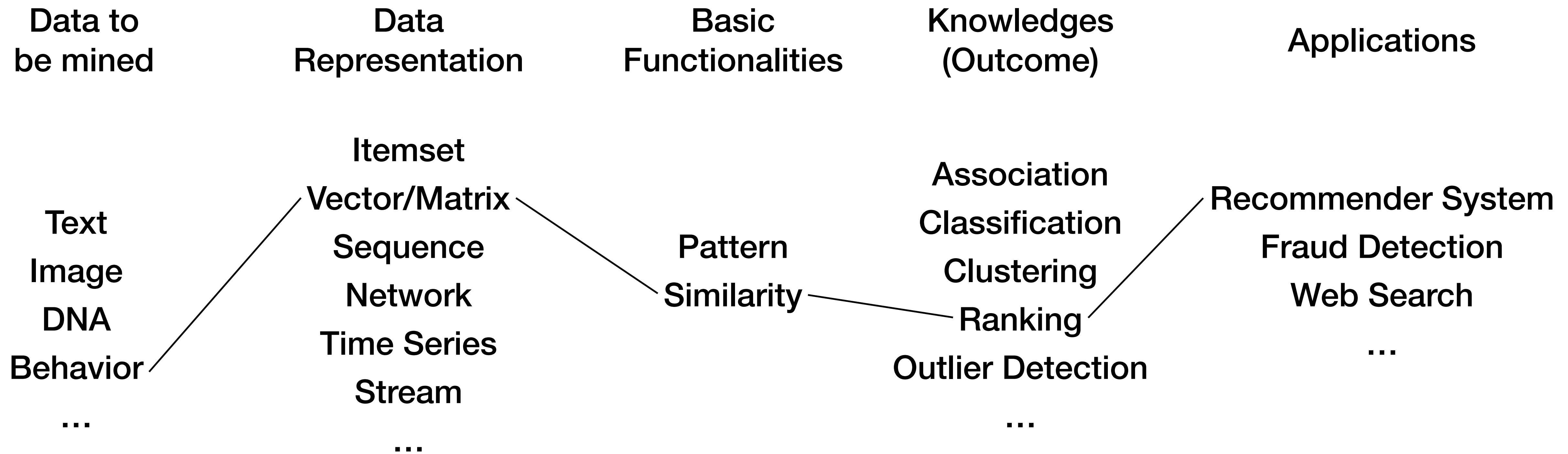


# Example: Amazon “Frequently Bought Together”

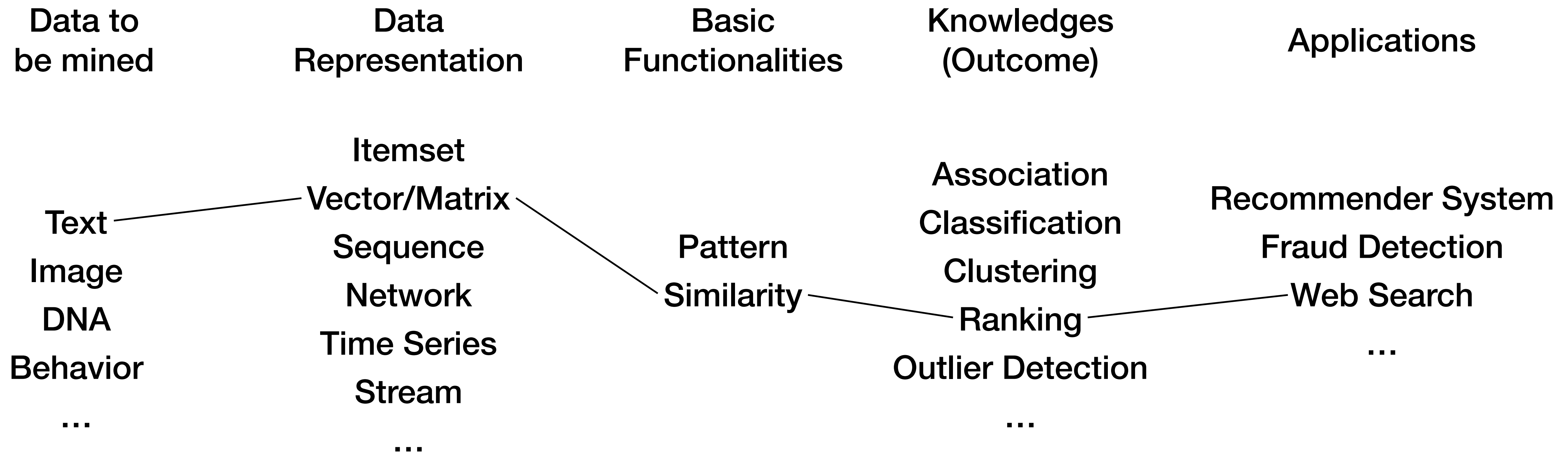




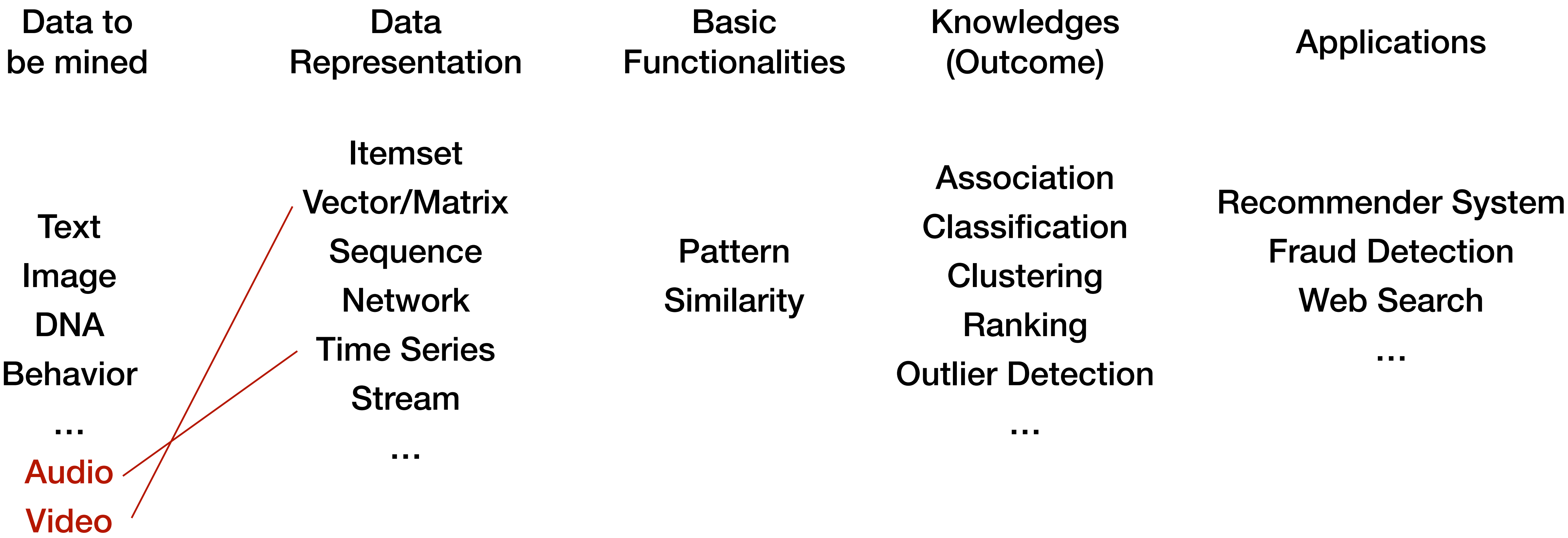
# Example: Netflix Movie Recommender System



# Example: Text Retrieval



# Multimodal Data are Just Another Data Genre



# Multimodal Data are Just Another Data Genre

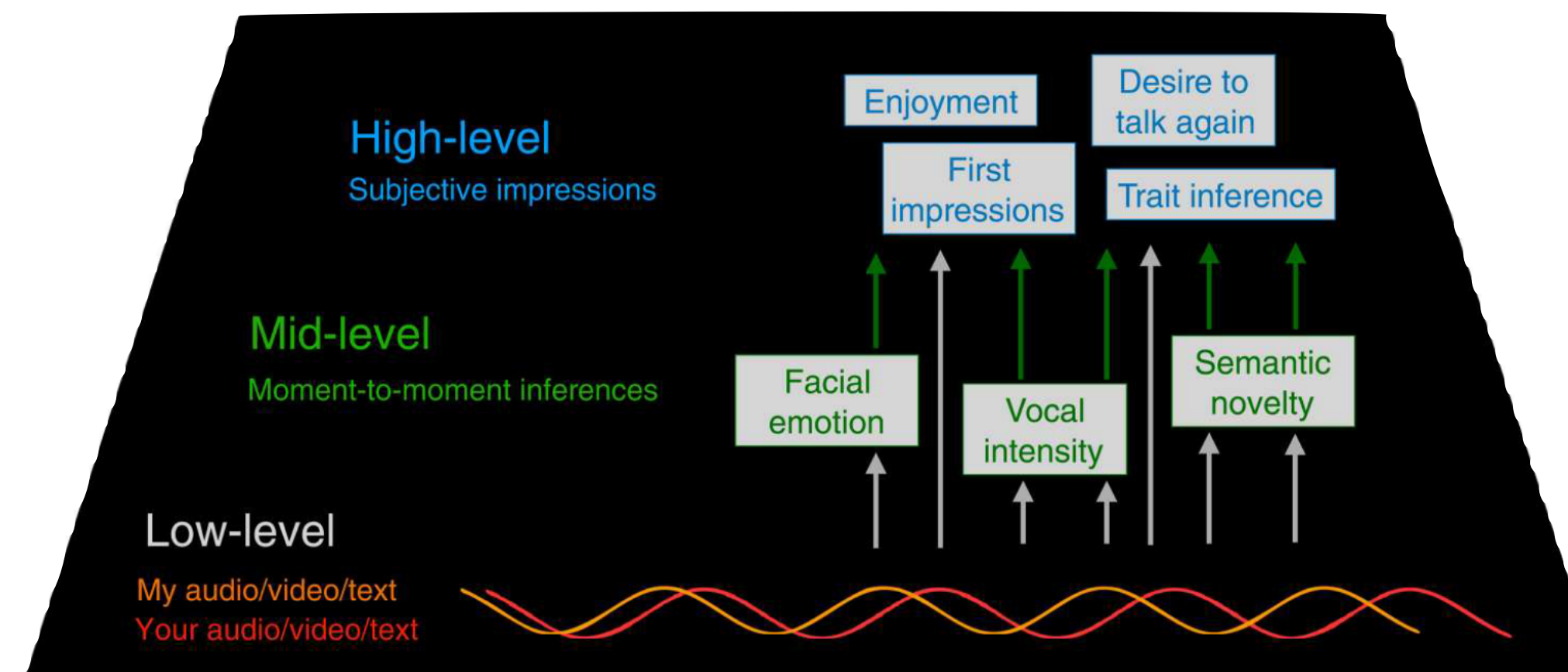
- We can still abstract multimodal data as the data representation that we are familiar with.
- It allows us to build features bottom-up, and can be naturally followed with
  - Exploratory analysis
  - Visualization
  - Correlation/prediction analysis
- For example:
  - Extract the *loudness* of one's speech at every 0.02 second, and make it a *time series*.
  - Classify one's *facial expression* with emotion labels, and *vectorize* its distribution.
  - Identify *nodding yes* and *shaking no* (head movement) as an event sequence.

# The CANDOR Dataset

- 1656 unscripted conversations over video chat, recorded in 2020
- Available data
  - Pre- and Post-conversation survey
    - Perceptions of their conversation partners, their feelings about the overall conversation, their personality, etc..
  - Video recording
    - Computationally extracted features
    - Modalities: Text, Audio, and Video

# Structured as an “vertically integrated” framework

- “low-level” mechanical features of conversation,
  - e.g. turn-taking.
- “mid-level” information streams,
  - e.g. semantic exchange, psycholinguistic markers, and emotion expressions,
- “high-level” judgments reported after conversation,
  - people’s enjoyment and the impressions they formed of their conversation partners.



# Low-level features

- Closest to raw signals in the audio, video, and text of a conversation recording
- Often vary on a nearly continuous time scale.
- Require some degree of inference:
  - Extracting vocal markers with signal processing (pitch, volume/loudness, etc.)
  - Automated transcription with Automated speech recognition (ASR)
  - Head movement, eye gaze, etc.
- Concrete, specific, and objective properties from which higher-order inferences are derived.



# High-level features

- Subjective judgments about their conversations.
- Formed on a coarse time scale (e.g. conversation level)
- Reflected in the survey responses:
  - measures of liking, enjoyment, and conversational flow,
  - evaluations of one's partner's social status, intelligence, and personality.
  - ...
- The broadest range of information and context are incorporated in making such judgements, and thus distinguished as high-level inferences.



# The rich mid-level features

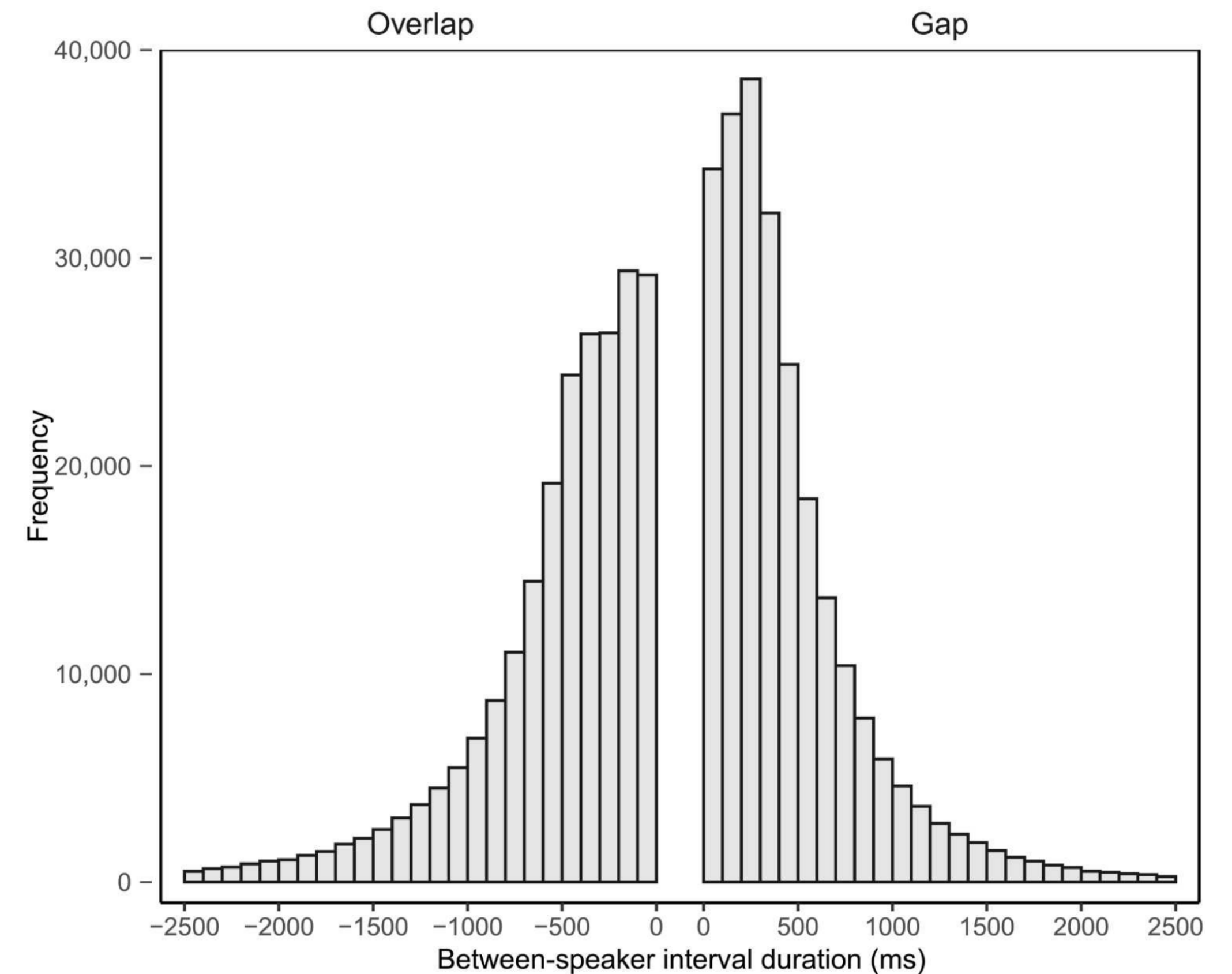
- numerous features related to subjective perceptions of the interaction that typically vary on an intermediate time scale (e.g. turn level)
- Usually computed using **algorithmic tools** that were trained to attend to specific aspects of speech, sound, and movement **to infer a psychological content**, for example:
  - happy facial expression,
  - increasing intensity in one's voice,
  - a timely change of subject.
- Characterized by their use of a narrower scope of context and antecedent reference.
  - *“A hitch in the voice, a sad glance away; all these signals, essential for shared understanding between humans, will go unnoticed by a machine that knows only language.”*

# The Turn-Taking System (low-level)

- Turn exchange: the way people manage to pass the floor back and forth in an orderly and efficient manner;
- Turn duration: how long speakers talk before they turn over the floor;
- Back-channeling: the active engagement that listeners display while speakers are talking
  - “mhm,” “yeah,” and “exactly”

# Turn exchange

- Median between-speaker turn interval was +80 ms and was distributed approximately normally.
- Consistent with previous literature.
- Takeaway: reproduce findings in existing literature as sanity check!



# Turn duration — What do we mean by “turn”?

Fatima:

Hello,

my

name

is

Fatima.

I'm

from

Egypt.

Eduardo:

Hi.

Nice

to

meet

you.

Audiophile	
Hello,	
	Hi.
my name is Fatima. I'm	
	Nice
from	
	to
Egypt.	
	meet you.

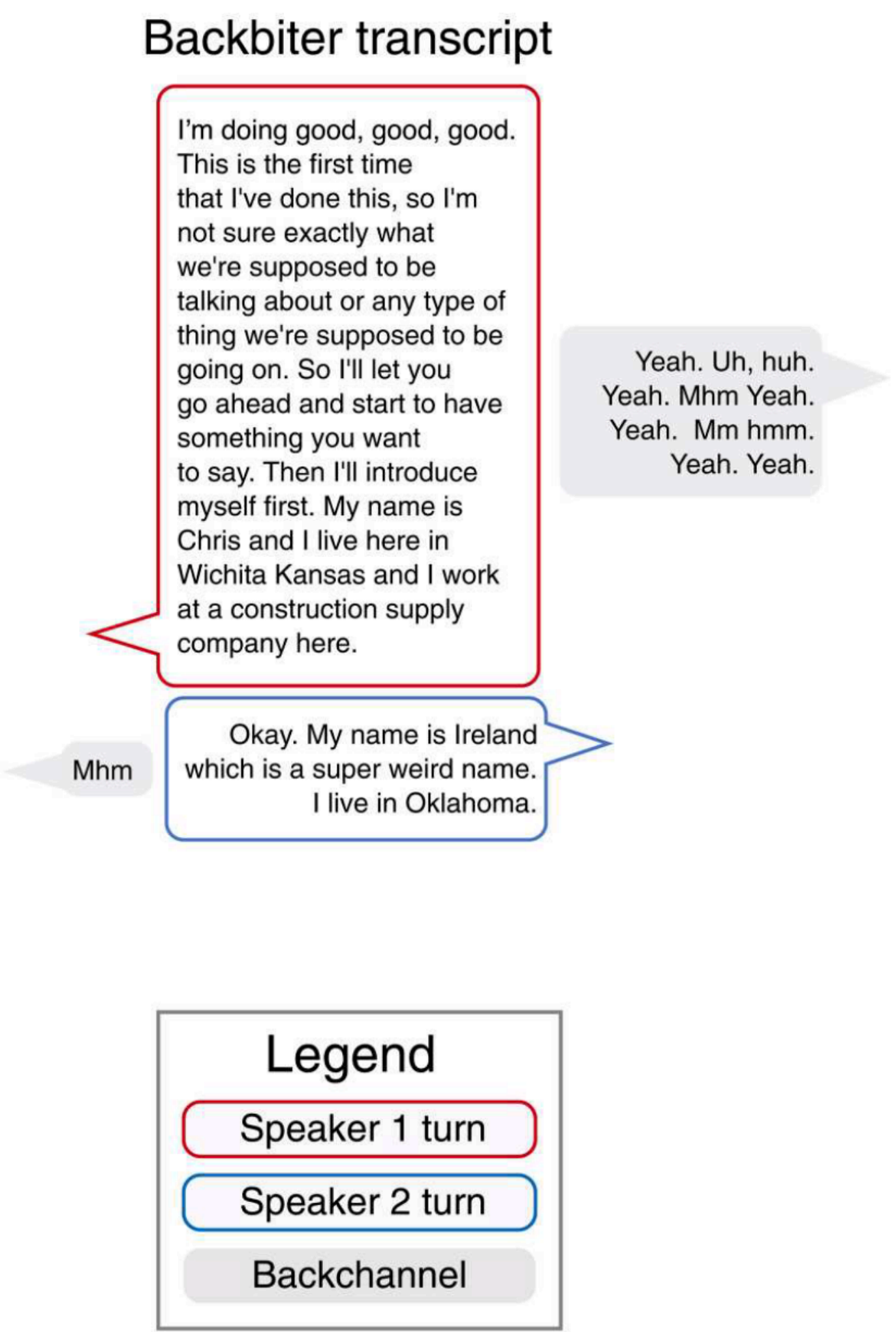
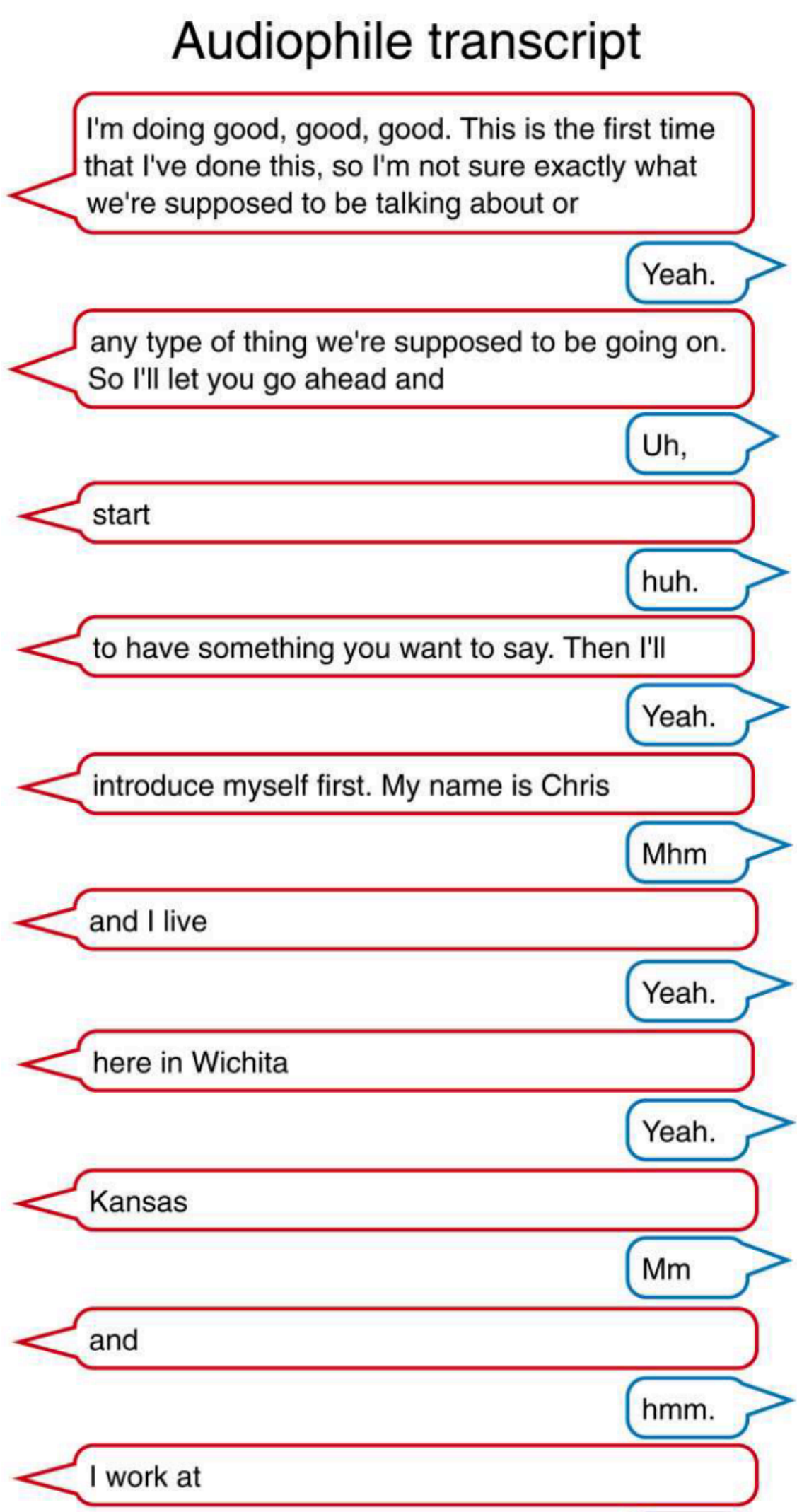
Cliffhanger	
Hello, my name is Fatima.	
	Hi.
I'm from Egypt.	
	Nice to meet you.

**Table 1. Basic comparison of Audiophile and Cliffhanger turn models.** Speakers’ mean and median turn durations were four to five times greater for the cliffhanger turn model compared to those of Audiophile. This suggests more broadly that analytic decisions about transcript segmentation will play a key role in the empirical investigation of speaking duration, an understudied topic.

Model	Turn duration		Number of words		Average turns
	Mean	Median	Mean	Median	
Audiophile	2.22	0.92	6.40	2	440.70
Cliffhanger	8.52	5.81	17.81	9	159.41

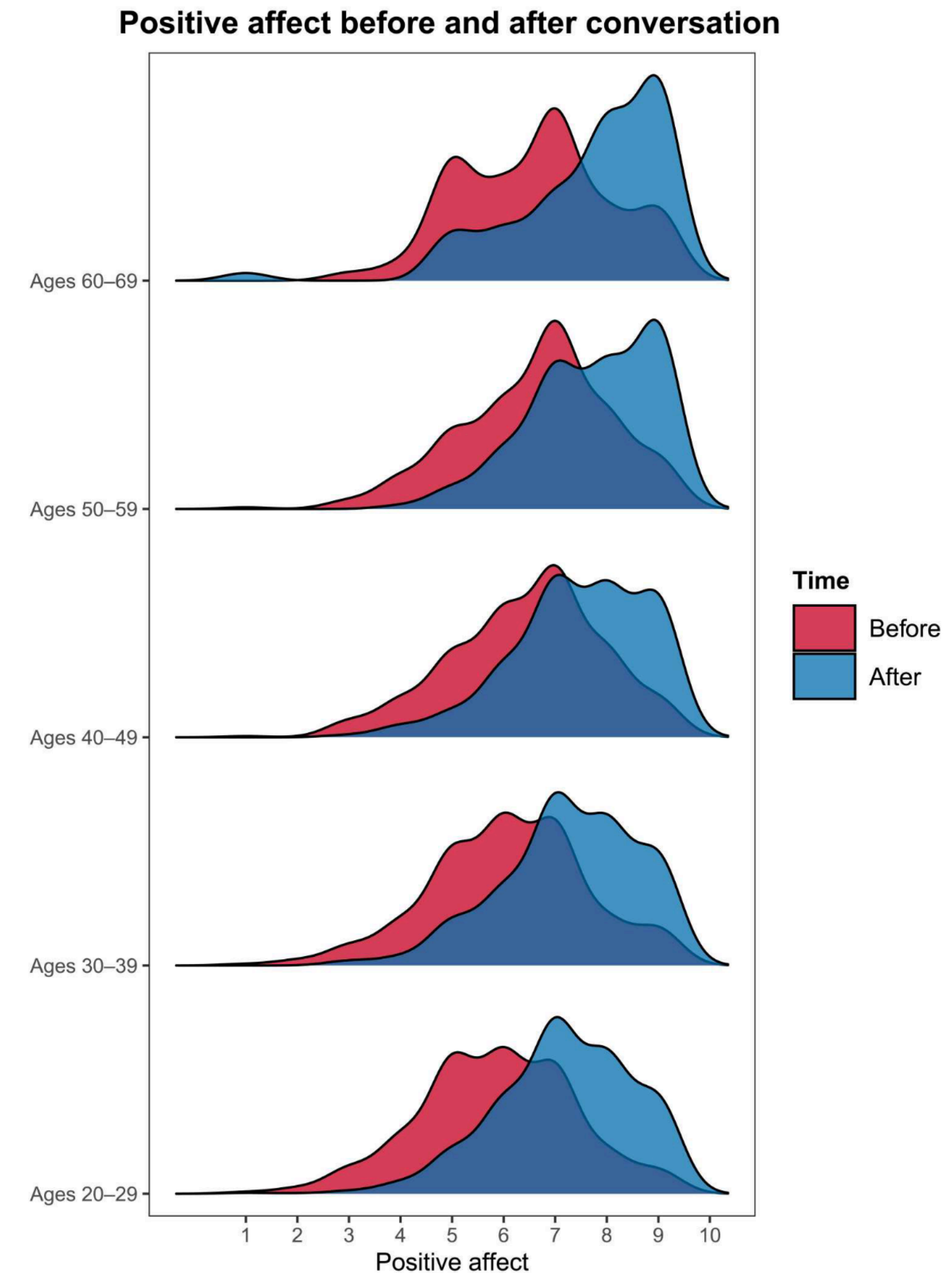


# Back-channelling



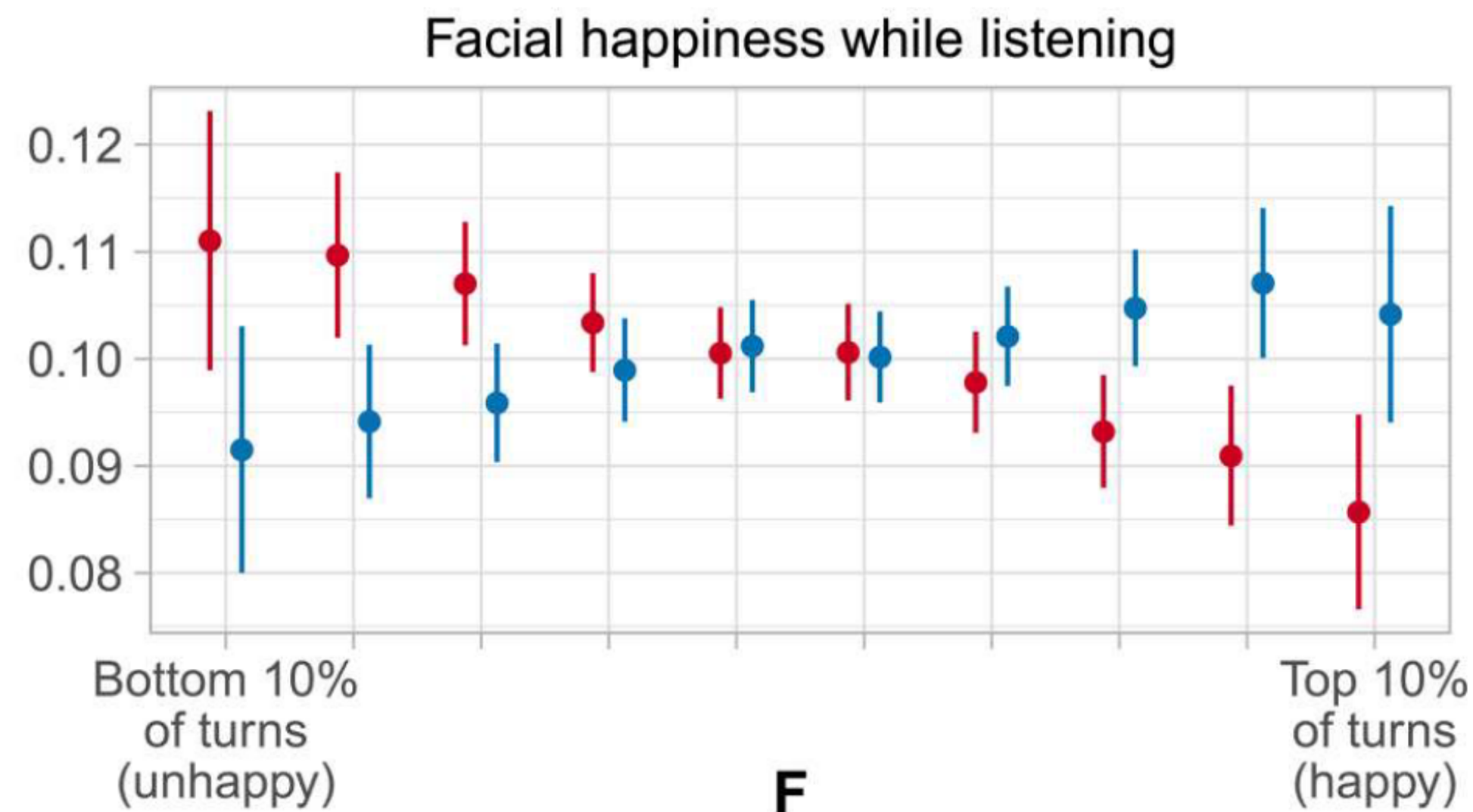
# Conversation and well-being (high-level)

- Respondents were asked to report their mood immediately before (red) and after (blue) their conversation.



# Now is the fun part!

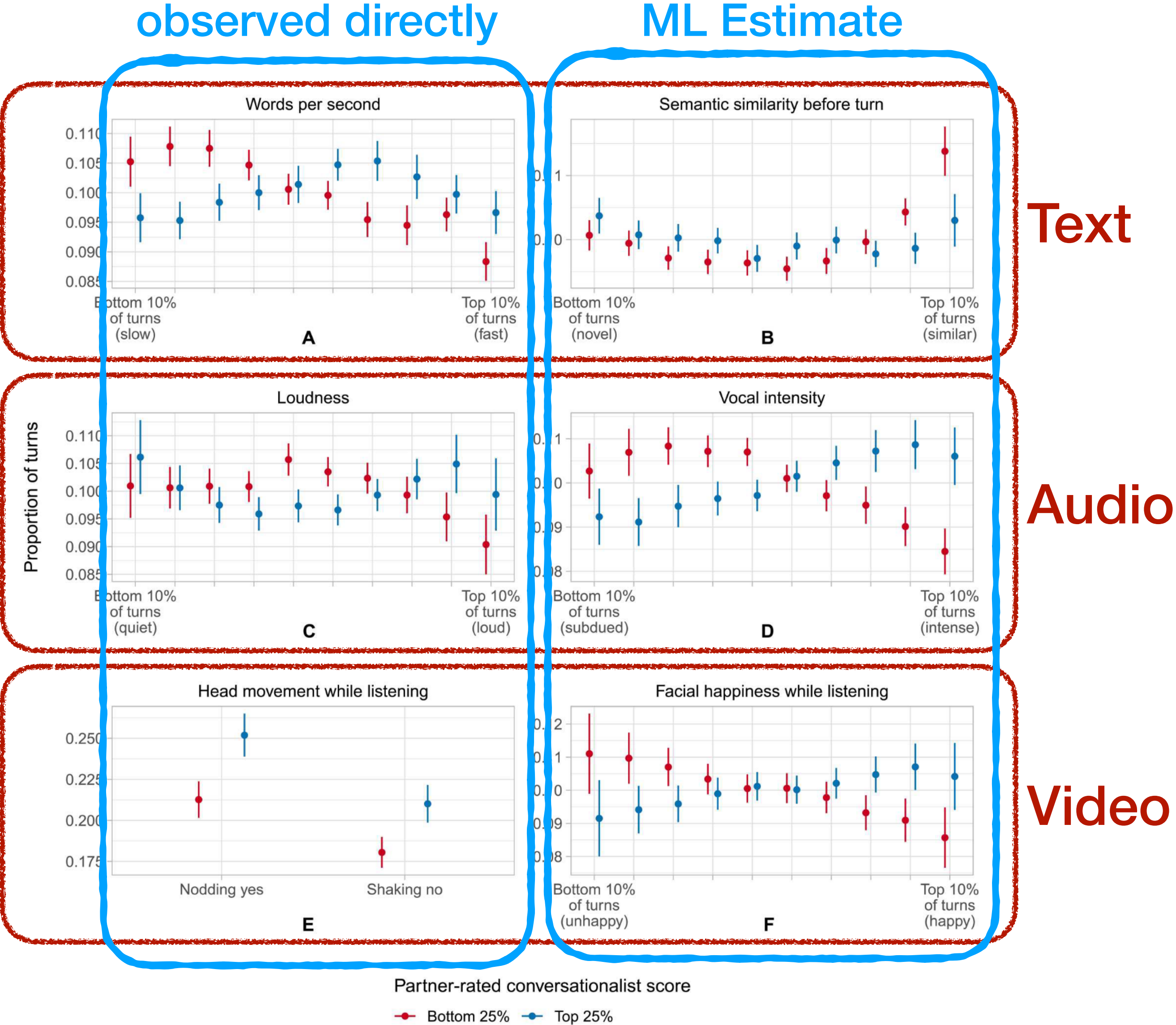
- We explore this *rich “middle layer”* of interaction and associate it with *high-level impressions* by exploring an open question in conversation research: *What distinguishes a good conversationalist?*
- E.g. good conversationalists exhibited significantly more facial happiness expressions while listening.





# Behavior Patterns of Good and Bad Conversationalists (fig. 8)

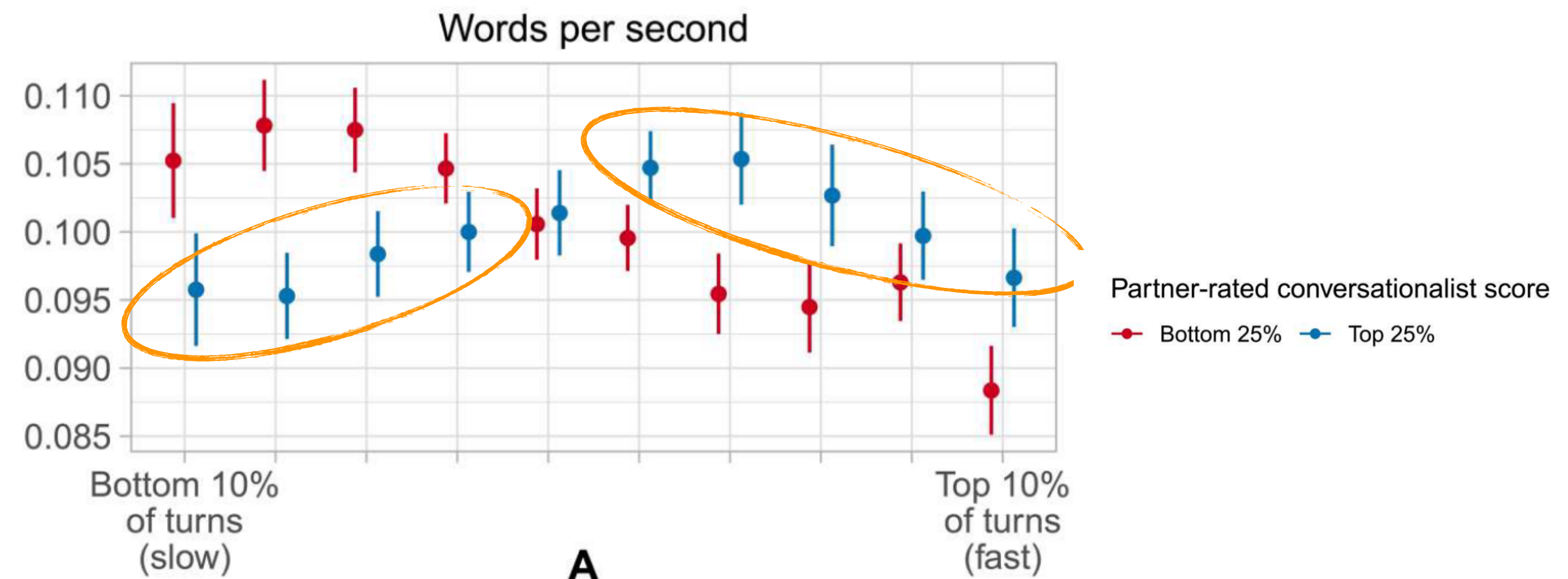
- Visualized with frequency plots, as opposed to linear analysis.





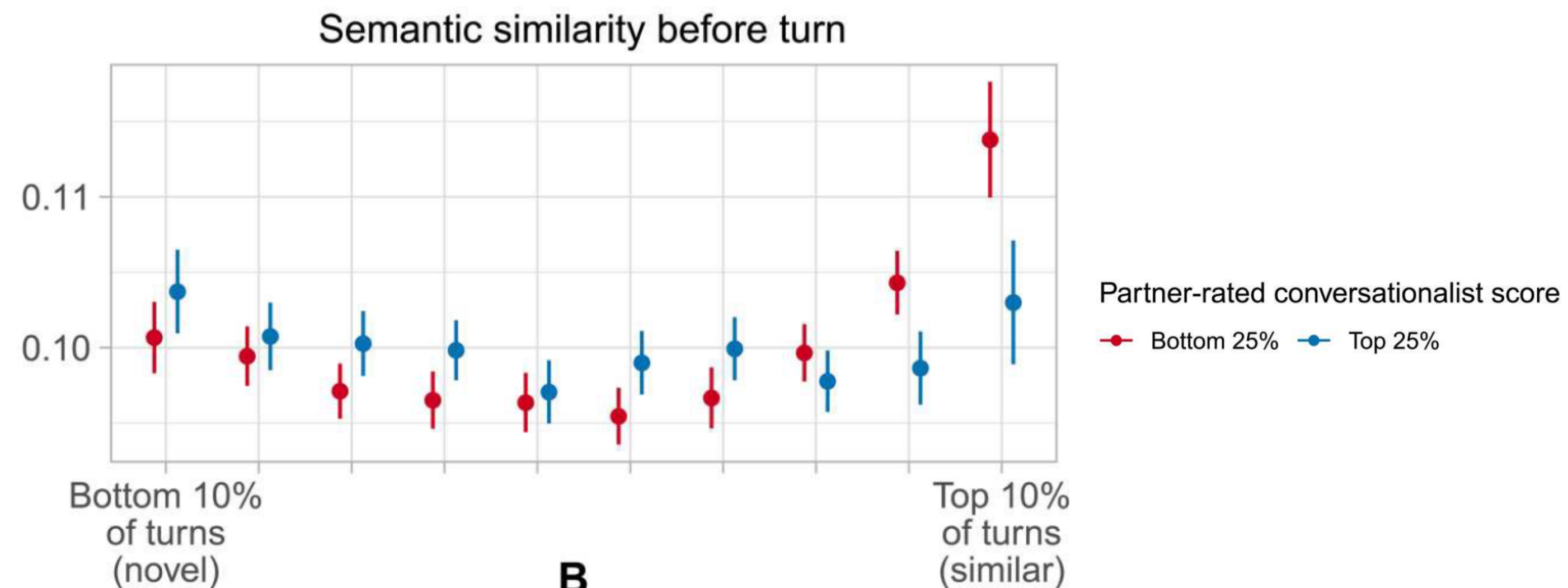
# Speech rate (word per minute)

- Good conversationalists spent more of their turns speaking quickly.
- Bad conversationalists spent a greater proportion of turns speaking slowly.



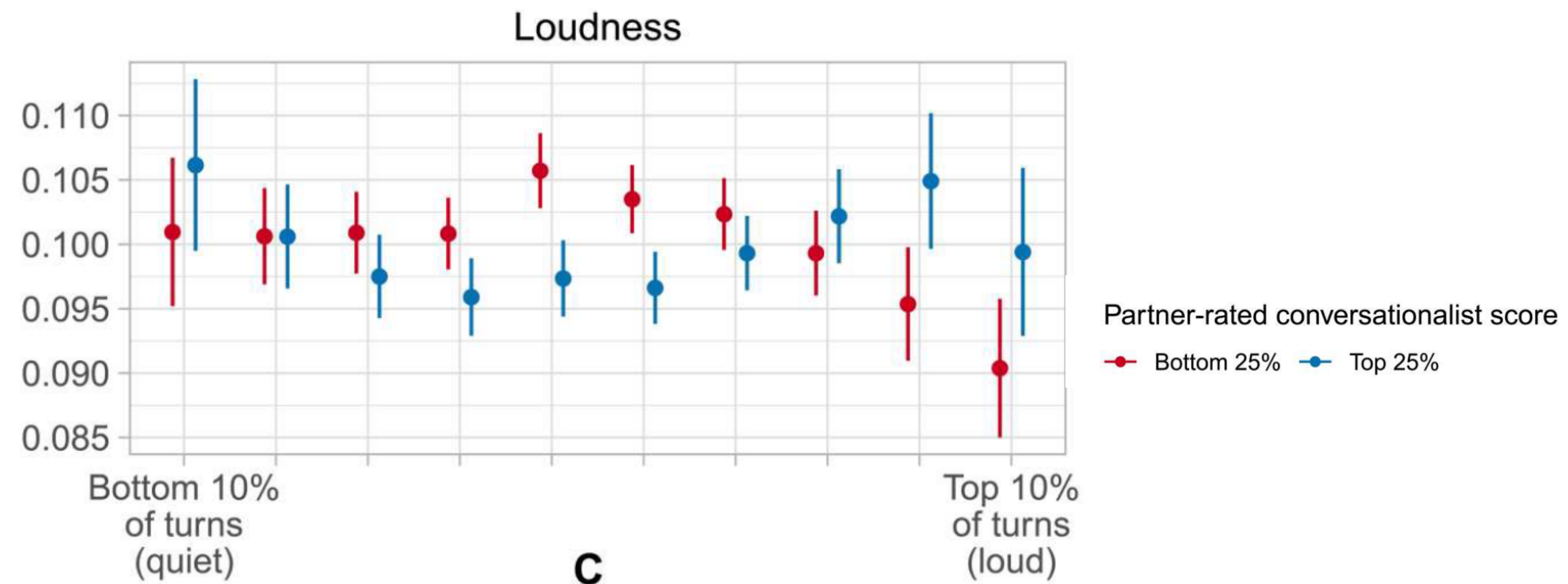
# Semantic similarity/novelty

- Use Text Embedding as vector representation of the spoken content.
- Use cosine similarity of the embedding vectors (this turn vs. previous turn) as a measurement of similarity/novelty.
- Good and bad conversationalists differ significantly.
- But good conversationalists also do not always add more novelty.



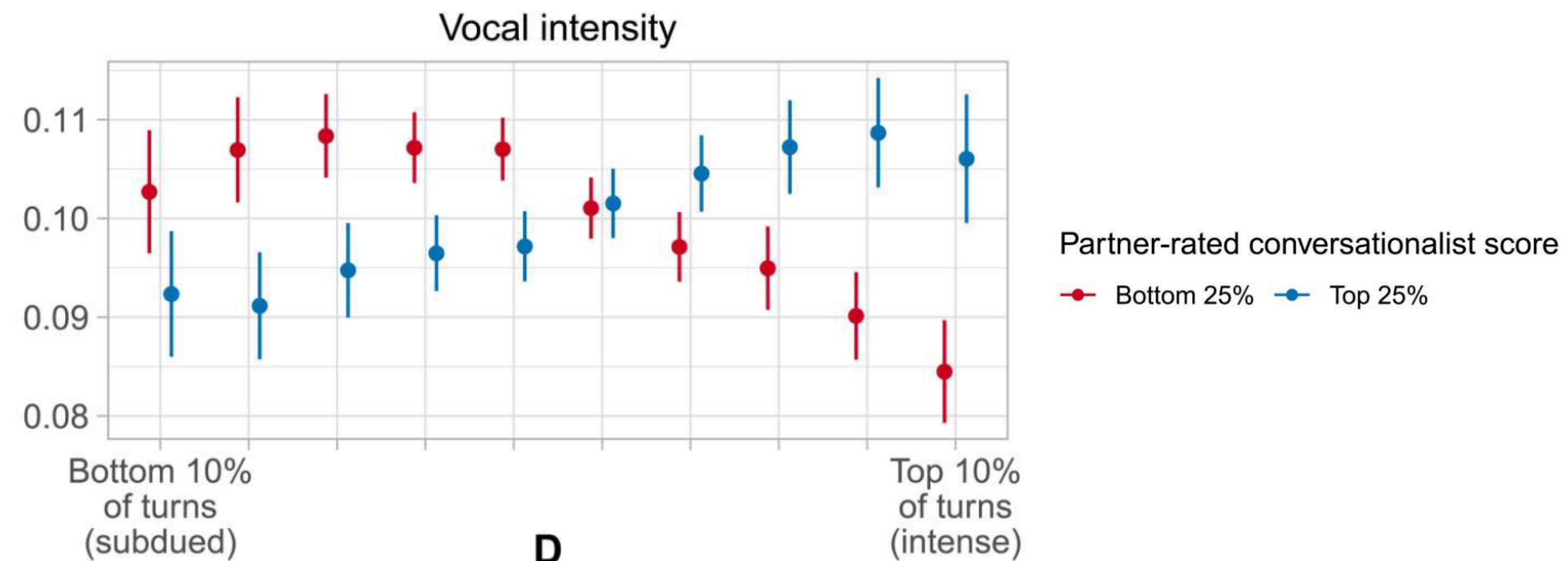
# Loudness

- No significant difference in the mean turn loudness.
- Bad conversationalists take more turns of **medium loudness**,
- Good conversationalists take more turns with **either lower or higher average loudness**



# Vocal intensity

- A vocal intensity classifier is trained with the RAVDESS dataset.
  - RAVDESS: recordings of trained actors who were prompted to read simple statements with either “normal” or “high” emotional intensity (used as labels).
  - Features include summary statistics of common prosodic features, such as
    - mean, max, and SD of fundamental frequency (F0), volume (log energy), etc.
  - The classifier is then applied to each turn in the CANDOR dataset
- people rated as good conversationalists spoke with greater intensity than bad conversationalists

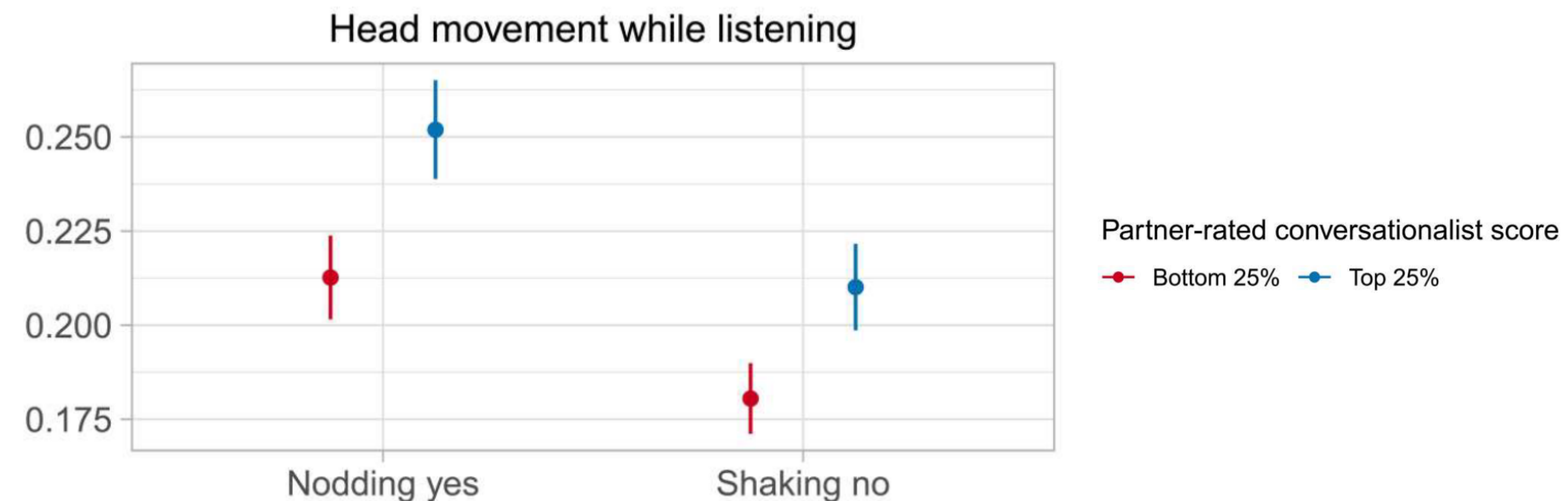




# Head Movement

Question: Is there better way to recognize nod and shake?

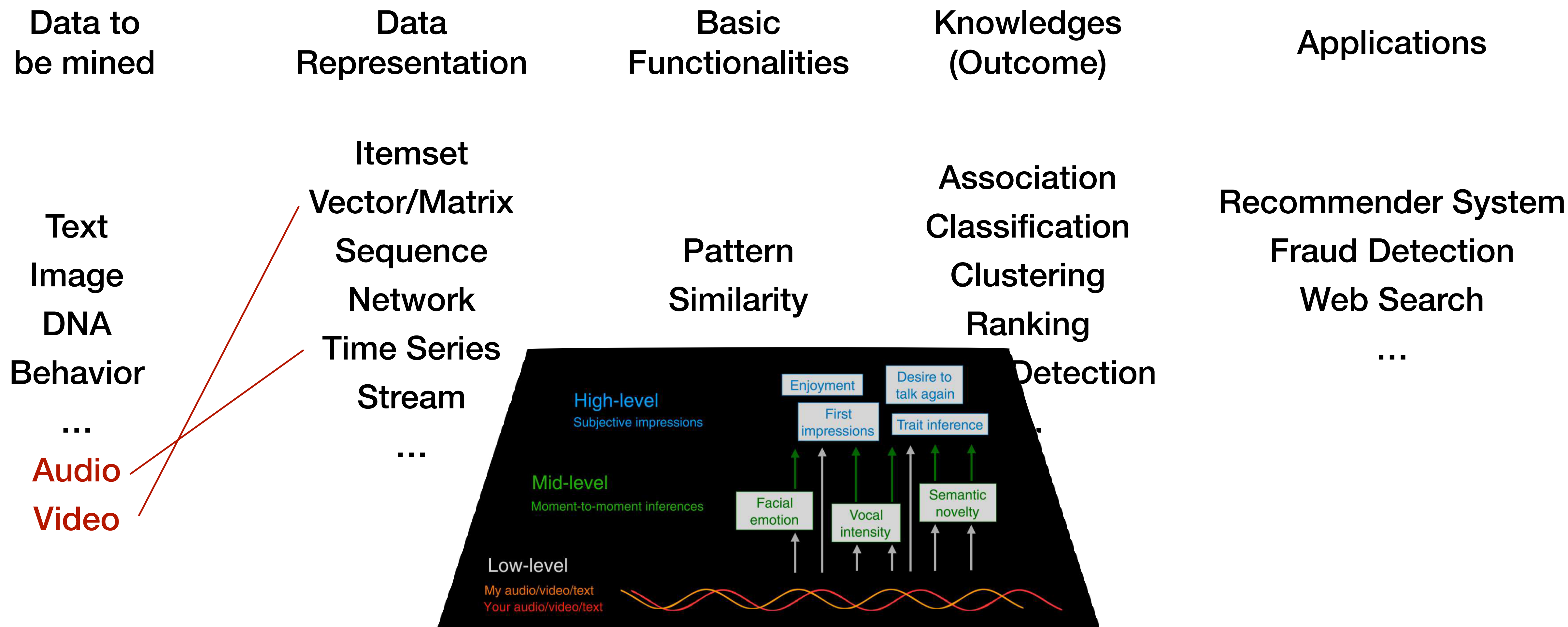
- Facial recognition followed with heuristic rules
  - Check if over a 2-s period, (i) at least 10% of a participant's face (ii) crossed its beginning position at least twice.
  - If yes, record “nod” if it occurs along vertical axis, record “shake” for horizontal.
- Good conversationalists were significantly more engaged not only in nodding “yes”, but also “shaking” no.



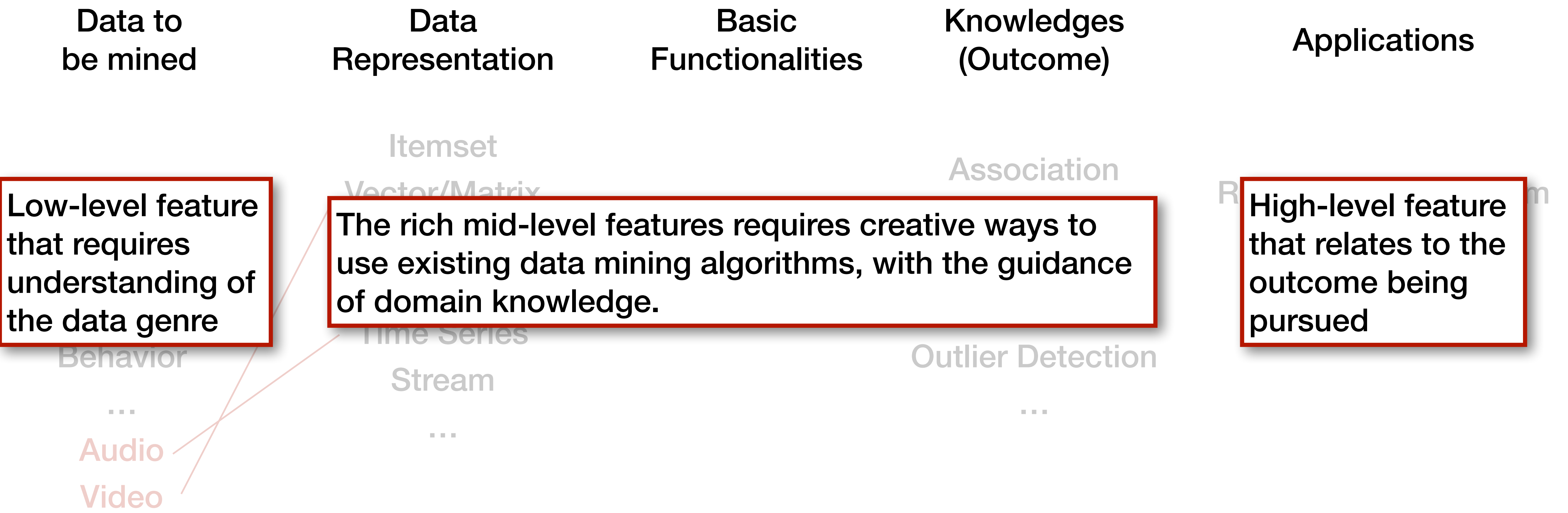
# The rest of the CANDOR paper discusses

- A qualitative glance of the topics in the conversation.
- Practical considerations and limitations of the dataset.
- A detailed Supplementary Materials

# Low-,Mid-,High-level features meets data mining pillars



# My opinion



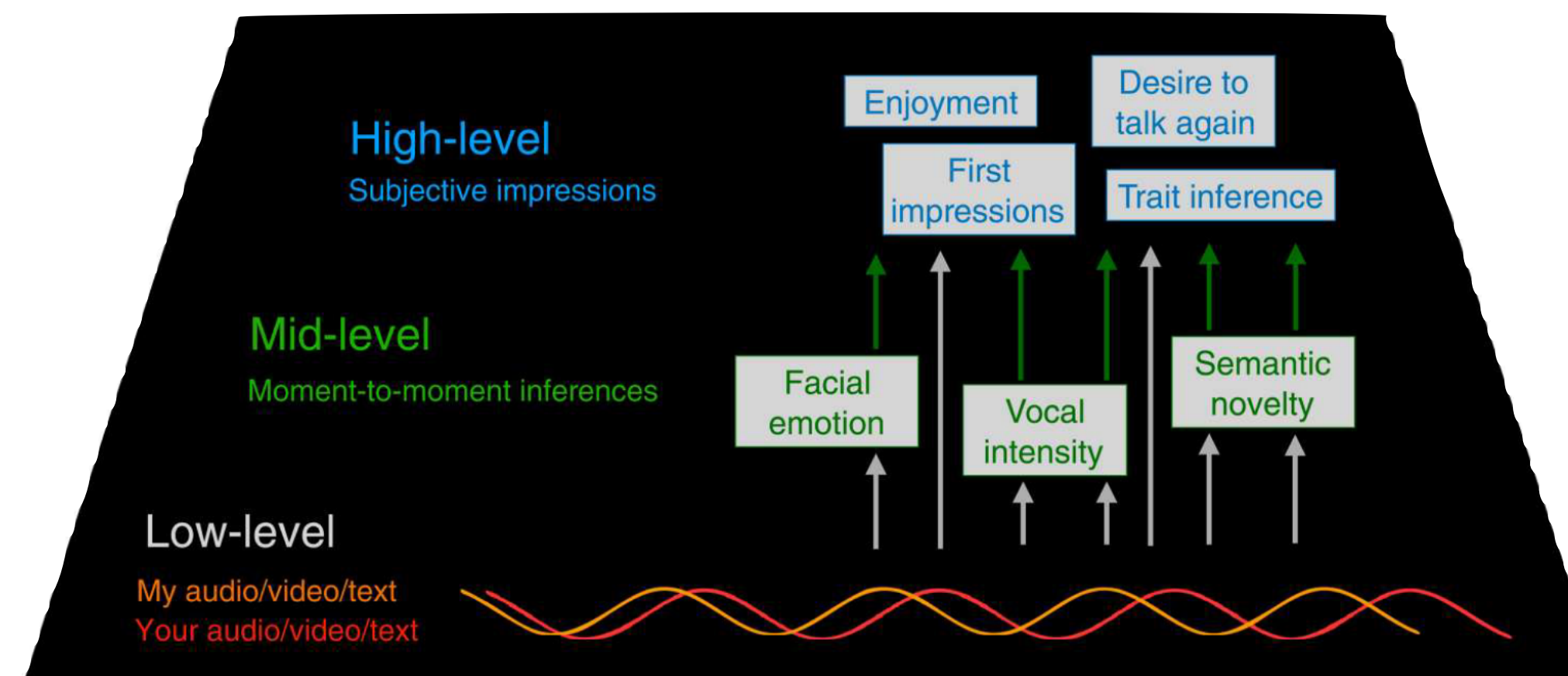


# What you should know

- Multimodal data is a new data genre that fits into the data mining pipeline.
- Roughly categorized as low-, mid-, high-level features.
- Mid-level features offers the richest opportunities.
- Conversation analysis is centered on the turn-taking mechanism.
  - Data mining techniques and domain knowledge should go hands-in-hands.

# Discussion: From Conversation to Classroom

- What is equivalent to a “turn” in classroom?
- What would be the low-, mid-, and high-level features in the classroom setting?
- What can or cannot be applied?



# Demo: Audio Feature Extraction with OpenSmile

# Assignment - Option 1

- Reproduce a emotion intensity classifier by following the approach described in the CANDOR paper
  - Using the RAVDESS dataset, and use “emotional intensity” as the label.
  - Compute a series of acoustic summary statistics—mean, maximum, and standard deviation for fundamental frequency, log energy, and voiced and unvoiced duration—as the features.
  - Train a logistic regression model (w/ or w/o feature normalization and regularization)

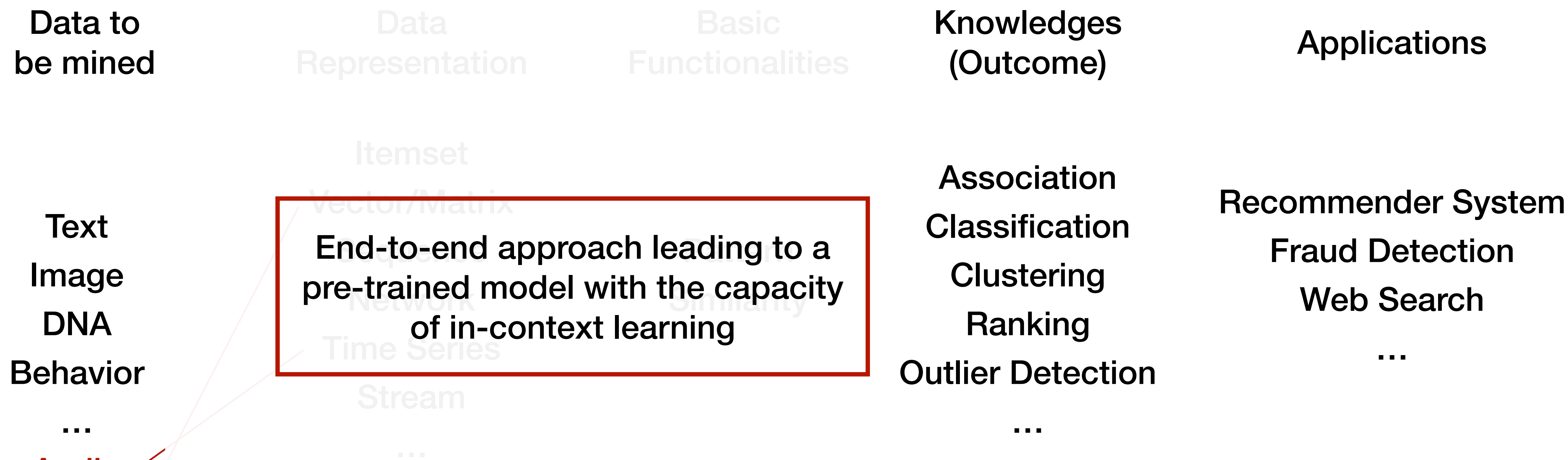
# Assignment - Option 2

- Analyze how acoustic features contribute to emotion recognition, including:
  - Using the RAVDESS dataset, and use "emotional intensity" as the label.
  - Fundamental frequency (F0) as a measurement of pitch
  - Intensity/amplitude measurements (related to perceived loudness)
  - Frequency perturbation (jitter) and amplitude perturbation (shimmer)
  - Additional features such as spectral measures (e.g., MFCCs, formants) and voice quality metric
- Easier than Option 1 as it does not require training models.

# Bonus Slides

How does LLM change the game?

# Multimodal Foundation Models






we no longer care about the data representation and basic functionalities, we pre-trained the model and expect that it can perform (generalize) well on various task, through either fine-tuning or in-context learning

# (Large) Audio-Language-Model

- Audio processing is challenging due to data scarcity, but language is abundant.
- ALM can learn a shared latent space between audio-language modalities.
- In the post-LLM era, LALMs integrate audio encoders with LLMs, unlocking deeper audio perception and reasoning.
  - Capable of captioning (ASR), reasoning, and open-ended QA.
  - Examples:
    - Gemini, GPT-4o,
    - [open access] Qwen-audio, Qwen-audio-chat, Qwen-audio-instruct,
    - [open source] GAMA, SALMONN, GAMA, etc.






# Capability of Large Audio Language Models






recognize the speech, and translate it into English

OK. The sentence “大家好我们来讨论一下公司周末团建的地点吧大家有什么意见呢” translated into English is “Hello, let’s discuss the location for our company’s weekend team building. What are your opinions?”



recognize the speech, and translate it into English







OK. The sentence “I think we can go to the countryside and find a beautiful natural area.” translated into English is “I think we can go to the countryside and find a beautiful natural area.”



recognize the speech, and translate it into English

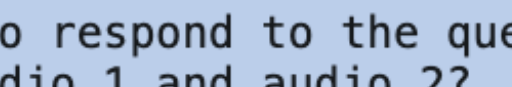
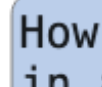

OK. The sentence “わかりました。田舎に行くことに同意します。田舎暮らしを体験してみるのもいいかもしれません。” translated into English is “I understand. I agree to go to the countryside. It’s also good to experience country life.”

(a) Multilingual ASR and Translation



What are the emotions in audio 1 and audio 2, respectively?

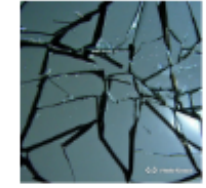



The tone of the woman in audio 1 is cheerful, possibly talking about something pleasant. The tone of the woman in audio 2 is fearful, possibly discussing something terrifying.



How to respond to the questions in audio 1 and audio 2?

In audio 1, when the woman asks if you're okay, you can respond with, "Thank you for your concern, I'm fine now, just felt a bit uncomfortable all of a sudden."  
In audio 2, considering the woman's fearful tone and discussing something terrifying, you can respond to her with either, "I'm scared, I don't know what to do" or "I'm feeling really upset, I don't know how to face this situation."

(b) Multiple Audio Analysis





Recognize the sound information and assume that the sound is happening around the user. Please provide the user with some handling suggestions.

The sound is of breaking glass. Here are some suggestions for handling the situation:  
1. Turn off the source of the sound: If the sound is coming from a broken glass item, turn off the source of the sound as soon as possible to prevent further damage.  
2. Clean up the broken glass: Use gloves and a broom or vacuum cleaner to clean up the broken glass. Be careful not to cut yourself.  
3. Assess the damage: If the broken glass is in a location that is not easily accessible or if it poses a safety hazard, contact a professional for assistance.  
4. Cover the area: Cover the area where the broken glass is located to prevent injury or further damage.  
5. Avoid walking on broken glass: If the broken glass is on a floor, do not walk on it. Instead, use a broom or vacuum cleaner to clean it up or cover it with a piece of cardboard or tape.


(c) Sound Understanding and Reasoning




# Capability of Large Audio Language Models






Write a poem about the above audio



Birds are singing in the trees,  
A gentle breeze is blowing free,  
The rustling leaves are music,  
As nature's symphony goes on.  
The sun is shining bright and warm,  
A peaceful atmosphere is sown,  
In the midst of nature's beauty,  
We can find serenity and duty.  
The birds take flight with ease,  
Their chirping voices fill the trees,  
A symphony of nature's grace,  
A moment to pause and embrace.  
The wind whispers through the branches,  
As the leaves dance in the branches,  
A peaceful moment in time,  
A chance to let go and unwind.  
So let us take a moment here,  
To appreciate the beauty so dear,  
In the midst of nature's wonder,  
We can find a sense of wonder.





(d) Audio-Motivated Creative Writing



what is it


A solo piano piece.






Describe this music in detail

This is a solo piano piece. It is fast tempo with a playful and humorous mood. The piece is in the key of C major and is in 4/4 time signature. The melody is played on the right hand with the left hand playing chords. The piece is a classic ragtime style.






Give me some recommendations for similar music.

some similar music to "Solo Piano Ragtime" could be:

1. "The Entertainer" by Scott Joplin
2. "Maple Leaf Rag" by Scott Joplin
3. "Ave Maria" by Franz Schubert
4. "Für Elise" by Ludwig van Beethoven
5. "The Blue Danube" by Johann Strauss II



(e) Music Appreciation



change "joyful" to "interesting"

ok, here is the edited audio





change the person name to "John"

ok, here is the edited audio





Edit out the word "what", so I can turn it into a meme

The word "what" starts at 3.04 seconds and ends at 3.23 seconds.



(f) Speech Editing with SWRT

# Assignment - Option 3 [challenging]

- Evaluate audio-capable foundation models for emotion intensity recognition
  - Using the RAVDESS dataset, and use "emotional intensity" as the label.
  - Implement and evaluate one or more audio foundation models such as Whisper and Qwen-Audio
  - For Whisper, you will need to
    - Use Whisper's pretrained model as a feature extractor,
    - Extract embeddings from one of Whisper's encoder layers,
    - Add a classification head on top of the Whisper embeddings,
    - Fine-tune to predict emotion intensity levels
  - For Qwen-Audio, you will need to
    - Practice deploying the Qwen-Audio in a local environment, write prompt to obtain the predicted emotion intensity.