# Towards Neural Encoding of Higher-Order Cognition with Large Language Models

Sean Mata

Advisor: Professor Peter Ramadge, Professor Uri Hasson

Submitted in partial fulfillment

of the requirements for the degree of

Bachelor of Science in Engineering

Department of Electrical and Computer Engineering

Princeton University

May 2025

I hereby declare that I am the sole author of this thesis.

I authorize Princeton University to lend this thesis to other institutions or individuals for the purpose of scholarly research.

_____

Sean Mata

I further authorize Princeton University to reproduce this thesis by photocopying or by other means, in total or in part, at the request of other institutions or individuals for the purpose of scholarly research.

_____

Sean Mata

# Abstract

Large Language Models (LLMs) have demonstrated remarkable language understanding capabilities, significantly advancing neural encoding models. While recent progress has improved LLMs' reasoning abilities, they still face limitations in social contexts that require nuanced understanding of human cognition. This thesis investigates whether inference-time techniques like chain-of-thought reasoning—which externalize an LLM's reasoning process—can be leveraged for neural encoding of human reasoning and higher-order cognitive processing. Specifically, we focus on theory of mind and social reasoning within conversational contexts, examining the detection of subtext (unspoken beliefs, intentions, and desires) in communication.

Towards this we present two complementary projects: First, we develop a novel Theory of Mind-inspired Tree of Thought approach for dialogue generation that explicitly models subtextual reasoning. Our method produces diverse, contextually appropriate outputs, though automatic selection mechanisms currently struggle to outperform single-step generation. Through perturbation analysis and evaluator consistency tests, we gain insights into model confidence and the distribution of plausible responses. Second, we apply LLM-derived embeddings to neural encoding tasks, successfully modeling brain activity associated with both word- and sentence-level language processing using banded-ridge regression.Our findings suggest the potential for modifying sentence-level inputs to target the encoding of higher-order thinking in high-frequency Electrocorticography (ECoG) data. Together, these projects highlight both the promise and current limitations of LLMs in modeling social reasoning.

We conclude with proposals for future work, including improving training data, context structuring, evaluation metrics, and thought representation. Our findings point to new directions in aligning language models with human-like communicative competence and mental state inference.

# Acknowledgements

An uncountable number of positive influences have shaped my journey to this point. For me, this thesis marks not just the culmination of a project, but also the nearing end of my undergraduate years—years that have defined much of my young adult life.

Within the context of this project, I am deeply grateful to Itamar Jalon for his mentorship and unwavering support. Without his guidance, I would have been lost in how to even begin tackling this work. Our regular meetings were often a highlight of my week, filled with discussions about neuroscience, as well as general life and the future. Thank you to Peter Ramadge for your invaluable support and for your guidance in managing both the research and the writing process. I am also especially grateful towards Professor Uri Hasson for welcoming me into his lab and facilitating the research opportunities that made this thesis possible. To all the members of the Hasson Lab, thank you for making me feel welcome, answering my many questions, and creating an environment in which I could learn and grow.

Beyond this project, there are many more people I wish to thank. At Princeton University, I am deeply grateful to everyone who helped me build a sense of home within this "orange bubble." I'm proud of the community we created—rich with stories, memories, and shared experiences. From travels across states and countries—New York, Chicago, Puerto Rico, Mexico, and more—to countless campus events spanning sports, music, and community and awareness-building, to the joy of watching and performing in dance shows, each moment added depth and color to my time here. I look forward to the friendships that will continue to grow long after graduation, and I hope to cross paths with many of these friends again in the years to come.

To my family: thank you for your constant love and support. I miss you every time I find myself living in a different city or country. To my parents, no matter where

I go, I will always be your baby. To my siblings, I will always be there to support you.

Finally, to my friends from Oxford, secondary school, and throughout the years: thank you for the countless moments of growth, learning, and laughter we shared. We navigated so many phases of life together, and I look forward to continuing to grow and experience new chapters with you and all my dear friends.

# Contents

# List of Figures

# List of Symbols

$\mathcal{D} = \{D_1, D_2, \ldots, D_J\}$: Collection of dialogues, where each dialogue $D_j$ is a sequence of utterances between two speakers.

$D = \{u_1, u_2, \ldots, u_N\}$: A dialogue between two speakers, represented as a sequence of $N$ utterances.

$c_i = \{u_{i-w}, \ldots, u_{i-1}\}$: Context for utterance $u_i$, consisting of the $w$ preceding utterances, or $\min(i-1, w)$ if $i - 1 < w$.

$w$: The number of preceding utterances included in the context for generating each response. This can be referred to as the context length or width.

$C = \{c_{w+1}, c_{w+2}, \ldots, c_N\}$: Set of all contexts used for response generation, corresponding to turns $i = w + 1$ to $N$.

$R = \{r_{w+1}, r_{w+2}, \ldots, r_N\}$: Model-generated responses for turns $i = w + 1$ to $N$, conditioned on context $c_i$.

$L = \{l_{w+1}, l_{w+2}, \ldots, l_N\}$: Human-written reference responses for turns $i = w + 1$ to $N$.

*Note:* For turns $i \leq w$, the model does not generate responses due to insufficient context.

$S$: Set of scores produced by the evaluation metric, $\{s_1, s_2, \ldots, s_N\}$.

$Q$: Set of human quality annotations, $\{q_1, q_2, \ldots, q_N\}$.

$p_\theta$: A pre-trained language model (LM) with parameters $\theta$.

$x, y, z, s, \ldots$: Language sequences, where $x = (x[1], \ldots, x[n])$ represents a collection of tokens, with each $x[i]$ being a token.

$p_\theta(x)$: The probability of sequence $x$ under the pre-trained model $p_\theta$, given by:

$$p_\theta(x) = \prod_{i=1}^{n} p_\theta(x[i] \mid x[1 \ldots i]).$$

$\bar{X}, \bar{Y}, \bar{Z}, \bar{S}, \ldots$: Collections of language sequences, represented in uppercase with a bar.

$p_\theta(y \mid \text{prompt}_{IO}(x))$: The probability of generating output $y$ from input $x$ wrapped in a prompt with task instructions and/or few-shot input-output examples.

$p_\theta^{\text{prompt}}(\text{output} \mid \text{input})$: Simplified notation for the probability of output given the input with a prompt.

prompt: A prompt that contains task instructions and/or few-shot input-output examples used to generate the next response.

$T = \{\tau_0^{(m)}\}_{m=1}^{M}$: The collection of subtexts $\tau_0^{(m)}$, indexed by $m$ from 1 to $M$, where each $\tau_0^{(m)}$ represents a potential subtext idea based on a theory of mind taxonomy.

$G(p_\theta, s, k)$: Thought generator that generates $k$ candidate thoughts based on the pre-trained language model $p_\theta$ and tree state $s$.

$s = [x, z_1, \ldots, z_i]$: A tree state, where $x$ is the initial input, and $z_1, \ldots, z_i$ are the previous thoughts generated.

LLM- Large Language Model

ToM -Theory of Mind

# Chapter 1

# Introduction and Background

Large Language Models (LLMs) have recently exploded in popularity due to their ability to model language and their emerging, but imperfect, capacity for reasoning [2]. One application of LLMs is in human neuroscience, where they have been used to model neural signals from the literal content of words spoken [3]. This thesis aims to leverage recent advances in LLM inference-time reasoning to not only model the literal content of words during language production and comprehension but also subtext, meaning possible unspoken thoughts and inferences that arise during communication. Towards investigating this, this thesis presents two projects. The first applies an in-context learning technique, tree of thought, to a novel domain by aiming to improve dialogue generation over a naturalistic dataset by reasoning over potential subtextual thoughts — unspoken beliefs, desires, or intentions — inspired by theory of mind. The second builds on prior neural encoding projects in the Hasson lab, developing models that leverage both word-level and sentence-level embeddings to encode neural activity. The rest of this section will give a summary of related work and an outline of this report.

## 1.1 Reasoning versus language abilities in LLMs

It is under debate whether LLMs can truly reason in an abstract way to solve problems or whether they learn narrow-non-transferable skills for task-solving [4]. Nevertheless, they have had increasingly strong performance on reasoning benchmarks such as math, logic and reasoning [5] as well as language [6, 7]. This leads to the question of whether intelligence can be gained solely through mastery of language. This is related to a common fallacy conflating language and thought ability, where the ability to use language is neither a necessary nor sufficient condition for reasoning (Mahowald, Ivanova et al. coined the difference between knowledge of lingustic rules and understanding of language as formal linguistic competence and functional linguistic competence [8]). In humans, language and reasoning have been suggested to use different brain networks [9, 10], and individuals with aphasia who exhibit severe linguistic deficits still have other intact cognitive abilities such as social and mathematical reasoning [11]. Nevertheless, language can still bootstrap cognition, facilitating processes such as inner speech[12] and enabling faster development of social reasoning[13]. However, it seems unlikely that language models can acquire robust reasoning ability solely through next-token prediction during pre-training -motivating the development of various inference-time techniques aimed at enhancing LLM reasoning.

## 1.2 Inference time generation methods

Inference time generation methods are a method to improve the performance of LLMs without the large amounts compute and data needed for pre-training. This can be used alongside other methods of post-training improvement such as fine-tuning, reinforcement learning from human feedback and pure reinforcment learning in models such as o1 by OpenAI [14], R1 by DeepSeek [15], and QwQ by Qwen Team[16].

Inference-time methods come in many flavours, they can explore a single reasoning path often by optimizing a single prompt, such as in Chain of Thought [17, 18], multiple reasoning paths with or without multiple agents as in multi-agent debate [19] or boosting of thoughts [20] or can include an iterative or memory mechanism as in Reflexion [21] or buffer of thoughts[22]. However, there is no technique that consistently outperforms others in all tasks, so the implemented method must be chosen according to the specific task [23].

In this project, the strategy implemented is tree of thought [24]. This has the benefit of having multiple reasoning chains to track different ideas, as well as the flexibility to modify the tree's complexity and structure as needed for deeper or shallow levels of reasoning. Following an analysis of bias, a majority voting mechanism was also implemented [25].

## 1.3 Theory of mind in LLMs

Theory of Mind (ToM) refers to the cognitive capacity to infer the mental states—such as beliefs, desires, and intentions—of others, allowing us to predict and interpret their behavior. Traditionally, ToM has been modeled as a predictive model of others' actions based on their mental states. However, this approach often struggles with flexible generalization to novel situations and is limited in its use for complex planning tasks that involve many possible actions. A more powerful alternative frames ToM as a causal model—an abstract, structured representation of mental states that supports efficient reasoning, planning, and decision-making over a wide space of possibilities [26].

One domain where this richer model of ToM becomes essential is pragmatic language use. For instance, speakers frequently tailor their words to achieve specific emotional or social outcomes, such as using indirect or softened language to avoid

offending someone. This type of communication requires balancing honesty, social norms, and an estimation of the listener's mental state. The Rational Speech Act (RSA) framework captures this complexity by modeling communication as a Bayesian reasoning process. Within RSA, speakers choose utterances by anticipating how listeners will interpret them, and listeners interpret utterances under the assumption that speakers are cooperative and informative. ToM is also central in interpersonal affect regulation, where individuals deliberately attempt to influence others' emotional states—calming them down, cheering them up, or helping them regain control. This kind of planning involves both empathy and strategic reasoning, highlighting the real-world complexity of applied ToM.

The extent to which large language models (LLMs) can perform ToM reasoning is the subject of ongoing debate [27, 28]. To address this, researchers have introduced a number of benchmarks to systematically evaluate ToM abilities in LLMs [29, 30, 31, 32]. While LLMs show strong results on certain tasks, their overall ToM performance remains inferior to that of humans, and is often accompanied by signs of overfitting [33].

Current datasets for evaluating Theory of Mind (ToM) in large language models (LLMs) suffer from notable limitations. Rich, human-authored narratives that naturally involve ToM reasoning are both rare and costly to curate. Consequently, researchers often rely on synthetic, template-based datasets that are overly simplistic and lack ecological validity. These datasets frequently include explicit mental state language (e.g., "Amy thinks that..."), which makes the reasoning task unrealistically easy and reduces the diagnostic value of the evaluation. Moreover, they tend to focus narrowly on action prediction, failing to capture the broader and more nuanced reasoning processes involved in real-world ToM. One dataset that attempts to address these shortcomings is SimpleToM [34], which includes both types of reasoning: the recognition of mental states (explicit ToM) and the application of that knowledge

to predict actions or responses (applied ToM). They found that while LLMs tend to perform well on explicit ToM tasks, their performance on applied ToM remains weak—highlighting a critical gap and a central challenge relevant to our domain. Additionally, ToM performance in LLMs can vary significantly depending on the scenario, highlighting the need for methods that enhance LLMs' capacity for robust, context-sensitive reasoning rather than relying solely on narrow or superficial cues.

Several methods have been proposed to enhance Theory of Mind (ToM) performance at inference time [33, 35], but many of these approaches rely on few-shot examples or rigid assumptions that hinder scalability. While techniques such as custom chain-of-thought prompting can guide models toward more accurate ToM reasoning, they are often fragile and require extensive manual tuning [34]. Structured frameworks have also been introduced to support ToM reasoning. For instance, Wu et al. [36] developed COKE, a system that reasons over a knowledge graph by instantiating ToM as a collection of manually verified cognitive chains. These chains represent human mental activities in specific social contexts, including their behavioral and affective responses. More recently, Kim et al. [37] proposed a similar approach to the one in this project, propagating a series of hypothetical thoughts to explain and infer human actions with a method inspired by the sequential Monte Carlo algorithm.

## 1.4 Use of Language Models to study processing in the brain

A computational model's neural plausibility can be evaluated by assessing whether its internal representations align with patterns of brain activity elicited by corresponding stimuli. If a model consistently predicts neural responses to novel inputs, this suggests it captures key aspects of the brain's representational structure. This approach has been used in decoding brain activity in the visual pathway, using techniques

such as motion energy models, Bayesian decoders, and deep learning frameworks [38]. Similarly, large language models (LLMs), which have achieved remarkable success in modeling human language, raise the question: how closely do their internal representations mirror those of the human brain?

To quantify representational similarity between LLMs and the brain, researchers typically employ two types of models: encoding models, which assess how well LLMs predict neural activity in response to linguistic stimuli, and decoding models, which infer the stimulus that produced a given neural response [39]. Empirical studies have revealed overlaps in representational structure—particularly between contextual embeddings from LLMs and decontextualized word embeddings. For example, fMRI studies [40] and intracranial recordings [41, 42] indicate that contextual embeddings better align with brain representations than static embeddings.

To deepen our understanding of this shared structure, Tuckute et al. [39] propose two key strategies for leveraging LLMs in the study of brain function. The first involves systematically probing LLMs by varying their architecture, training data, and behavioral outputs to isolate factors that align with brain activity. The second involves using LLMs as tools for in silico neuroscience experiments, allowing researchers to simulate and test hypotheses derived from empirical data.

These approaches have been extended to higher-level cognition. For instance, one study used fMRI to measure brain responses to 1,000 diverse sentences and showed that an encoding model based on GPT could predict the magnitude of the neural response to each sentence [43]. Interestingly, for sentences describing others' mental states, the model did not explain significantly more variance (beyond surprisal), suggesting that language-related brain regions are not modulated by social content. This finding supports the view that the brain's language network is functionally distinct from the theory of mind network.

One intriguing direction for future research is whether in-context learning—a

method known to improve reasoning in LLMs—could be used to simulate or even enhance models of human reasoning. Such approaches may offer novel tools to investigate the mechanisms of cognitive processing in the brain.

### 1.4.1 Sentence embedding encoding

When working with fMRI signals, sentence-based encoding and decoding is more common due to the longer time duration of the signals. Early research demonstrated the ability to encode sentences using embeddings [44], and further analysis has shown that syntactic structure, rather than lexical-semantic content, is the primary contributor to the similarity between artificial neural network representations and brain responses in the language network [45]. Sentence embeddings are typically generated by averaging the word embeddings in the sentence or using the representation of the last sentence token (often the final period token ".") as a sequence summary.

## 1.5 Report Overview

In the following chapters, we will review the relevant literature and methodological background for natural language processing using transformers and dialogue modeling metrics, go over the methodology and discuss experimental results on the naturalistic 24/7 Conversations dataset, including future work and current limitations. Finally, we go over the method and results for the sentence encoding project and explore strategies for integrating it with the dialogue modeling work presented in this report.

# Chapter 2

# Methodological Background

## 2.1 Natural Language Processing using Deep Neural Networks

Natural language processing involves many different tasks which includes classification, text generation and sentence pair matching. A variety of deep neural network models can be used, including recurrent neural networks, long short-term memory networks, and most commonly now transformers. Modern transformer-based models based on the original architecture by Vaswani et al., 2017[1] give us the best performance on the relevant tasks for this project which requires text generation and sentence similarity scoring.

An overview of the transformer architecture is given in Figure 2.1. While a detailed account is covered in textbooks [46], the following section will give an overview of the key aspects of the models that we will repeatedly refer to.

**Tokenisation**

In NLP, a token is a single unit of text, such as a word, subword, or character, depending on the tokenization method used. For example, in the sentence "I love

Figure 1: The Transformer - model architecture.

Figure 2.1: Transformer Architecture from Attention is All You Need (Vaswani et al., 2017)[1]

NLP!", the tokens could be ["I", "love", "NLP", "!"] (word-level tokenization) or ["I", " lo", "ve", " NLP", "!"] (subword tokenization like in Byte Pair Encoding). Tokenization is a fundamental step in processing text for machine learning models. Different models may have different tokenization schemes with different vocabulary sizes and elements.

**Static Embeddings**

Static embeddings map each unique word in a vocabulary to a dense d-dimensional vector that capture meaning and similarity. This can be thought of as an improvement

over one-hot encoding, where each word is represented as a large sparse vector with no meaningful relationship between each vector. In static embeddings, words with similar meanings have similar vectors and tend to cluster together in the embedding space. Furthermore, certain directions within this space can even capture specific relational meanings, for example, the vector difference between king and queen is similar to that between man and woman. The key attribute of static embeddings is that each word has a single fixed vector, regardless of its surrounding context. Common methods of creating these embeddings include Word2Vec and GloVe.

**Contextual Embeddings**

By contrast, with contextual embeddings, such as those learned by masked language models like BERT, each word w will be represented by a different vector each time it appears in a different context. The embeddings are made contextual through the self-attention mechanism.

The attention mechanism then measures how much focus each token should give to every other token. This is done by computing the similarity between the Query $(Q)$ of a token and the Keys $(K)$ of all tokens through the dot product $QK^T$. In the decoder only part this computation is masked, so can only look at the past tokens. The result is then scaled by $\frac{1}{\sqrt{d_K}}$ to normalize the values, helping stabilize gradients during training. Then the softmax function normalises the attention scores to ensure they sum to 1, converting raw similarity scores into a probability distribution. The result of this determines how much weight each token should assign to others in the sequence, and finally these computed attention weights are applied to the Value $(V)$ vectors. Tokens with higher attention scores contribute more to the final representation of a given token. The weighted sum of values produces a new representation for each token, capturing contextual dependencies dynamically. This mechanism can be applied multiple times sequentially and for multiple heads. So when we refer to

"contextual embeddings" this refers to weighted value vector after attention is applied using the $Q$ and $K$ matrices.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right) V \qquad (2.1)$$

**Next token prediction models**

For generation of dialog I will use a decoder-only transformer e.g. DeepSeek, Mistral, GPT. In next-token prediction models like decoder-only transformers, text is generated one token at a time based on previous context. These models are also trained autoregressively. After processing input tokens, the hidden state of the last token in the final layer $h_N \in \mathbb{R}^{1 \times d}$ is passed through a linear layer using an unembedding matrix $W \in \mathbb{R}^{d \times V}$ (called the language modeling head) to produce *logits* $z_i \in \mathbb{R}^{1 \times V}$ — raw scores for each token in the vocabulary.

$$z_i = h_N W$$

These logits are converted to probabilities using the softmax function, producing $P(\text{next token} \mid \text{context})$:

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{j=1}^{V} e^{z_j}}$$

The model can then either pick the most likely token (greedy decoding) or sample from this distribution (temperature sampling/ to generate the next token.

**Masked word prediction**

As opposed to autoregressive decoder models, masked language models like BERT operate by randomly masking words in sentences and training the model to predict these missing words using surrounding context. This approach falls under denoising

training objectives, where inputs are deliberately corrupted through masking, substitution, reordering, deletion, or insertion, and the model learns to reconstruct the original input by minimizing cross-entropy loss between its predictions and the true values of the missing tokens.

This bidirectional approach enables learning rich contextual word representations applicable to various NLP tasks such as text classification, named entity recognition, and sentence similarity assessment. Often called encoder-only models, they excel at producing contextualized embeddings for interpretative tasks rather than generating text, making them particularly valuable for methods like BERTScore.

## 2.2 Text similarity metrics

For this project, we need a reference-based metric(s) to compare the similarity of the model generated response to the human reference response. In practise, human evaluation (typically from crowd-workers over experts for dialog) is used a gold standard for text and dialog evaluation with various strategies used for data collection and different questions asked to evaluators [47]. However, collecting this data from humans is time-intensive and expensive, so it is typically only used for a final evaluation, while automatic evaluation metrics are used during development to optimise model design and the choice of hyperparameters.

The ideal automatic metric maximally correlates with human judgement, however most metrics struggle to achieve this goal [48] hence metric-creation is still an open problem itself[49][50]. To improve accuracy or cost a hybrid approach may also be used, potentially combining multiple automatic metrics [51][52] or human and automatic metrics [53][54]. There have been several reviews on the several text similarity metrics [51] [55], and here we summarize their findings.

Formally, we can describe dialog evaluation metrics as follows: Given a dialog con-

text $c = \{c_1, \ldots, c_m\}$, model response $r = \{r_1, \ldots, r_n\}$, and human-written reference response $l = \{l_1, \ldots, l_k\}$, the goal is to learn a function:

$$f : (c, r, l) \rightarrow s \tag{2.2}$$

that evaluates the generated response. (A metric is said to be reference-free if it can be expressed as a function of only $c$ and $r$.)

These metrics are assessed by comparing them to human judgment. Concretely, a human annotator (or several annotators) scores the quality of a given response conditioned on the dialog context:

$$(c, r) \rightarrow q. \tag{2.3}$$

Given the scores produced by a particular metric, $S = \{s_1, s_2, \ldots, s_n\}$, and the corresponding human quality annotations, $Q = \{q_1, q_2, \ldots, q_n\}$, we can measure the performance of the metric by calculating the correlation between $S$ and $Q$.

There are some important features and challenges to consider in the selection of this metric. The first is that the surface form representation is less important that the semantic meaning of the text, for example ""The cat sat on the mat.", "A mat was where the cat chose to sit.", ""The cat settled on the mat." , "The mat became the cat's resting place." should all receive similar similarity scores. Furthermore, the metric should be robust without changing greatly word small changes in word order or synonyms. Additionally, metrics can typically be evaluated at the utterance level as well as the system level by computing the average utterance-level score across many utterances [56], and having a strong correlation at the system-level is still useful even if the metric per utterance is noisy or has a lower correlation with human judgement.

Achieving these goals can be complicated, especially because of the long-tailed nature of language. You can not train a model on every possible combination of

sentences because there are too many, which is why supervised trained metrics tend to overfit the dataset they were trained on [57]. Additionally, it is important to consider all metrics will have limiations and optimising a particular metric can lead to unintended effects [58][59].

In this section we go over the various options for evaluation of natural language text generation for dialogue, and explain the workings of BERTScore -the final metric selected. In the following sections, existing metrics can be broadly categorized into using n-gram matching, edit distance, embedding matching, or learned functions.

### 2.2.1 N-gram matching

N-gram matching metrics are based on word or character based overlap between the prediction and reference. The most popular example is BLEU, with other metrics including METEOR AND NIST. While originally used in translation to detect the textual similarity when changing language, they can also be used for evaluating textual similarity within the same language.

As an example, the BLEU score for a corpus of candidate translation sentences is a function of the n-gram word precision over all the sentences combined as well as a brevity penalty computed over the whole corpus [46]. Here n-gram refers to a sequence of token-length n. N-gram precision would be the percentage of n-gram tokens in the candidate translation that also occur in the regerence translation. The harmonic mean of the unigrams, bigrams, up to 4-grams is used in the calculation of BLEUscore. Since it is a word-based metric, it is very sensitive to word tokenization, making it impossible to compare different systems if they rely on different tokenization standards. It also cannot detect semantic equivalence, so two sentences with the same meaning but mostly different characters can receive an incorrect low score. Because of this, it is not generally used for sentence similiarity evaluation [60] and it is not a good fit for our uses.

### 2.2.2 Edit distance

Minimum edit distance/word error rate is another common metric with its origin in automatic speech recognition. However, since it is commonly used in speech recognition this equally penalises the completely wrong word in a sentence as well as a close synonym. Therefore, it is also misses capturing some of the idea of the semantic meaning of sentences. This lack of semantic meaning makes it unsuitable for this project.

### 2.2.3 Trained metrics

Various metrics are trained to optimise correlation with human judgement. The input to these models can include other different metrics such as character n-grams in BEER which used a regression model or ADEM which uses the context and the reference response as inputs to a recurrent neural network trained to predict appropriateness ratings by human judges. These methods usually have the highest correlation to human judgement [61]. However, these methods require expensive human judgments as supervision for each dataset and also generally have poor generalisation to new domains and data. These metrics have been shown to not be robust [57]. For example with ADEM, in 48.66% of cases the predicted score increased when the generated response was reversed. In 86.93% of cases the predicted score increased when the generated response was replaced with a dull dummy response.

### 2.2.4 Perplexity

Perplexity is a metric which was first introduced by Jelinek et al. in 1977 while working on automatic speech recognition [62], which was its original use case [63]. However, now it is most commonly now it is used to decide how to sample from a language model for text generation [64][65].

Given the probability assigned by the language model to a sequence $\mathbf{x}$ as:

$$p_\theta(\mathbf{x}) = \prod_{i=1}^{n} p_\theta(x[i] \mid x[1 \ldots i-1])$$

Perplexity (PPL) can be derived from this probability as;

$$\text{PPL}(\mathbf{x}) = p_\theta(\mathbf{x})^{-\frac{1}{n}} = \exp\left(-\frac{1}{n}\sum_{i=1}^{n} \log p_\theta(x[i] \mid x[1 \ldots i-1])\right)$$

where $n$ is the sequence length. Perplexity is highly interpretable, because it quantifies how "surprised" the model is by the given sequence—lower perplexity means the model assigns higher probability to the observed sequence. Perplexity can range from a theoretical value of 1, where the language model is completely certain of its prediction to infinity, where the language model is completely uncertain.

One example of its use in sampling is the maximum mutual information (MMI) scoring function [66] [67] has been used to predict such "boring" responses. MMI employs a pre-trained backward model to predict source sentences from given responses, i.e., $P(Source|target)$. The probability of $P(Source|Hypothesis)$ can be used to rank the quality of hypotheses. Intuitively, maximizing backward model likelihood penalizes the bland hypotheses such as "that's interesting" or "ok", since the ability to discern the source sentence from such a target sentence is low, because these hypothesises can be associated with many possible queries. Another use case has been to evaluate the coherence of two consecutive sentences [68] for .

Perplexity is not as commonly used for evaluation of dialogue models. One reason for this is that this metric must be model-dependent, and it is better to avoid situations where the model generating the response also evaluates its quality because this would lead to bias [48]. Perplexity is also not comparable between language models with different vocabularies [63]. Moreover, perplexity is slightly different to the other metrics because it does not directly measure textual or semantic similarity between

two samples. Even if two texts have similar perplexity values, they could be semantically very different. It could still be used to evaluate which subtext or thought (or no thought) most likely results in a given human reference, however it does not give us a direct measure of how well a model response matches the reference response. Hence, I did not prioritize implementing this metric.

### 2.2.5 Embedding based metrics

As discussed in Section 2.1, token (word) embeddings are learned dense token (word) representations that can be represented in a continuous vector space where semantic similarity is encoded so that tokens (words) with similar meanings are closer in the vector space. Embedding based metrics tend to have the highest correlation to human judgment based on their ability to match semantic meaning between candidate and reference sentences. Hence, this is what we chose to use in the rest of this report.

The general method of computing similarity is between two tokens is using the cosine similarity between two embedding vectors $v$ and $w$ as:

$$\text{cosine}(v, w) = \frac{v \cdot w}{\|v\| \|w\|}$$

As mentioned in pervious section, these embeddings can either be static and learned from large text corpora such as Word2Vec or GloVe, or they can be from contextual models like BERT. This is important for measuring semantic similarity of texts, since contextual embeddings allow the same word to have a different representation depending on its surrounding words, allowing the embeddings to capture nuanced meanings, polysemy, and relationships within the sentence context. For example, the word "bank" will have different embeddings in "river bank" versus "financial bank," reflecting its specific usage in each sentence. Static embeddings will measure the similarity of the word "bank" in these two sentences as exactly the same. Hence, using

contextual embeddings for this task tends to do better[69].

The specific metric chosen for this project is BERTScore based on its superior performance for measuring textual similarity[70][69]. This metric is based on BERT, which is an encoder that creates contextual embeddings. Encoder-only models like BERT are known for producing high-quality, robust embeddings, making BERTScore an ideal choice for this task.

To explain how BERTScore works in detail, the BERTSCORE algorithm first passes the reference sentence $x$ and the candidate sentence $\hat{x}$ through BERT, computing a BERT embedding for each token $x_i$ and $\hat{x}_j$.

Given this sequence of tokens in a reference sentence $x = \langle x_1, \ldots, x_k \rangle$ and a candidate sentence $\hat{x} = \langle \hat{x}_1, \ldots, \hat{x}_l \rangle$, the similarity is calculated using a greedy matching algorithm.

Each token in $x$ is matched to a token in $\tilde{x}$ to compute recall, and each token in $\tilde{x}$ is matched to a token in $x$ to compute precision (with each token greedily matched to the most similar token in the corresponding sentence). BERTSCORE provides precision and recall (and hence $F_1$):

$$R_{\text{BERT}} = \frac{1}{|x|} \sum_{x_i \in x} \max_{\tilde{x}_j \in \tilde{x}} (x_i \cdot \tilde{x}_j)$$

$$P_{\text{BERT}} = \frac{1}{|\tilde{x}|} \sum_{\tilde{x}_j \in \tilde{x}} \max_{x_i \in x} (x_i \cdot \tilde{x}_j)$$

$$F_{\text{BERT}} = \frac{2 \cdot P_{\text{BERT}} \cdot R_{\text{BERT}}}{P_{\text{BERT}} + R_{\text{BERT}}}$$

Additionally, BERTScore includes baseline scaling. From the formula mentioned above, scores would have the numerical range of cosine similarity between -1 and 1, however scores are observed in a smaller range in practice. To improve readability, BERTScore rescales its output with respect to an empricial lower bound $b$. This is

computed by finding the average score between 1M candidate-reference pairs made by grouping two random sentences in the Common Crawl monolingual datasets. This lower bound is then used to rescale BERTScore linearly. The rescaled BERTScore is given as:

$$\hat{R}_{\text{BERT}} = \frac{R_{\text{BERT}} - b}{1 - b}$$

# Chapter 3

# Methodology

In this chapter, we will formalize the dialog generation task and describe the methods implemented to improve on this task. Then I will outline the generation strategies used to create a baseline score and the tree-of-thought and majority voting strategies implemented to aim to improve upon this score.

## 3.1 Problem setup

Given a dataset of conversational turns, we aim to construct a set of dialogues between two speakers. Each sample consists of a context, represented as a sequence of previous dialog turns $c_1, \ldots, c_m$, and a label, which is the true next utterance written by a human, denoted as $l_1, \ldots, l_k$.

The objective is to improve the ability of a pre-trained language model $p_\theta$ to generate a model response $r_1, \ldots, r_n$ that closely matches the human label response, $l_1, \ldots, l_k$. This similarity is measured using a metric such as BERTScore, which produces a set of scores $S = s_1, s_2, \ldots, s_n$ comparing the generated response to the human-written label.

Our approach is to enhance the model's generation process by conditioning its response on reasoning paths which are guided by prompts that consider a set of po-

tential subtexts, $\tau_0^{(i)}{}_{c=1}^{M}$, which should help the model generate output that reflects the thought process the human speaker likely followed. By incorporating these reasoning paths into the model's response generation, we aim to produce more coherent and contextually appropriate responses that better match the human label. If successful, this would not only improves the quality of the model's output in natural dialogue settings but also demonstrates the model's ability to understand and reflect the underlying mental states involved in conversation.

## 3.2   Generation strategies

This work explores two distinct next-utterance generation approaches. The first is a straightforward baseline that conditions only based on previous context. The second seeks to improve upon this baseline by adopting a tree of thought reasoning framework, in which the model explicitly reflects on the next speaker's possible mental states and intentions before generating a response.

### 3.2.1   Single step utterance generation

This section describes our procedure for generating single-step utterances using large language models, and introduces the notation we use throughout. This notation based on the original Tree of Thought paper [24], with adaptations to suit the specific context of our work.

We begin by formalizing the general language generation process, before describing how we apply this framework to obtain a baseline generation score without any in-context learning. Let $p_\theta$ denote a pre-trained LM with parameters $\theta$, and lowercase letters $x, y, z, s, \ldots$ to denote a language sequence which is a collection of tokens, i.e., $x = (x[1], \ldots, x[n])$ where each $x[i]$ is a token. Then $p_\theta(x) = \prod_{i=1}^{n} p_\theta(x[i] \mid x[1 \ldots i])$. We also use uppercase letters with a bar $\bar{X}, \bar{Y}, \bar{Z}, \bar{S} \ldots$ to denote a collection of

language sequences.

We can turn a problem input $x$ into an output $y$ using a language model by wrapping input $x$ with a prompt containing task instructions and/or few-shot input-output examples. We can write this as $y \sim p_\theta(y \mid \text{prompt}_{IO}(x))$. We can simplify this notation further by saying that $p_\theta^{\text{prompt}}(\text{output} \mid \text{input}) = p_\theta(\text{output} \mid \text{prompt}(\text{input}))$

Now we can describe the specific implementation approach. The relevant variables here are the context $c$ and the prompt prompt . We can tie this together as:

$$r_1^{(i)} \sim p_\theta^{\text{prompt}} \left( r_1^{(i)} \mid c \right) \tag{3.1}$$

Figure 3.1 depicts what this generation process would look like where the context $c$ consists of a single preceding utterance. We will create our baseline based on the generation procedure outlined in Equation 3.1. The prompt will be manually iterated to produce the best performance while the effect of changing context length will be included in the results.



Figure 3.1: An illustration of the single-step generation procedure, where the model generates a response conditioned using a context length of 1 (only on the previous utterance in the conversation.)

### 3.2.2 Tree of thought



Figure 3.2: A diagram representing the Tree of Thought and evaluation methodology for utterance generation. Orange highlighted cells represent the actual states searched through using the LLM. Yellow highlighted cells represent the set of potential model responses. Red arrows represent LLM generation while the dashed arrows represent comparison using BERTScore.

In this work, we adapt the original Tree of Thought (ToT) framework to the domain of dialogue generation. To structure our method clearly, we follow the organization of the original ToT paper, which frames a specific instantiation of ToT around four key design questions: (1) How should the reasoning process be decomposed into intermediate thought steps? (2) How should potential thoughts be generated from each state? (3) How should states be heuristically evaluated? (4) What search algorithm should be used to navigate the tree? In the next section, we describe how we adapt and extend each of these components to address the challenges specific to dialogue generation.

At a high level, our method generates intermediate thoughts based on a fixed

collection of subtexts of size $M$, which guide the generation of the final dialogue response. If all branches were explored, this would results in a tree of height 2 and width $N$ as depicted in Figure 3.2, with the final layer corresponding to $M$ candidate responses. However, in practice, the state evaluator $V(p_\theta, S)$ and search algorithm will be used to identify and follow the most promising branch without exploring the whole tree. All the explored states are highlighted in orange in Figure 3.2.

## 1. Thought Decomposition

While Chain of Thought samples thoughts without explicit decomposition, Tree of Thought designs and decomposes intermediate thought steps using the problem properties. In the following section, I will describe how the model uses a collection of subtexts $T = \{\tau_0^m\}_{m=1}^M$ to guide the reasoning process along specific paths. However, I will not focus on the decomposition of these thoughts, as we use a simple, low-depth tree structure, making further decomposition less necessary. This is because there is only one intermediate thought step before the model generates its response. Future work could explore increasing the depth of the tree to capture higher-order theory of mind or, alternatively, simplifying the thought process, as proposed in Chain of Draft [71].

## 2. Thought Generation

We require a thought generator $G(p_\theta, s, k)$ that generates $k$ candidate thoughts based on the pre-trained language model $p_\theta$ and tree state $s$. Given a tree state $s = [x, z_1 \ldots z_i]$, there are two strategies in the original paper to create a generator to create $k$ candidates for the next thought step:

(a) **Propose thoughts sequentially using a "propose prompt":**

$$[z^{(1)}, \ldots, z^{(k)}] \sim p_\theta^{\text{propose}}(z_{i+1}^{(1\ldots k)} | s).$$

24

(b) **Sample i.i.d. thoughts from a CoT prompt**:

$$z^{(j)} \sim p_\theta^{\text{CoT}}(z_{i+1}|s) = p_\theta^{\text{CoT}}(z_{i+1}|x, z_1 \ldots z_i), \quad (j = 1 \ldots k).$$

Proposing thoughts sequentially using a propose prompt works better when the thought space is more constrained. However, dialogue presents an inherently open thought space, suggesting that sampling i.i.d. thoughts from a CoT prompt would be more appropriate. The standard Tree of Thought approach would sample diverse reasoning paths independently like this, although sampling thoughts in this way may not lead to results that are interpretable from a theory of mind perspective.

To address these challenges, we instead initialize the first layer of tree states $S_0 = s_0^{(1\ldots k)}$ by adding to the context $c$ with a collection of subtexts $T = \{\tau_0^{(i)}\}_{i=1}^M$, which reflect different mental states or intentions the human speaker might have had. Specifically, each initial state is defined as $s_0^{(i)} = (c, \tau_0^{(i)})$. This can be interpreted as creating a manual thought generator $G_0(p_\theta, s_0 = c, k = M)$ which generates $M$ new nodes given a context $c$ based on the set of potential subtexts $T$, however this distinction is a matter of perspective rather than mechanism.

With the issue of creating thought diversity handled, the primary thought generator then does not need to calculate $k$ new candidates given a tree state, instead it effectively sets $k = 1$ and is primarily concerned with propagating a node forward to produce intermediate thought(s) and then the final model response.

The final tree structure uses one intermediate thought which results in a tree structure of width $M$ and height 2, which is summarized in Equation 3.2 and 3.3 below. Note that the prompts, and therefore the thought generators, for producing the intermediate thought and final model response are distinct because they create fundamentally different types of content. This contrasts with the original Tree of Thought framework, which does not explicitly distinguish between the generation of

intermediate thoughts and final output, instead treating all nodes as being of the same type.

$$z_1^{(i)} \sim p_\theta^{\text{prompt1}}\left(z_1^{(i)} \mid \tau^{(i)}, c\right) = G_1(p_\theta, s_0^{(i)}, 1) \tag{3.2}$$

$$r^{(i)} = z_2^{(i)} \sim p_\theta^{\text{prompt2}}\left(r_2^{(i)} \mid z_1^{(i)}, \tau^{(i)}, c\right) = G_2(p_\theta, s_1, 1) \tag{3.3}$$

In this implementation discussed in the results, a set of six subtexts are used, based on a framework of self versus other for feelings, goals and information seeking within a larger taxonomy of theory of mind reasoning [72].The selection of these subtexts are discussed in more detail in Section 1. The exact prompts are specified in Appendix XYZ.

## 3. State Evaluator

Given a series of different thoughts we also need a state evaluator $V(p_\theta, S)$ to evaluate the progress they make towards solving the problem. Traditionally, this is used as a heuristic for a search algorithm to determine which states to keep exploring and in which order. We use the language model to reason over these states/thoughts to choose which one to generate a response for. I implemented both strategies that the original paper provides to evaluate states either independently or together:

(a) **Value each state independently:**

$$V(p_\theta, S)(s) \sim p_\theta^{\text{value}}(v|s), \quad \forall s \in S,$$

where a value prompt reasons about the state $s$ to generate a scalar value $v$ (e.g., 1-10) or a classification (e.g., sure/likely/impossible) that could be heuristically turned into a value.

(b) **Vote across states:**

$$V(p_\theta, S)(s) = \mathbb{1}[s = s^*],$$

where the "best" state $s^* \sim p_\theta^{\text{vote}}(s^*|S)$ is voted for based on deliberately comparing different states in $S$ using a vote prompt. This is useful when problem success is harder to directly value.

We additionally define two reference state evaluators for comparison. The first is the random evaluator which can act as a baseline and the second is a perfect (oracle) evaluator which would represent an ideal upper bound.

(c) **Random state evaluator:**

$$V_{\text{random}}(S)(s) = \frac{1}{|S|}, \quad \forall s \in S,$$

where each state $s \in S$ is selected uniformly at random. This gives an equal value to each state, so that one is selected randomly at each level of the tree.

(d) **Perfect state evaluator:**

$$V_{\text{perfect}}(S)(s) = \mathbb{1}[s = s^*],$$

where:

$$s^* = \arg\max_{s \in S} \text{BERTScore}(\text{complete}(s))$$

and complete($s$) denotes the best possible full solution reachable from state $s$.

This always chooses the state that will end up maximising the BERTScore. While defining this function is very hard (and using LLMs to approximate this is a focus of this thesis), by propagating each branch of the tree forward to completion and backtracking we can deduce what the ideal state evaluator would have chosen.

## 4. Search Algorithm

Finally, within the Tree of Thoughts (ToT) framework, one can use different search algorithms depending on the tree structure. The original paper mentions the depth first and breadth first algorithms as starting points. Given the simple tree structure here, one of the main constraints is the accuracy of the state evaluator $V(p_\theta, S)$. In particular, it is likely easier for the language model to evaluate the likelihood of a particular mental state directly using a subtext instead of a thought generated by that subtext, because it will potentially match what it has seen in its training distribution better.

Hence, the final algorithm we implement can be thought of a breadth first search over the first layer, including (context+subtext) with a greedy rollout of the most promising state to produce a thought based on the subtext and then the final model response as in Equations 3.2 and 3.3.

The algorithmic implementation of this for this tree structure is included below.

---

**Algorithm 1** Greedy Rollout After Breadth-First Search

---

**Require:** Input $x$, LM $p_\theta$, subtexts $T = \{\tau_0^{(i)}\}_{i=1}^N$, thought generator $G()$, state evaluator $V()$.

1: $S_0' \leftarrow \{[x, \tau_0^{(i)}] \mid \tau_0^{(i)} \in T\}$
2: $V_0 \leftarrow V(p_\theta, S_0')$
3: $s_0^* \leftarrow \arg\max_{s \in S_0'} V_0(s)$
4: $s_1 \leftarrow G_1(p_\theta, s_0^*, 1)$
5: **return** $G_2(p_\theta, s_1, 1)$

---

## 3.3   Ranking and State Evaluation Metrics

For being able to characterise the ranking of the thoughts as well as the consistency of the voting process, there are various metrics we will use.

**Ranking thoughts**

**Average Reciprocal Rank**

To be able to rank the thoughts we use the Average Reciprocal Rank (ARR) which is very similar to the Mean Reciprocal Rank used to evaluate the effectiveness of ranking systems such as search engines and recommendation systems. The main difference is that ARR is calculated per thought instead of per-query. It is given as:

$$ARR_i = \frac{1}{N} \sum_{k=1}^{N} \frac{1}{r_{ik}} \tag{3.4}$$

Where $ARR_i$ is the Average Reciprocal Rank of thought $i$ and $r_{ik}$ is the rank of thought $i$ in trial $k$. This creates a score that varies from worst at $1/N$ to a best score of 1. While doing a sum of ranks would treat 1st vs 2nd vs 3rd as linear differences, MRR or calculating the number of times a thought is ranked first would only consider the rank of the 1st ranked thought. Hence, ARR provides a middle ground where being 1st is prioritized, but having a good ranking such as always ranking second will still have a good score.

**Measuring consistency in voting**

When prompting the LLM with the subtext thoughts in different orders, I observed that the voted thought depended on the order of thoughts in the prompt. Ideally, the selection mechanism would be perfectly consistent, however there must be some bias or uncertainty in the voting process that decreases this consistency. We can quantify the consistency of the voting process using two complementary measures:

**1. Maximum Agreement Score**

For each sample we define $V_{i,j} = \{v_{i,j}^{(1)}, v_{i,j}^{(2)}, \ldots, v_{i,j}^{(L)}\}$ as the set of votes received by thought $j$ in sample $i$ across $L$ shuffles, where: $v_{i,j}^{(\ell)} = \begin{cases} 1 & \text{if thought } j \text{ was selected in shuffle } \ell, \\ 0 & \text{otherwise.} \end{cases}$

Then the count matrix $C_{i,j}$ represents the total number of times thought $j$ was selected for sample $i$ across all shuffles:

$$C_{i,j} = \sum_{\ell=1}^{L} v_{i,j}^{(\ell)}$$

(Note that the number of shuffles can be recovered by calculating $L = \sum_{i=1}^{k} c_i.$). The maximum agreement score for sample $i$ is then:

$$\max_j C_{i,j}.$$

This score ranges from $T/L$ (minimal consistency) to $L$ (perfect consistency). Higher values indicate stronger robustness to ordering effects.

## 2. Evenness Index

The evenness index provides a normalized measure of vote distribution uniformity for sample $i$. It ranges from 0 to 1, where a value of 1 indicates the distribution is uniform and highly inconsistent, while a value of 0 represents perfect consistency. It is defined as:

$$E_i = \frac{H_i}{\log_2(L)} = \frac{-\sum_{i=1}^{k} p_i \log_2 p_i}{\log_2(L)} \tag{3.5}$$

Where $E_i$ is the evenness index for sample $i$, $H_i$ is the Shannon entropy of the vote distribution at sample $i$ and $p_j = \frac{C_{i,j}}{L}$ is the probability that thought $j$ was selected at sample $i$. The denominator $\log_2(L)$ normalizes the entropy by the maximum possible entropy over $L$ shuffles, providing a measure of how evenly the votes are distributed across the thoughts.

## 3.4  Majority Voting

Using the measure of consensus described in the previous section, we can create a quasi-mixture of experts voting approach by aggregating scores/choices over different order of subtexts and selecting the most popular vote.

Effectively, this creates a meta-state evaluator, that takes in the output of the state evaluator $V(p_\theta, S)(s)$ over $L$ shuffles. The voted subtext thought for sample $i$ using count matrix $C_{i,j}$ by seeing what the most popular thought is would be:

$$\arg\max_j C_{i,j}.$$

# Chapter 4

# Experimental Results

## 4.1 Experimental Setup

### 4.1.1 Dataset curation

The first step of the experimental setup is setting up the dataset. For a general dialogue dataset with a series of utterances between two speakers, the first step is to segment it into a collection of dialogues, where each dialogue $D = \{u_1, u_2, \ldots, u_N\}$ is represented as a sequence of $N$ utterances between two speakers. After each dialogue is created, we can extract a series of contexts $C = \{c_{w+1}, c_{w+2}, \ldots, c_N\}$ which are used to generate model responses $R = \{r_{w+1}, r_{w+2}, \ldots, r_N\}$ to compare to human reference responses $L = \{l_{w+1}, l_{w+2}, \ldots, l_N.$

To construct the sets $C$ and $L$, a sliding window of size $w$ is applied over each dialogue. For each turn $i$ from $w + 1$ to $N$, we collect the $w$ utterances immediately preceding $u_i$ as the context $c_i$, and define the corresponding label $l_i$ as the ground-truth utterance $u_i$. Since at least $w$ previous utterances are required to form a complete context, the model only generates responses for turns where $i > w$. This ensures that predictions are all conditioned on the same amount of prior context for fairness.

## 24/7 Conversations

To implement this setup for the "24/7 Conversations" dataset, a series of preprocessing steps are required. To give some context, the "24/7 Conversations" dataset contains 750 hrs of continuous electrocorticography (ECoG) data from epilepsy patients engaging in free daily life conversations alongside audio recordings using a microphone to capture all the speakers in the room. These recordings have already been preprocessed by members of the Hasson lab using an Automatic Speech Recognition system to generate a time-stamped transcript. For one of the patients, patient 798, we have this data segmented by sentence with speaker labels.

To prepare the dataset as a collection of dialogues $\mathcal{D}$, we first transform the transcript into a sequence of utterances. Multiple consecutive sentences spoken by the same speaker are merged into a single utterance. After constructing the collection of dialogues $\mathcal{D} = \{D_1, D_2, \ldots, D_J\}$, we apply steps which aim to filter out interactions that lack meaningful turn-taking and exchange or are not useful for modeling conversational dynamics.

The first filtering step targets one-sided conversations, which are identified based on the ratio of speaker turns. Specifically, we exclude any dialogue in which a single speaker contributes more than 75% of the utterances, as these are less likely to have meaningful two-way exchanges. This step is visualized in Figure 4.1. In this dataset, we observe several instances where one participant—typically a doctor—dominates the conversation, often explaining something while the other speaker listens passively. These are some of the interactions we specifically seek to exclude.

Finally, we apply a length-based filter to the reference responses. Many utterances in naturalistic dialogue are extremely short—often backchanneling cues such as "mhm," "yeah," or "okay"—which are challenging to model and offer limited value for evaluating response generation. On the other hand, extremely long utterances introduce their own modeling difficulties. Therefore, in our analysis, we restrict at-

33

**Histogram of percentage of sentences spoken
by the domininant speaker in each conversation**

Figure 4.1: A histogram of the percentage of sentences spoken by the dominant speaker in each conversation

tention to responses containing between 5 and 25 words. Then, for patient 798, the final dataset consists of 827 labels across 318 dialogues, the average number of turns per dialogue being 13.3. These statistics correspond to the case where the maximum number of context turns used in the analysis is $w = 15$. A smaller value of $w$ would yield more datapoints, as fewer turns would be excluded for lacking sufficient preceding context.

(The code for this filtering is included in Appendix A.3)

### 4.1.2 Large Language Model Configuration

When selecting a language model for this project, we considered various models with different trade-offs, such as ease of use, performance on specific benchmarks, and flexibility for future fine-tuning. After careful evaluation, we chose the Llama 3.x series of models [73], which offers several key advantages. Notably, at the time of development it was one of the most powerful models with publicly available weights, providing the flexibility for potential future fine-tuning. Additionally, Llama 3.x is

available in multiple sizes—1B, 7B, and 80B parameters—trained on similar datasets, providing both adaptability across different use cases and the ability to assess the impact of model size on performance.

In particular, we use the instruction-tuned versions of Llama 3, which are specifically fine-tuned to improve the quality of generations by understanding more complex requests and producing more relevant responses. Generally, models are fine-tuned for specific tasks such as question-answering, summarization, and coding. Without this fine-tuning, the model would only predict the next token based on statistical patterns in the training data, leading to less rich and context-aware dialogue. The instruction-tuned models also enable us to incorporate a "system prompt," which helps shape the style and content of the model's responses.

The hyperparameters, such as temperature, that control the sampling behavior for generating the next token, are detailed in Appendix A.5.

## 4.2    One step utterance generation

This section will describe the results of the single step utterance generation method described in Section 3.2.1. The system prompt for this generation is given in Appendix A.1. The code for generation is given in Appendix A.4.

### 4.2.1    Sentence generation improves with model size and context length-Up to a Limit

Figure 4.2 presents the average BERTScore across all samples, varying by model size and context length (measured as the number of preceding utterances). Broadly, we observe that BERTScore improves with both increasing model size and longer context. However, the gains from additional context diminish beyond approximately 10 preceding turns. This could suggest that speakers typically do not reference dialogue

more than 10 turns prior as often, or that excessive context may dilute the model's attention to relevant parts of the conversation.

With regard to model size, both the 8B and 70B parameter models outperform the 1B model as context increases, with the 70B model slightly edging out the 8B.

This score is computed by averaging the BERTScore across randomly shuffled candidate-reference label pairs, which approximates the baseline generation method used in the original BERTScore paper [69] (the primary difference being that the original paper compares 1M pairs).

Two baseline scores are included as dashed lines in Figure 4.2. The red line represents the average BERTScore when comparing each human label $l_a$ to another randomly selected label $l_b$ from the dataset. Computing the score over randomly shuffled candidate- reference pairs of labels in this way approximates the method used to create a baseline in the original BERTScore paper [69] (the main difference being that the original paper compares 1M pairs). The blue line shows the average BERTScore between human labels and random generations from an 8B model without context. These baselines offer reference points: for instance, we expect models to perform no worse than the random generation baseline of 0.638.

Interestingly, the 8B and 70B models only begin to surpass the pairwise random-label baseline after a context length of 5, while the 1B model never exceeds it.

To evaluate the effectiveness of in-context prompting methods in later sections, we adopt the BERTScore of the 8B model with 10 turns of context as our model baseline because this configuration yielded the highest observed score.

Figure 4.2: Average sentence similarity scores between model generated and human reference turns of conversation for differing amounts of context and different model sizes

## 4.3 Thought generation for theory of mind reasoning

This section will go over the results of the multi-step sentence generation based on the thought generator by evaluating all possible branches in the tree of thought. By characterizing the performance of each branch, we can characterize the problem of choosing the best thought in more detail. Overall, I find that the model response based on the best subtext outperforms the baseline, however randomly selecting a subtext to generate based off of will not improve the baseline, hence being able to find the best ranked thought is important.

### 4.3.1 The best generated responses based on subtext thoughts improve upon baseline

Figure 4.3 presents the distribution of scores for each ranked thought, with thoughts ordered from highest to lowest score within each sample. The box plot illustrates how the top-ranked, second-ranked, third-ranked, and subsequent thoughts perform across all samples. Notably, selecting the model response associated with the first or second highest-ranked subtext consistently outperforms the baseline. In contrast, the third-ranked thought performs comparably to the baseline on average, while lower-ranked thoughts generally underperform. These results highlight a clear opportunity to improve upon the baseline by selecting higher-ranked generations.



Figure 4.3: Box plot showing the BERTScores for all ranked ToM subtext generations (1st, 2nd, 3rd best, etc.) for each sample.

Figure 4.4: Box plot showing the BERTScore distribution for each different subtext over all samples

## 4.3.2 Selecting Responses Based on a Single Subtext Under-performs Baseline

While choosing a highly ranked generation can outperform the baseline, consistently selecting responses based on a single subtext tends to slightly underperform. This trend is evident in Figure 4.4, which presents the average BERTScore for model responses categorized by subtext. From a Theory of Mind perspective, this result is logical: applying the same inferred mental state across all interactions overlooks the nuances of individual interactions and lacks conversational specificity to the context. As a result, biasing the selection mechanism toward any single subtext would not necessarily lead to any improved performance.

## 4.3.3 Self-Reflective and Goal-Oriented Thoughts Rank Higher

We use the Average Reciprocal Rank (ARR), defined in Equation 3.4, to evaluate the relative ranking performance of different subtexts. As shown in Figure **??**, there is some variation across subtexts, though none deviate drastically from the expected

ARR for a uniform distribution of ranks, given by $\frac{1}{M} \sum_{k=1}^{M} \frac{1}{k} = 0.408$. Notably, subtexts that prompt the model to think about *goals* tend to achieve higher ARR scores, while those focused on *decision-making* or *information processing* perform slightly worse. Across all subtexts, prompting the LLM to consider the *upcoming* speaker's state of mind from their own perspective generally yields better rankings than prompting it to reflect on the *previous* speaker's state.



Figure 4.5: Bar chart showing the Average Reciprocal Rank for each theory of mind subtext. The possible range is between 0 and 1, with a perfect score of 1 being achieved by having a rank of 1 over all samples

### 4.3.4 Sampling multiple times from a varied distribution can still lead to improvement in BERTScore

Having shown that a thought generation approach based on a set of subtexts can improve text similarity scores, we next explore how perturbing this method affects performance.

Figure 4.6 summarises this analysis. To produce this graph, we perturb several key components involved in dialogue generation—the context $c$, the list of subtext

Figure 4.6: BERTScore outcomes for different generation strategies with various perturbation conditions. Dark blue bars represent the average best score per sample, light blue bars the worst and thick black bars the score from randomly sampling a subtext each sample. The yellow bar shows the average score with a simple label. The black dashed line is the model baseline.

thoughts $T$, and the intermediate generated thought $z$. The only major variable that remains unchanged is the prompt.

To perturb the context $c$, we use the Random Generation (RandomGen) condition. As in Section 4.2, this involves removing the context from the LLM input and instead sampling random text generations. We sample $M = 6$ generations to match the multi-sample setup used when conditioning on different subtexts.

For the Random Subtext condition, we replace the theory-of-mind-based subtexts with a set of unrelated, casual thoughts that are not clearly tied to mental states or intentions. The set includes: "Hm. Where was I going?" "Mmm. I could have toast later." "I need to remember to call my mom later." "What's the weather like today?" "I need to remember to pick up groceries later." "Do fish ever get thirsty?"

In the Random Thought condition, we retain the original theory-of-mind subtexts, but instead of generating thoughts conditioned on the subtext and context, we substitute the intermediate thought $z$ with a completely random text sample.

41

We also include a simple baseline label—"Yeah, that's what I see too."—as a final comparison. Since this case does not involve multiple generations, it is plotted as a single yellow bar in Figure 4.6.

The figure shows the following: the light blue bars represent the average worst score across samples for each condition, while the dark blue bars show the average best score. The thick black horizontal lines indicates the expected score if a sample were selected at random. Lastly, the dashed black line represents the baseline model response from Section 4.2.

First, note that sampling multiple times from a distribution with some variance and choosing the best outcome will, by definition, outperform the mean of the distribution. This section evaluates that intuition in practice by modifying the generation process to inject variety.

The results confirm that all generation strategies exhibit variance in output quality, meaning that sampling multiple times and choosing the best result consistently improves BERTScore. Interestingly, both the theory-of-mind and random subtext conditions yield similar distributions, with the random subtext condition exhibiting slightly more variance—resulting in a higher best-case score.

Surprisingly, the Random Thought condition outperforms both subtext-based strategies. One plausible explanation emerges from examining the performance of the Simple Label condition, which closely matches the model baseline despite its simplicity. This suggests that basic or generic statements such as "Yeah, that's what I see too" can score well in similarity metrics. By introducing randomness in thought generation, the model may "get distracted" and fall back to more simple, generic responses. These responses would cluster near the average score, thereby lowering the worst-case performance but potentially boosting the maximum due to the occasional strong match.

Given these findings, for this strategy to have any use, it is crucial to be able to

Figure 4.7: Bar chart showing the Average Reciprocal Rank for each random idea subtext. The possible range is between 0 and 1, with a perfect score of 1 being achieved by having a rank of 1 over all samples

interpret performance changes in relation to underlying phenomena, rather than the simple effect of sampling multiple times from a noisy distribution. In this context, further analysis is needed to determine whether improvements in generation quality can be meaningfully linked to real-world mental states.

### 4.3.5 Using random ideas in subtext does not have clear trends in scores

Given our perturbation analysis, we can examine the random subtext and theory of mind subtext conditions in more depth for interpretability.

Looking at Figure 4.7 and comparing it to 4.5, we can see that the random subtext condition has no clear trend in performance based on the random subtext ideas. There is also more deviation from the uniform ranking ARR value.

Additionally, when computing the average scores across all samples in Figure 4.8, the results are comparable to those for the theory of mind subtexts shown in

Figure 4.8: Box plot showing the BERTScore distribution for each different subtext over all samples

Figure 4.4. One likely reason for this is that the large amount of conversational context—specifically, the inclusion of 10 preceding turns—allowed the model to produce a sensible reply, even when the subtext and resulting thought were largely ignored. When inspecting the response often they were mainly related to the context and not directly to the subtext thought.

## 4.4 State Evaluator

This section will review all the selection mechanisms implemented. In summary, no selection mechanism was able to select subtexts based on the conversational context to improve the generation BERTScore.

### 4.4.1 Choosing the best ranked thought improves BERTScore

To be able to contextualize our LLM-based voting mechanisms, the results from the random and perfect state evaluators are included as a reference for comparison here. A scatter plot of the scores from both state evaluators against the baseline

Figure 4.9: Box plot showing the BERTScores for all ranked random idea subtexts (1st, 2nd, 3rd best, etc.) for each sample.

generation score is included in Figure 4.10. We can see that generating a model response based on the perfect state evaluator which chooses the subtext resulting in the best outcome has a BERTScore which is almost always better than baseline single-step generation. However, when selecting a generation based on a random subtext there is no improvement in performance. This is consistent with the data in Figure 4.3 which shows that either the first or the second ranked subtext must be used on average to improve on the baseline.

### 4.4.2 Both LLM-based state evaluators are not effective in selecting the best subtext

Figure 4.11 depicts the scatter plot of the scores from both LLM-based state evaluators against the baseline generation score. Visually, it seems that the voting mechanism has no noticeable change on the distribution of scores, with the distribution appearing much more similar to the random state evaluator compared to the perfect state evaluator in Figure 4.10.

Figure 4.10: Scatter plots of the BERTScore based on model response from the best ranked thought (left) and a randomly selected thought (right) against baseline response

The results in Figure 4.12 confirm these observations. All selection mechanisms apart from the perfect state evaluator perform marginally worse than the baseline. It can be seen that the perfect state evaluator also has a slightly less spread distribution.



Figure 4.11: Scatter plots of the BERTScore based on model response from the voting over all possible subtexts (left) or scoring each subtext individually (right) against the baseline response.

Figure 4.12: Violin Plot of the BERTScore for different State Evaluators. The dark blue line represents the mean of the distribution while the black lines represent the datapoints 1.5 times the interquartile range (IQR) from the first and third quartile.

## 4.5 Evaluation over multiple shuffles

The following results are obtained by applying the voting scoring method and voting choice method 10 times each with the order of subtexts to choose randomized each time.

### 4.5.1 Evaluation by Scoring and by Voting are biased towards certain numbers

One reason that the state evaluation mechanism is not working well could be because of bias, which can be an issue when using LLMs as evaluators [74][75]. An obvious way this appears is a preference for certain numbers.

In Figure 4.13 we can see that for example the LLM perfers not to choose the first item in the list regardless of which subtext it represents. Similarly, in Figure 4.14 we

can see that the scoring mechanism has a strong preference for the digits 2,4 and 8. With a couple other digits having frequencies one or multiple orders of magnitude less, since the graph is plotted using a logarithmic scale.



Figure 4.13: The frequency each numerical index in the list is chosen when the state evaluator votes over states



Figure 4.14: The frequency specific digits appear when the state evaluator scores each option. Calculated over 10 shuffles.

### 4.5.2 Evaluation by Scoring is more consistent than by Choice

We can also quantify the consistency of the voting mechanism by using the metrics defined in Section 3.3, evenness index, and maximum agreement score.

Figures 4.15 and 4.16 show the dstribution of these metrics over all the samples. It is clear that both voting mechanisms have high levels of bias; for example, for most samples, the maximum agreement on the best subtext over the different shuffles did not each a majority of more than 4. However, these plots do show that the evaluation method by scoring independently is more consistent than voting over all the possible options, since the average entropy is closer to 0 and the number of maximum agreements also increases.



Figure 4.15: The continuous density plot of evenness index of all samples for different shuffles. An evenness index of 0 would indicate perfect consistency, while an index of 1 would indicate a perfectly random/uniform evaluation process.

## 4.6  Majority Voting

Since we implemented a shuffling mechanism, we can also create a meta-state evaluator by selecting the most popular subtext after a certain amount of shuffles.

Figure 4.16: The histogram of maximum number of agreements on the best subtext over different shuffles for all samples. A score of 10 would indicate perfect consistency while a uniform/random evaluator would score 1.67 on average.

### 4.6.1 Majority voting does not improve performance

In the scatter plots in Figures 4.17 and 4.18 we can see that sampling the state evaluator multiple times does not have a noticeable improvement on the distribution of scores.

### 4.6.2 High consistency on evenness index leads to lower variation of BERTScore of output

Figure 4.19 shows the distribution of scores for the majority vote meta-state evaluators as well as the distributions of the samples evaluated with high consistency (which are also highlighted in the scatter plots in Figures 4.18 and 4.17). While neither the majority vote meta-State Evaluators or the high consistency examples have a greater similarity to the human reference compared to the baseline, the evenness index metric of the scores or voting is indicative of the variance of the final BERTScore distribution. That is, the distribution of high consistency examples plotted as judged by evenness

Figure 4.17: A scatter plot of the BERTScore of the majority vote meta-state evaluator voting by voting and by scoring. General points are in grey while coloured points refer to high consistency examples with number of agreements ¿3.



Figure 4.18: A scatter plot of the BERTScore of the majority vote meta-state evaluator voting by voting and by scoring. General points are in grey while coloured points refer to high consistency examples with evenness index0.55.

index has a lower variation of BERTScore compared to all other conditions.

Figure 4.19: Violin Plot of the BERTScore for the majority vote meta-State Evaluators based on scoring and voting. Also plotted are the distributions of high consistency examples highlighted in the scatter plots in Figures 4.18 and 4.17

## 4.7 Overall findings

In this project, we further demonstrated that large language models are capable of generating sensible and coherent dialogue. Our experiments showed that both model size and context length positively correlate with similarity to human-generated references.

We introduced a thought generator based on a theory-of-mind taxonomy, which produced diverse response sets—some outperforming and others underperforming the baseline. This motivated the implementation of a selector mechanism designed to choose the most promising subtext to guide response generation. To better understand the generator's behavior, we conducted a perturbation analysis, exploring how changes to different components influenced the distribution of model outputs. While BERTScore appeared to be a reasonable metric for assessing textual similarity, further investigation is needed to determine whether it captures sufficient semantic nuance,

and whether the generated thoughts and responses genuinely align with the theory-of-mind taxonomy provided.

Finally, we developed three state evaluation methods to select the best subtext for response generation. However, none of these approaches outperformed the baseline model that generates responses in a single step. Interestingly, we observed inconsistency and bias in the voting and scoring evaluators, though the consistency of an evaluator across multiple data shuffles—measured using the evenness index—appeared to reflect the spread in BERTScores of the model outputs. This suggests that LLMs might have some capacity to internally assess or express confidence in their generated responses.

# Chapter 5

# Discussion

While the ultimate goal of automatically improving dialogue generation using a theory of mind-based in-context learning approach was not fully achieved, this project still provides valuable insights—both for future improvements and for our broader understanding of conversational modeling.

To better understand the performance and shortcomings of the system, we can break down the implementation into its key components. Given the complexity and number of moving parts, this component-wise reflection is essential for identifying areas for refinement.

## 5.1   Future work

### 5.1.1   Language Model Capabilities

First, the capabilities of the language model itself, denoted $p_\theta$, may have been a limiting factor. It is possible that the model's social reasoning abilities are insufficient for accurate inference in this domain. This challenge is compounded by the inherent noisiness and ambiguity of social reasoning tasks. If the LLM has not been trained extensively on such data—likely, given the difficulty of curating high-quality theory-

of-mind datasets—it may lack the inductive biases required for these tasks. Although the weights of the LLaMA 3.x models are publicly available, their training data is not, making it difficult to verify this hypothesis. Potential next steps include fine-tuning the model on curated social reasoning datasets, such as SocialIQA [76], or designing a new benchmark of simpler, well-scaffolded theory-of-mind tasks to evaluate and adapt the model progressively.

### 5.1.2 Context Representation

Second, the context $c$ provided to the model may not have been rich or structured enough. While it included recent conversational history, many other types of information are relevant for generating socially coherent dialogue: background world knowledge, participant personalities, shared memories, and conversational goals. A more structured or explicitly abstracted representation of context may better support theory-of-mind reasoning, aligning with thoughts in other work [26].

### 5.1.3 Evaluation Metrics.

Third, the evaluation metric—BERTScore—may lack the granularity necessary to capture nuanced social reasoning or alignment with human intent. Although BERTScore correlates with textual similarity, it does not assess propositional content and does not have extremely high correlations with human judgment in other tasks. Alternative metrics could include: using a pipeline to convert responses and labels to propositional form, to remove any reliance on surface represetnation of text before using a text similarity like BERTScore. Additionally, encoding accuracy from the sentence embedding of the generated text could be used as a score of how good a generation is. Then BERTScore as a metric could also be validated by seeing how well it correlated with encoding accuracy. Perplexity could also be used as a metric to capture likelihood of a response given context. Finally, human evaluations canbe used to validate

the use of these automatic scores.

### 5.1.4   Thought Representation and Tree of Thought.

Another area for exploration is the structure of the thought process itself. It remains unclear whether the intermediate thoughts generated by the LLM accurately reflect the model's internal reasoning [77, 78, 79, 80]. Doing more perturbations of thoughts could help show this.

An especially compelling dataset to apply this method to is ESConv [81], which features emotional support conversations between a "speaker" describing their problems and a crowdworker "listener" providing support. Each response from the listener is annotated with a communicative strategy, selected from a set that includes Questions, Self-disclosure, Affirmation and Reassurance, Providing Suggestions, Other, Reflection of Feelings, Information, and Restatement or Paraphrasing. This makes ESConv a rich resource for evaluating whether social reasoning capabilities—such as those enabled by a tree-of-thought approach—can be leveraged to predict the underlying social strategies used by the listener.

Furthermore, one potential direction is to simplify the thought process, for instance by reducing the tree depth from two to one, and conditioning the final response more directly on the intermediate reasoning steps. Prompt engineering—particularly for generating or interpreting theory-of-mind subtext—could also be refined.

### 5.1.5   Filtering to Target Turns that Require Reasoning

An important open question in conversational analysis is determining when participants engage in more automatic, reflexive responses versus when they employ deliberate reasoning. Our current predictions may capture turns that elicit automatic replies, which are inherently more challenging to model due to their low cognitive load and variability. Incorporating more sophisticated filtering methods could help isolate

conversational turns that genuinely require thought or reasoning. For instance, the high BERTScore of the generic label "Yeah I see that too" in Section 4.3.4 suggests that such low-effort responses may be disproportionately influencing our evaluation metrics, indicating a need for refined filtering criteria.

## 5.2   Limitations

First, the project assumes that conversations can be segmented into rigid, alternating turns. While this simplifies implementation, it does not reflect the fluid nature of real-world dialogue, which often includes interruptions, overlaps, or simultaneous speech. Moreover, it overlooks how interlocutors build shared context over time. The common ground theory [82] describes how shared knowledge accumulates and simplifies communication—for example, initial explanations may be detailed, but future references become more abbreviated. While this project began with segmenting dialogues into discrete two-person turns, future work could explore continuous, overlapping, or multi-party dialogue representations.

Second, the model is fundamentally unimodal—it operates on language alone. However, in human conversation, responses are influenced not only by the preceding words but also by a range of visual, situational, and embodied cues. Listeners naturally integrate information such as visual scene understanding [83], the goals and perspectives of the speaker [84], the physical context or affordances of objects in the environment [85], and paralinguistic signals such as facial expressions or body language [86, 87]. While some of these aspects may be indirectly encoded in language, much of this contextual richness remains implicit and inaccessible to language-only models. As such, future modeling efforts could benefit significantly from incorporating multimodal inputs or grounding the model in real-world context.

# Chapter 6

# Related Explorations-Neural Encoding using Sentence and Word Embeddings

This chapter presents the second project of the thesis: exploring how to perform sentence-level encoding of intracranial neural recordings. Unlike fMRI, intracranial recordings offer much higher temporal resolution—often far exceeding the duration of a single sentence. This mismatch in timescales makes sentence-level alignment particularly challenging.

Previous work on this dataset has primarily focused on encoding neural activity using token-level embeddings. However, incorporating sentence embeddings could provide a window into higher-order cognitive processes by capturing more abstract semantic representations. In this section, we investigate the feasibility of combining word and sentence embeddings to model neural activity using a banded ridge regression framework.

## 6.1 Data collection and pre-processing

The original dataset and pre-processing method is provided by Zada et al. (2025) [88]. To summarise, a 30 minute audio story (podcast) was played to nine participants while undergoing electrocorticographic monitoring for epilepsy. A simple preprocessing pipeline was implemented to extract the high-gamma band power per electrode. Alongside the neural recording, the story was also manually transcribed with each word timestamped at high temporal resolution.

The stimulus presented to participants was a segment from the podcast "This American Life" entitled "So a Monkey and a Horse Walk Into a Bar: Act One, Monkey in the Middle" released on November 10, 2017. The original audio and transcript are freely available online (https://www.thisamericanlife.org/631/transcript). We manually transcribed the story and timestamped words at high temporal resolution.

Raw electrode data underwent the following preprocessing pipeline, which is quoted from the Zada et al [88]. The steps include "removing bad electrodes, downsampling, despiking and interpolating high amplitude spikes, common average re-referencing, and notch filtering. First, we visualized the power spectrum density of each electrode per subject. From this, we were able to annotate unusual electrodes that did not conform to the expected 1/f pattern, had a consistent oscillatory pattern, or showed other unusual artifacts. We found 31 such electrodes and marked them as "bad" (identified in the accompanying metadata). The source of these artifacts may be due to several factors, including excessive noise, epileptic activity, no noise, or poor contact. For data acquired with a sampling rate greater than 512 Hz, we downsampled to 512 Hz to match the sampling rate across subjects. We then applied a despiking and interpolation procedure to remove time points that exceeded four quartiles above the median of the signal and refill it using pchip interpolation. For re-referencing, we subtracted the mean signal across all electrodes per subject from each of their

individual electrode time series. Finally, we used notch filters at 60, 120, 180, and 240 Hz to remove power line noise. This pipeline produces a "cleaned" version of the raw signal."

## 6.2   Method: Ridge regression

Here, I will describe the original method in the paper for ridge regression to decode neural activity from token-level embeddings.

Let $X \in \mathbb{R}^{n \times p}$ be a feature matrix with $n$ samples and $p$ features, $y \in \mathbb{R}^n$ a target vector, and $\alpha > 0$ a fixed regularisation hyperparameter. Ridge regression defines the weight vector $b^* \in \mathbb{R}^p$ as:

$$b^* = arg \min_b ||Xb - y||_2^2 + \alpha ||b||_2^2. \tag{6.1}$$

The simplest way of casting this problem would be to have $n$ be the number of words and $p$ be the number of electrodes, so that the target matrix $Y$ has length of number of electrodes. Each entry in $X$ would contain the word embedding for the word at each new word onset, while $Y$ would contain the voltage signals recorded from the electrodes placed on the brain's surface in Volts. ECoG electrodes measure the aggregate electrical activity of thousands to millions of neurons and the recorded will mostly reflect local field potentials (LFPs), which arise from the synchronous activity of many neurons. Additionally, if the word consists of more than one token, we average all of them to get the embedding per word.

However, this method is unideal because the electrode data (512Hz) is much higher frequency than the spoken words ( 2Hz). To solve this, we epoch the electrode data around the onset of each word, so that $X$ now has shape (number of words, number of ECoG electrodes x number of lags) and $Y$ has shape (number of ECoG electrodes x number of lags). The epoch period chosen is from -2 seconds to +2 seconds relative to

word onset. This is equivalent to fitting a separate regression model for each electrode and each epoch.

Before the model can be fit, there are there are two additional elements in the pipeline we define in Scikit-learn [89]. First, $Y$ and the regressors in $X$ were mean-centered using StandardScaler. Second, because the design matrix was wider than it is long, we used the kernel method to solve the ridge regression in its dual form [90]. Specifically, we used a linear kernel for each feature space separately before fitting the model.

The penalty term $b$ is learned using an inner cross-validation setup within the training set using the himalaya python library [91]. Then an outer cross-validation loop is used to evaluate the encoding model with $k = 2$. We evaluated the performance of each encoding model on the held-out fold by correlating the model-predicted ECoG signal $Y_{\text{preds}}$ with the actual ECoG signal $Y_{\text{test}}$ in the test fold. Finally, we averaged the two correlations for each fold to obtain one correlation value denoting the encoding performance for each feature space, electrode, and lag.

## 6.3   Method: Banded Ridge regression

We use a banded ridge regression model to regress to get the ECoG signal in $Y$ with two different feature spaces with their own regularisation parameters as follows.

$$b^* = \arg\min_b \left\| \sum_{i=1}^{m} X_i b_i - y \right\|_2^2 + \sum_{i=1}^{m} \alpha_i \|b_i\|_2^2 \tag{6.2}$$

The variables are similar to Equation 6.1, where $X_i$ represents the feature matrix corresponding to the $i$-th group, $b_i$ denotes the coefficient vector for the $i$-th group, $\alpha_i$ is the regularization parameter specific to the $i$-th group and $m$ is the total number of feature groups.

In practice for constructing X, the shape of $X_1$ and $X_2$ are both (number of words,

number of ECoG electrodes x number of lags). Because a sentence contains multiple words, the sentence embedding will be the same for each word in that sentence.

The embeddings used are:

- **Static word embeddings:**

    – Word2Vec embeddings (300-dimensional)

- **Contextual word embeddings:**

    – GPT-2 (13 layers of embeddings, each word embedding is 768-dimensional)

- **Sentence embeddings:**

    – SentenceTransformer `intfloat/e5-mistral-7b-instruct` (4096-dimensional)

## 6.4 Results

### 6.4.1 Encoding with context and sentence embeddings outperforms static embeddings, but not contextual embeddings

Figure 6.1 depicts the encoding performance of a banded ridge regression with joint contextual and static embeddings compared to a ridge regression with either only static or contextual embeddings. It is clear that the joint embedding model improves on the static embeddings however compared to the contextual embeddings the performance appears very similar.
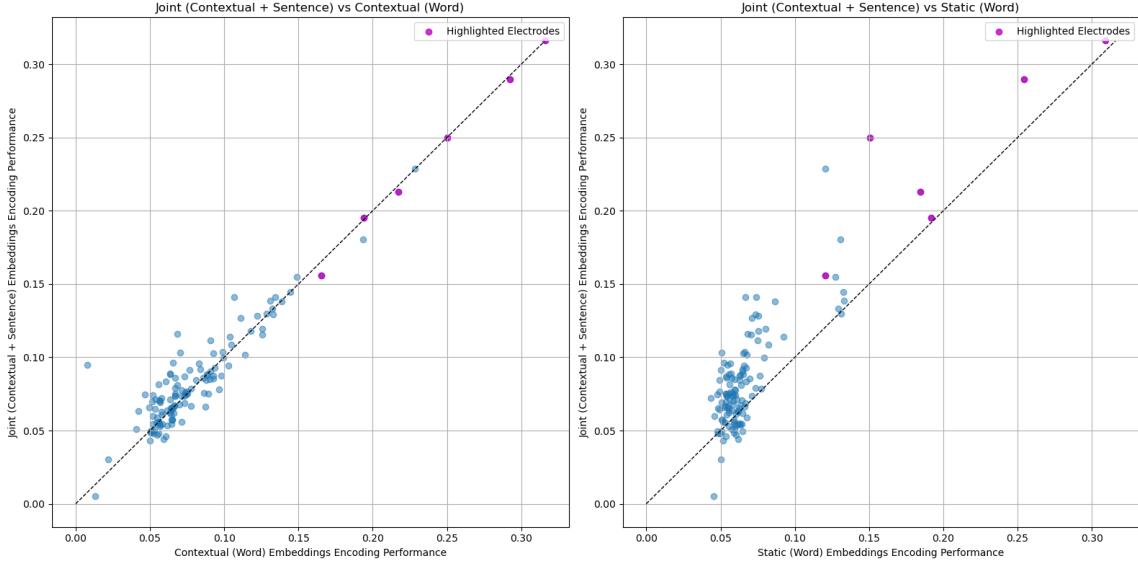
Figure 6.1: Scatter plots of encoding performance (correlation) for joint (sentence+contextual word) embeddings against static and contextual embeddings

## 6.4.2 Sentence Embeddings take up variance in Joint Context+Sentence encoding models

While using banded ridge regression with contextual and sentence embeddings does not improve upon the performance of ridge regression using only contextual embeddings, examining how the variance is split in the joint model can still be ueful.

Figures 6.5 and 6.6 and show that variance is split to a large extent between contextual and sentence embeddings. Figure 6.5 shows the encoding performance with either feature space separately, which shows that by themselves the overall encoding performance decreases. However, there are some electrodes where the sentence embeddings have a higher relative encoding performance and others where contextual embeddings have a much higher encoding performance on the right of the graph.

Plotting the encoding performance on the brain in Figure 6.6 we can see that the electrodes where the contextual embeddings do best are around XYZ brain areas. However, the correlation for sentence embeddings throughout the brain is generally more diffuse.
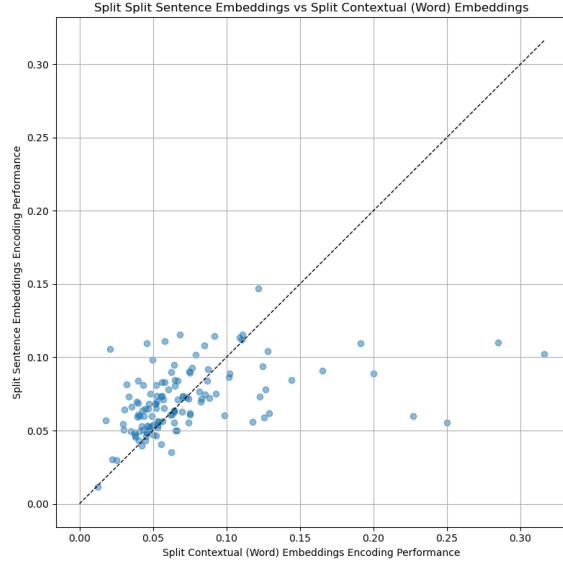
63

Figure 6.2: The encoding performance using either the contextual embedding or sentence embedding feature space plotted against each other



Figure 6.3: Correlation scores per electrode mapped onto the brain for joint (contextual word + sentence), contextual word and sentence embeddings

### 6.4.3 Encoding with static and sentence embeddings outperforms static embeddings, but not contextual embeddings"

In contrast to Figure 6.1 we compare the encoding performance of joint (static + sentence) embeddings compared to contextual and static embeddings in Figure 6.4

as opposed to joint (contextual+sentence) embeddings.

In Figure 6.4 we can see that the joint (static+sentence) model performs worse than model based on just contextual embeddings, however this joint model improves performance on models trained with just static embeddings. Most of this increase in performance happens in a range between a correlation of 0.05 to 0.09 for the static model, which increases to a range of 0.05 to 0.13. The performance for high correlation electrodes tends to be similar for both static and sentence embeddings.

6.4.



Figure 6.4: Scatter plots of encoding performance (correlation) for joint (sentence+static word) embeddings against static and contextual embeddings

### 6.4.4 Joint Sentence+static embeddings enable encoding for best electrodes over a greater period of time compared to only static embeddings

Another way of plotting our results is by collapsing the electrode dimension by selecting a subset and averaging, rather than collapsing the time dimension by choosing the maximum correlation over all time.

Figure 6.5: The encoding performance using either the contextual embedding or sentence embedding feature space plotted against each other
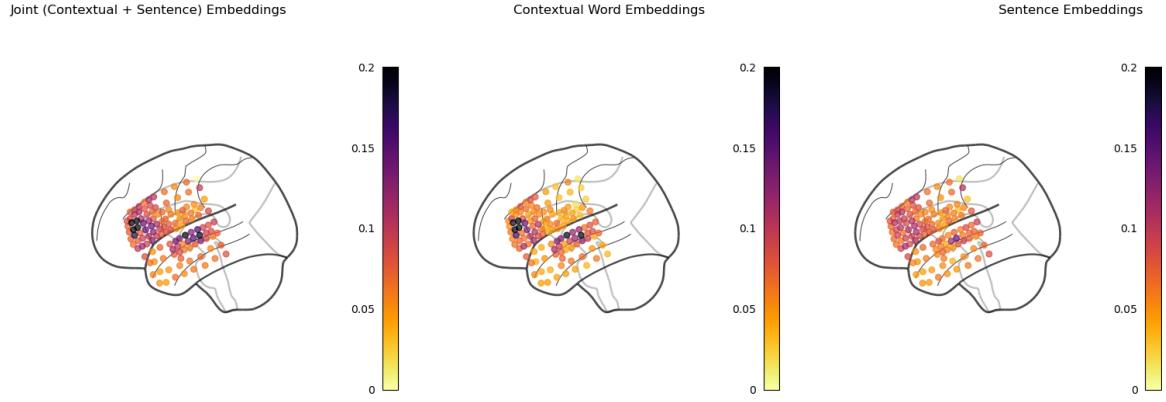


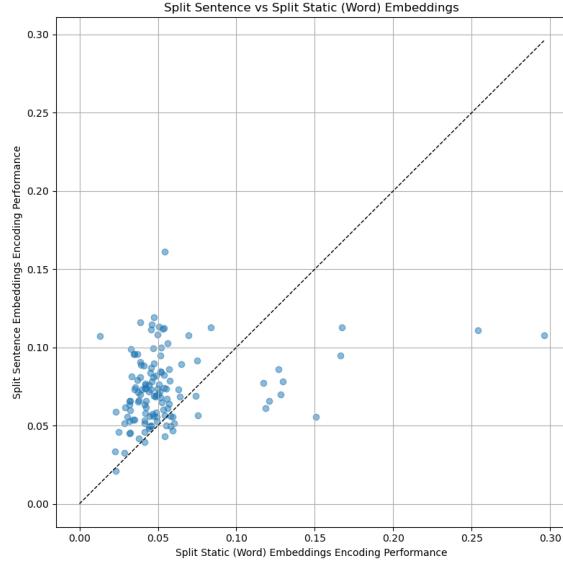Figure 6.6: Correlation scores per electrode mapped onto the brain for joint (static word + sentence), static word and sentence embeddings

Figures 6.7 and 6.8 show the correlation over time for a subset of the best-performing electrodes. From Figure 6.7 is clear that encoding using sentence embedding in addition to static embeddings increases the encoding performance over a greater period of time. In Figure 6.8 we can see this effect as well, but it is less noticeable because the contextual embeddings already perform better over a greater period of time compared to the static embeddings. Nevertheless, by looking at the split embedding spaces we can see that the word embeddings generally peak a few

100ms after word onset, but the encoding accuracy using sentence embeddings is more spread out over time.

To be able to relate the size of the outer blue envelopes in Figures 6.7 and 6.8 we also plot them overlapping alongside the results from contextual embeddings in Figure 6.9(a). We can see that the outer envelopes of contextual, contextual+sentence and static+sentence have similar profiles. While in Figure 6.9(b) we can see that encoding using either static or split-static embeddings has a more narrow profile.

Encoding Performance Over Time For Best Electrodes



Figure 6.7: The correlation of model and actual neural signal plotted over time and averaged over best performing electrodes. Black lines show the encoding done using a split feature space after banded ridge regression for static embeddings (left) and sentence embeddings (right). The blue line represents the joint embedding performance. Electrodes where the joint (static+sentence) correlation had a performance of more than 0.15 were selected.

## 6.4.5 Overall findings

In this project, we successfully implemented a model for sentence encoding. We found that static and sentence embeddings when combined can get encoding performance that approaches that of contextual embeddings, with the added flexibility of being able to change the content of the sentence embedding easily compared to contextual word embeddings.
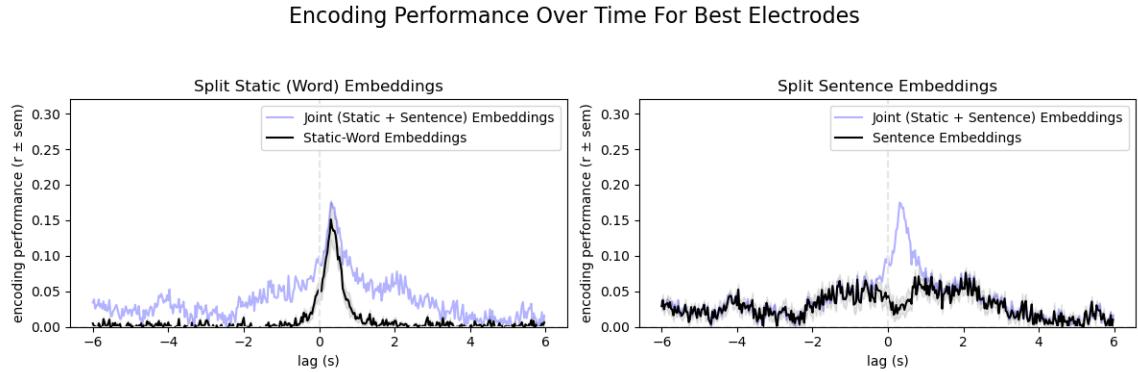
Figure 6.8: The correlation of model and actual neural signal plotted over time and averaged over best performing electrodes. Black lines show the encoding done using a split feature space after banded ridge regression for contextual embeddings (left) and sentence embeddings (right). The blue line represents the joint embedding performance. Electrodes where the joint (static+sentence) correlation had a performance of more than 0.15 were selected
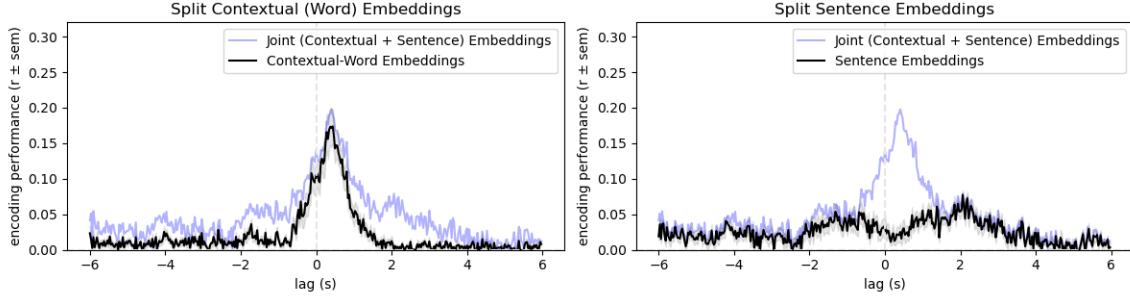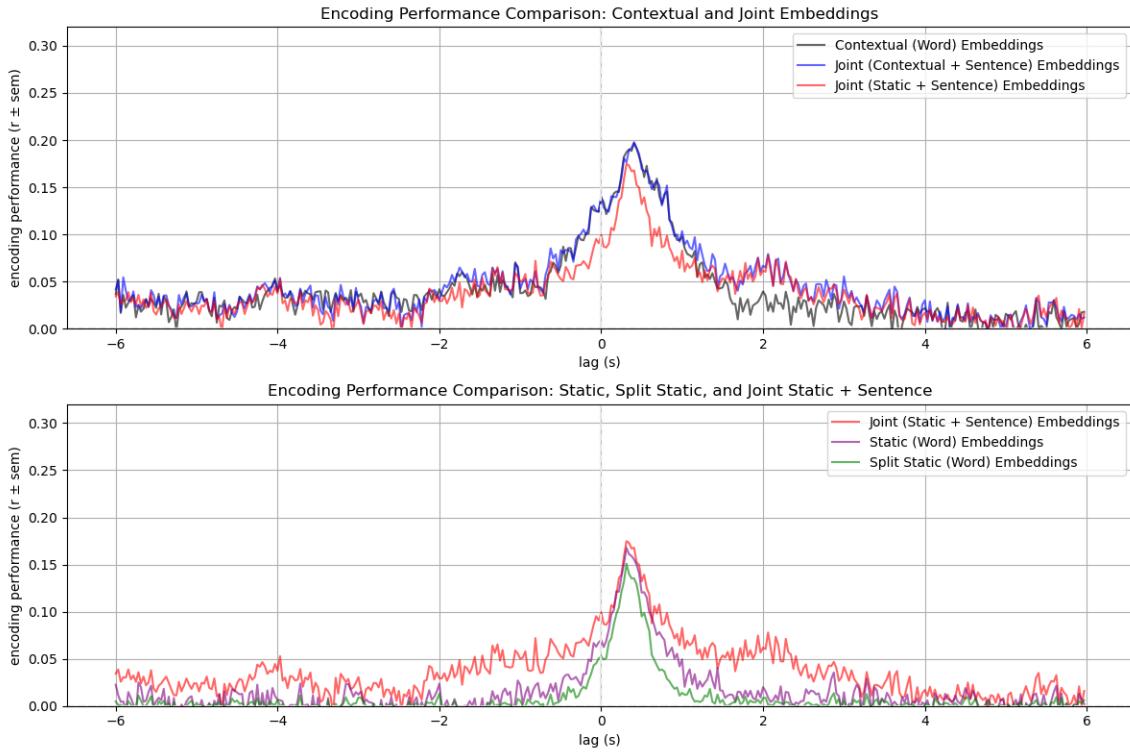


Figure 6.9: The correlation of model and actual neural signal plotted over time and averaged over best performing electrodes. Embeddings used in the model in order of envelop size are contextual, joint (contextual+sentence), joint (static+sentence) (plotted twice), static, and split static.

# Chapter 7

# Future Directions

Now that the two projects in this thesis have been established, it is worth exploring how they might be integrated to offer insight into the neural basis of theory of mind.

One direct way to bridge the conversational modeling project with neural data is by using the voted subtexts as labels for each utterance. These labels could serve as supervision signals in a neural classifier aimed at predicting the mental state most likely associated with each utterance. A crucial aspect of this task would involve distinguishing between moments when a person is engaged in mental state reasoning versus more routine conversational behavior—an open and intriguing question.

If the dialogue generation task proves challenging, an interesting alternative would be to adapt ideas from masked sentence modeling. Rather than predicting the final utterance in a conversation, the model could instead predict a missing utterance from the middle. This would allow the model to leverage both past and future context, potentially resulting better ability to reason over theory of mind.

In terms of evaluation, rather than relying solely on BERTScore, one could assess how strongly generated utterances drive activity in neural encoding models. After training on a subset of ground-truth data, such an approach could help identify "supernormal stimuli"—stimuli that elicit exaggerated responses from a given neural

system, as described by Barrett (2010) [92].

This project also lends itself naturally to an in-context learning framework, and its application need not be limited to social reasoning alone. Insights from psycholinguistics could enrich the modeling, particularly theories related to how people construct situation models through bridging inferences during discourse. Some accounts propose that literal and non-literal meanings are computed simultaneously within a parallel distributed framework—phenomena that could be explored by prompting large language models and examining their alignment with brain activity.

Another compelling area for exploration is inner speech. Inner speech serves diverse functions such as deliberation, planning, clarification, problem-solving, and even self-regulation or self-motivation (e.g., "You can do this," "Don't do that") [93]. Modeling these internal dialogues might provide a novel pathway for connecting language-based AI models with neural data, offering a richer picture of cognition and self-directed thought.

# Chapter 8

# Conclusion

This thesis explored the extent to which Large Language Models can perform socially informed reasoning during dialogue, focusing on the inference of subtextual thoughts—unspoken beliefs, desires, and intentions—through a theory-of-mind-inspired approach. By applying the Tree of Thought framework to naturalistic dialogue generation, we demonstrated that LLMs are capable of generating plausible, socially coherent responses. However, our results also revealed that existing evaluation metrics and selection mechanisms struggle to consistently identify the most appropriate subtext or response, highlighting limitations in both model expressivity and current assessment tools.

In the second part of this thesis, we extended LLM applications to neural encoding, showing how both word- and sentence-level embeddings can be used to model brain activity during language comprehension. This builds on prior work in the Hasson lab and offers new directions for integrating computational representations of language with neuroscience.

Across both projects, we found that while LLMs exhibit emerging capacities for modeling literal and inferred content, their reasoning processes remain opaque and prone to inconsistency. Evaluator bias, limited context representation, and inade-

quate scoring methods present ongoing challenges. Nonetheless, the promising signals observed in perturbation and consistency analyses suggest that LLMs may possess latent capacities for internal assessment and social inference that are yet to be fully unlocked.

Future progress will likely depend on more structured representations of context, improved datasets for training and evaluation, and new methods for interpreting model reasoning. Ultimately, this work contributes to the growing effort to bridge the gap between surface-level language modeling and deeper social cognition in AI systems, with implications for both computational neuroscience and the development of more human-aligned dialogue agents.

# Appendix A

# Code Implementation

## A.1   System prompt used

You will be given an excerpt of a conversation that includes the next speaker. It will be formatted as Speaker 1: message, Speaker 2: message. Reply as the next speaker in the conversation only. Format your reply as Speaker 1: message. Do not say you are a chatbot.

## A.2   Prompt Generation

```
1  {
2
3  import json
4  import random
5
6  #final JSON file format:
7  #prompts_subtext: a list of length s (number of subtexts) containing
       lists of length n (number of steps)
8  #subtexts: a list of length s containing the subtexts
```

```python
 9  #prompts_subtext_voting: a list of length 3 containing the prompts
        for the subtext generation (list of length s), voting (single
        string), and next speaker sentence generation (single string)

10

11

12

13  #numer of subtexts=s, number of steps=n

14

15  def exportJSONFile(data,fileName):
16      json_object = json.dumps(data)

17

18      with open(fileName, "w") as outfile:
19          outfile.write(json_object)

20

21

22  subtextPrompt="You will be given an excerpt of a conversation
        between two speakers. Format your output as {Output:}. Your task
        is to generate the subtext thought that the next speaker is
        having, based on the following idea:"

23

24  #length s
25  subtextIdeaList=["Do I agree with this?","What do they want from me?
        ","How are they feeling right now?","What are they avoiding
        mentioning?","Whats their body language saying?","Should I
        clarify my earlier point?","What do they really mean by this?"]
26  #subtextIdeaList2=["Why are they telling me this now?", "What
        triggered this topic in the conversation?","How does this relate
        to what they said earlier?","Does this connect to something
        theyve shared about themselves before?","How did they
        accomplish this, or how are they planning to?","What resources or
        skills did they use to solve this problem?","What are the
        implications of this for me or them?","What are they not saying
```

```
      that might be important?","What are they trying to achieve by
         saying this?","How does this align with their long-term goals or
         ambitions?"]
27 #subtextIdeaList3=["What are they about to say or ask?","Are they
         feeling frustrated, happy, or anxious?","Is this behavior typical
          for them, or is it out of character?","Are they hinting at a
         deeper issue without explicitly saying it?","Are they trying to
         persuade me?","Their background might explain why they see things
          this way."]
28 #subtextIdeaList4=["How do I feel about this?", "How does the other
         person feel about this?", "What do I want to happen next?", "What
          do they want to happen next?", "What do I need to know to make a
          decision?", "What do they need to know to make a decision?"]

29

30

31 randomSubtextList=["Hm. Where was I going?","Mmm. I could have toast
          later.","I need to remember to call my mom later.","Whats the
          weather like today?","I need to remember to pick up groceries
         later.","Do fish ever get thirsty?"]

32

33

34 subtextIdeaList_only4=["What do I want to happen next?", "What do
         they want to happen next?", "What do I need to know to make a
         decision?", "What do they need to know to make a decision?"]

35

36 prompt1=[]
37 prompt1Random=[]
38 selectedSubtexts=subtextIdeaList_only4
39 print(f"The number s of subtexts is {len(selectedSubtexts)}")
40 for i in range(len(selectedSubtexts)):
41     #try this later (include info in the content) prompt1.append({"
         role": "system", "content": subtextPrompt + " " + s})
```

```
42      prompt1.append(subtextPrompt + " " + selectedSubtexts[i])
43      prompt1Random.append(subtextPrompt + " " + randomSubtextList[i])

44

45

46 #promptVoting="Based on the subtext generated, select the thought
       that you think is most likely to be the subtext of the next
       speaker. Output a single number corresponding to the thought you
       choose. Output this number as {Output:num}"
47 promptVoting="Based on the subtext generated and the conversation
       history provided, select the thought that is most likely to be
       the subtext of the next speaker. Output a single number
       corresponding to the thought you choose. Output this number as {
       Output:num}. "

48

49 promptVotingScoring="Based on the subtext generated and the
       conversation history provided, output a score corresponding to
       the likelihood that the subtext is the true subtext of the next
       speaker. Output this score as a number from 1 to 10 as {Output:
       num}. "

50

51

52 #or can include a newly formatted prompt that just contains relevant
        info (does not include the last prompt)
53 #prompt2A="You will be given an excerpt of a conversation that
       includes two speakers. It will be formatted as {Speaker 1:
       message, Speaker 2: message}, followed by {Thought:}. Reply as
       the next speaker in the conversation only based on the
       conversation history and the subtext thought. Format your reply
       as {Speaker X:message}. Do not say you are a chatbot. Do not
       output any further thoughts."
54 prompt2A="You will be given an excerpt of a conversation between two
        speakers. It will be formatted as {Speaker 1: message, Speaker
```

76

```
    2: message}, followed by {Thought:}. Reply as the next speaker in
     the conversation only based on the conversation history and the
    subtext thought of the speaker replying. Format your reply as {
    Speaker X:message}. Do not say you are a chatbot. Do not output
    any further thoughts."

55

56 #systemPromptRigged="You will be given an excerpt of a conversation
    between two speakers. Format your output as {Output:}. Your task
    is to generate the subtext thought that the next speaker is
    having, based on the following idea:"
57 promptVotingRiggedChoice="You are a cognitive scientist. Based on
    the conversation history, which of these subtexts was the last
    speaker most likely thinking before they spoke? Choose the one
    that best fits the context. Output a single number corresponding
    to the thought you choose. Output this number as {Output:num}. "
58 promptScoringRiggedScoring="You are a cognitive scientist. Based on
    the conversation history and the given subtext, rate how likely
    it is that this subtext reflects what the last speaker was
    thinking before they spoke. Output a score from 1 (very unlikely)
     to 10 (very likely) as {Output:num}"

59

60

61 all_subtexts_parallel_list=[]# a list of length s (number of
    subtexts) containing lists of length n (number of steps)
62 all_subtexts_parallel_list=[prompt1,prompt2A]
63 all_subtexts_random_parallel_list=[prompt1Random,prompt2A]

64

65

66 #[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
    generate next speaker sentence)]
67 prompts_subtext_voting1=[prompt1,promptVoting,prompt2A]#choice
```

```python
prompts_subtext_voting2=[prompt1,promptVotingScoring,prompt2A]#
    scoring

prompts_subtext_voting3=[prompt1,promptVotingRiggedChoice,prompt2A]#
    voting_choice_rigged
prompts_subtext_voting4=[prompt1,promptScoringRiggedScoring,prompt2A
    ]#voting_scoring_rigged


myDict = {}
myDict["prompt_singleStep"] = "You will be given an excerpt of a
    conversation that includes the next speaker. It will be formatted
     as {Speaker 1: message, Speaker 2: message}. Reply as the next
    speaker in the conversation only. Format your reply as {Speaker
    1: message}. Do not say you are a chatbot."
myDict["prompts_subtext_multiStep"] = all_subtexts_parallel_list
myDict["prompts_subtext_multiStep_random"] =
    all_subtexts_random_parallel_list
myDict["subtexts"] = selectedSubtexts
myDict["prompts_subtext_voting1"]=prompts_subtext_voting1
myDict["prompts_subtext_voting2"]=prompts_subtext_voting2
myDict["randomSubtextList"] = randomSubtextList
myDict["prompts_subtext_voting3"] = prompts_subtext_voting3
myDict["prompts_subtext_voting4"]=prompts_subtext_voting4


exportJSONFile(myDict,"prompts/promptsC_only4.json")

promptOld="You will be given an excerpt of a conversation from the
    point of view of one speaker. Your task is to reply as the other
    speaker in the conversation. Reply as a human in the conversation
     not a chatbot."
```

```
89

90

91

92  #instead we want prompts to be structured like this:
93  #[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
        generate next speaker sentence)]
94  }
```

## A.3  Dataset filtering and setup

```
1   import pandas as pd
2   import numpy as np
3   import json
4   import matplotlib.pyplot as plt
5   from datasetUtilities import createDataset,exportJSONFile
6   import csv
7
8   df = pd.read_csv('decoding_df_798_final_notasks.csv')
9   print(df.columns)
10
11
12  sentence=df["sentence"]
13  corrupted=df["corrupted"]
14  speaker=df["speaker"]
15  conv_name=df["conversation_name"]
16
17
18
19  onsets=df["onsets"]
20  offsets=df["offsets"]
21
22  gap_time = np.zeros(len(onsets))  # Convert gap_time to a NumPy
```

```
         array
23  for i in range(len(onsets)):
24      if i == 0:
25          gap_time[i] = 0   # or some default value
26      else:
27          gap_time[i] = onsets[i] - offsets[i - 1]
28
29  print(gap_time[0:5])
30  print(len(gap_time))
31  N=len(speaker)#length of dataset
32  switchConvo = np.zeros(N, dtype=object)
33
34
35  #segmenting data into "two speaker segments", and also segmenting
        out corrupted segments
36  #get first two unique speakers
37  lastTwoSpeakers, index=np.unique(speaker,return_index=True)
38  lastTwoSpeakers=lastTwoSpeakers[np.argsort(index)]
39  lastTwoSpeakers=list(lastTwoSpeakers)[:2]
40
41  print(lastTwoSpeakers)
42
43  for i in range(len(sentence)):
44      curSpeaker=speaker[i]
45      #keep track of last two speakers, add new segment if one of them
         changes and update lastTwoSpeakers
46      if (curSpeaker not in lastTwoSpeakers):
47          prevSpeaker=speaker[i-1]
48          lastTwoSpeakers.clear()
49          lastTwoSpeakers.append(prevSpeaker)
50          lastTwoSpeakers.append(curSpeaker)
51          switchConvo[i]=1
```

```python
52          #print("switch speaker" + str(i))
53          #print(curSpeaker)
54          #print(lastTwoSpeakers)
55
56      #segment out corrupted data
57      if (corrupted[i]==1):
58          switchConvo[i]=1
59      if (i != 0):
60          if (corrupted[i-1]==1):
61              switchConvo[i]=1
62          if (conv_name[i]!=conv_name[i-1]):#add new segment if
    conversation as labelled in data changes
63              switchConvo[i]=1
64
65
66 #all the indicies where a new segment starts
67 start_indices = np.where(switchConvo == 1)[0]
68 start_indices =np.insert(start_indices,0,0)
69
70
71 # Create masks for each segment
72 masks = []
73 for i, start in enumerate(start_indices):
74     # Determine the end of the segment (either next start index or
    end of array)
75     end = start_indices[i+1] if i+1 < len(start_indices) else len(
    switchConvo)
76     # Create a mask for the segment
77     mask = np.zeros_like(switchConvo, dtype=bool)
78     mask[start:end] = True
79
80     #remove corrupted segments
```

```
81    if ( corrupted [ start ]==False ):
82        masks . append ( mask )
83

84

85  #segmenting data into segments based on masks
86  sentenceSegments = [ sentence [ mask ] for mask in masks ] # masked
        sentences
87  speakerSegments = [ speaker [ mask ] for mask in masks ] # masked speaker
         labels
88  gap_timeSegments = [ gap_time [ mask ] for mask in masks ] # masked gap
        times
89  overallConversation =[]
90  overallGapTimes =[]
91

92  #checking for corrupted segments
93  print ( "Checking for corrupted segments" )
94  corruptedSegments = [ corrupted [ mask ] for mask in masks ]
95  for i in range ( len ( corruptedSegments )):
96      if ( corruptedSegments [i ]. any ()):
97          print ( "missed corrupted segment !!!" )
98

99  #finding the percentage of the sentences spoken by the speaker that
        talks the most in each conversation
100 percentageSpeak =[]
101 for i in range ( len ( speakerSegments )):
102     percentageSpeak . append ( speakerSegments [i ]. value_counts ( normalize
        =True ). max ())
103

104 plt . figure ( figsize =(10 , 6))
105 plt . hist ( percentageSpeak , bins =30 , edgecolor ='black' )
106 plt . xlabel ( 'Percentage Spoken' )
107 plt . ylabel ( 'Frequency' )
```

```
108  plt.title('Histogram of percentageSpeak-percentage of sentences
         spoken by the speaker that talks the most in each conversation')
109  plt.show()
110  plt.savefig('percentageSpeak.png')
111
112  #joining sentences from each speaker into conversational turns, and
         formatting data to input to createDataset function
113  for i in range(len(sentenceSegments)):
114      #if conversation dominated by one speaker skip it
115      if (speakerSegments[i].value_counts(normalize=True).max() >
         0.75):
116          continue
117
118      sentenceSegment=list(sentenceSegments[i])
119      speakerSegment=list(speakerSegments[i])
120      gap_timeSegment=list(gap_timeSegments[i])
121
122      prevSpeaker=speakerSegment[0]#get initial speaker
123      curUtterance=""
124
125      conversationHistory=[]
126      gapTimeHistory=[]
127      for j in range(len(sentenceSegment)):
128          #if speaker changes, add utterance so far from prev speaker
         to conversation history
129          if (prevSpeaker!=speakerSegment[j]):
130              conversationHistory.append(curUtterance)
131              curUtterance=""
132              curGapTime=gap_timeSegment[j]
133              gapTimeHistory.append(curGapTime)
134          #add next sentence to current uttranace (add a space if not
         the first sentence)
```

```python
        if j!=0:
            curUtterance=curUtterance+ " " + sentenceSegment[j]
        else:
            curUtterance=sentenceSegment[j]
            gapTimeHistory.append(0)#add 0 gap time for first
    sentence


        #updating the previous speaker
        prevSpeaker=speakerSegment[j]
    if gapTimeHistory:
        gapTimeHistory.pop()  # Remove the last entry of
    gapTimeHistory
    overallConversation.append(conversationHistory)
    overallGapTimes.append(gapTimeHistory)




print(len(overallGapTimes[0]))
print(len(overallConversation[0]))

#shows how many conversations were filtered out
print("Num conversations before:")
print(len(sentenceSegments))
print("Num conversations after:")
print(len(overallConversation))

#plotting number of conversational turns in each conversation
conversationalTurns=[]
for i in overallConversation:
```

```python
165    conversationalTurns.append(len(i))
166
167
168
169 plt.figure(figsize=(10, 6))
170 plt.hist(conversationalTurns, bins=30, edgecolor='black')
171 plt.xlabel('Number of conversational turns')
172 plt.ylabel('Frequency')
173 plt.title('Histogram of number of conversational turns in each
        conversation')
174 plt.show()
175 plt.savefig('conversationLengths.png')
176
177 print("Average number of conversational turns:", np.mean(
        conversationalTurns))
178
179
180 #creating correctly formatted datasets
181 references=list()
182 labels=list()
183 references2=list()
184 labels2=list()
185
186 max_num_turns=15
187
188 references,labels=createDataset(overallConversation,1,max_num_turns
        -1)
189 references2,labels2=createDataset(overallConversation,2,
        max_num_turns-2)
190 references3,labels3=createDataset(overallConversation,3,
        max_num_turns-3)
191 references5,labels5=createDataset(overallConversation,5,
```

```
        max_num_turns -5)
192 references10 , labels10 = createDataset ( overallConversation ,10 ,
        max_num_turns -10)
193 references15 , labels15 = createDataset ( overallConversation ,15 ,
        max_num_turns -15)
194
195 referencesGapTimes15 , labelsGapTimes15 = createDataset ( overallGapTimes
        ,15 , max_num_turns -15)
196
197
198 references3_extra , labels3_extra = createDataset ( overallConversation
        ,3 ,3 -3)
199 referencesGapTimes3_extra , labelsGapTimes3_extra = createDataset (
        overallGapTimes ,3 ,3 -3)
200
201 print ( " Number of label sentences : " + str ( len ( labels ) ) )
202 print ( " labels identical : " )
203 print ( labels10 == labels15 )
204 print ( labels5 == labels10 )
205 print ( labels3 == labels5 )
206 print ( labels2 == labels3 )
207 print ( labels == labels2 )
208
209 exportJSONFile ( references , labels , " lab_data_no_tasks_one_turn . json " )
210 exportJSONFile ( references2 , labels2 , " lab_data_no_tasks_two_turn . json "
        )
211 exportJSONFile ( references3 , labels3 , " lab_data_no_tasks_three_turn .
        json " )
212 exportJSONFile ( references5 , labels5 , " lab_data_no_tasks_five_turn . json
        " )
213 exportJSONFile ( references10 , labels10 , " lab_data_no_tasks_ten_turn .
        json " )
```

```
214 exportJSONFile(references15,labels15,"lab_data_no_tasks_fifteen_turn
       .json")

215

216 exportJSONFile(referencesGapTimes15,labelsGapTimes15,"
       lab_data_no_tasks_gap_times_fifteen_turn.json")

217

218 exportJSONFile(references3_extra,labels3_extra,"
       lab_data_no_tasks_three_turn_extra.json")

219 exportJSONFile(referencesGapTimes3_extra,labelsGapTimes3_extra,"
       lab_data_no_tasks_gap_times_three_turn_extra.json")

220

221 print(len(labelsGapTimes3_extra))
```

```
1      import json

2

3 #formats the dataset for inference

4 #input is structured by conversation and then conversational turns

5 #turnCount is num of prev. conversational turns to include

6 #startPoint is used to make sure different datasets have the same
       references e.g. if startPoint=0, one turn would have more data
       than two turn (one sentence for each conversation)

7 def createDataset(dataset,turnCount,startPoint=0):

8     _prompts=list()

9     _references=list()

10

11     #for each conversation

12     for i in range (len(dataset)):

13         excerpt=dataset[i]

14         #for all sentence sequences we can use

15         if len(excerpt)<turnCount+1:

16             continue

17

18         for j in range (startPoint,len(excerpt)-turnCount):
```

```
19
20            prompt=""
21
22            #add current speaker to prompt string (assumes
    alternating speakers)
23            for k in range(0,turnCount):
24                speakerNum=k %2 + 1
25
26                prompt = prompt + "Speaker " + str(speakerNum) + ":
    " + str(excerpt[j+k])
27
28                if k!=turnCount-1:
29                    prompt=prompt + " "
30
31            #use next utterance as the reference
32            reference=excerpt[j+turnCount]
33
34            _prompts.append(prompt)
35            _references.append(reference)
36
37    return _prompts, _references
38
39
40 def exportJSONFile(_prompts,_references,fileName):
41    myDict = {}
42    myDict["references"] = _prompts
43    myDict["labels"] = _references
44
45    json_object = json.dumps(myDict)
46
47    with open(fileName, "w") as outfile:
48        outfile.write(json_object)
```

## A.4 Inference generation

```
1
2 import torch
3 import time, subprocess, argparse, json, re, math
4 import accelerate.utils
5 import numpy as np
6 import pandas as pd
7 from transformers import AutoModelForCausalLM, AutoTokenizer,
      AutoModel
8 from evaluate import load
9 from pathlib import Path
10 from datetime import datetime
11 import os
12 import math
13 import random
14 from inference import generateInference
15 from distutils.util import strtobool
16
17 def str2bool(v):
18     return bool(strtobool(v))  # Converts "true"/"false", "yes"/"no
      ", "1"/"0" to boolean
19
20
21 numSamplesTotal=200
22
23
24
25 def getSimilarityResults(_references,_labels, sys_prompt, model,
      tokenizer, bertscore, modelType="Llama8B", bertscore_model="
      distilbert-base-uncased",_temperature=1.0):
26     #perform inference and get similarity results through BERTscore
27
```

```
28    #SET UP PROMPTS
29    #Other prompt: You will be given an excerpt of a conversation
      from the point of view of one speaker. Your task is to reply as
      the other speaker in the conversation. Reply as a human in the
      conversation not a chatbot.
30    myinput=[]
31    for i in range (0,len(_references)):
32        myinput.append([{"role": "system", "content": sys_prompt}, {
      "role": "user", "content": _references[i]}] )
33    #getting next sentence and outputting similarity results
34    gen_text=generateInference(myinput, model, tokenizer, modelType,
      _temperature=_temperature)
35    print("Generated text: with tempeature " + str(_temperature))
36    results = bertscore.compute(predictions=gen_text, references=
      _labels, model_type=bertscore_model)
37
38    return results['f1'], gen_text
39
40 def getSimilarityResultsMultiStep(_references,_labels, sys_prompt,
      model, tokenizer, bertscore, modelType="Llama8B", bertscore_model
      ="distilbert-base-uncased"):
41    #perform inference and get similarity results through BERTscore
42
43    #SET UP PROMPTS
44    prompt1=sys_prompt[0]
45    prompt2=sys_prompt[1]
46
47    myinput=[]
48    for i in range (0,len(_references)):
49        myinput.append([{"role": "system", "content": prompt1}, {"
      role": "user", "content": _references[i]}] )
50
```

```python
   gen_text1=generateInference(myinput, model, tokenizer, modelType
   )#subtext thought
   gen_text1 = [re.sub("Thought:", "",i,3) for idx, i in enumerate(
   gen_text1)]


   myinput=[]
   for i in range (0,len(_references)):
       myinput.append([{"role": "system", "content": prompt2}, {"
   role": "user", "content": _references[i]  + "Thought: " +
   gen_text1[i]   }  ])



   gen_text2=generateInference(myinput, model, tokenizer, modelType
   )#final output (generated speaker sentence(s))
   gen_text2 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
   enumerate(gen_text2)]

   results = bertscore.compute(predictions=gen_text2, references=
   _labels, model_type=bertscore_model)

   return results['f1'], gen_text1, gen_text2

def getSimilarityResultsMultiStepRandomThought(_references,_labels,
   sys_prompt, model, tokenizer, bertscore, modelType="Llama8B",
   bertscore_model="distilbert-base-uncased",random_thought_list=
   None):
   #perform inference and get similarity results through BERTscore

   #SET UP PROMPTS
   prompt1=sys_prompt[0]
   prompt2=sys_prompt[1]

```

```
73    myinput =[]
74    for i in range (0,len(_references)):
75        myinput.append([{"role": "system", "content": prompt1}, {"
     role": "user", "content": _references[i]}] )
76
77    #gen_text1 = ["no actual thought"] * len(_references)  # Create
      an empty array of the same length as _references
78
79    if random_thought_list is not None:
80        if len(random_thought_list) >= len(_references):
81            gen_text1 = random.sample(random_thought_list, k=len(
     _references))
82        else:
83            gen_text1 = random.choices(random_thought_list, k=len(
     _references))
84    myinput =[]
85    for i in range (0,len(_references)):
86        myinput.append([{"role": "system", "content": prompt2}, {"
     role": "user", "content": _references[i]  + "Thought: " +
     gen_text1[i]   }  ])
87
88
89    gen_text2=generateInference(myinput, model, tokenizer, modelType
     )#final output (generated speaker sentence(s))
90    gen_text2 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
     enumerate(gen_text2)]
91
92    #gen_text2=gen_text1
93    results = bertscore.compute(predictions=gen_text2, references=
     _labels, model_type=bertscore_model)
94
95    return results['f1'], gen_text1, gen_text2
```

```python
96

97

98 def getSimilarityResultsVoteChoice(_references,_labels, sys_prompt,
       subtexts, model, tokenizer, bertscore, modelType="Llama8B",
       bertscore_model="distilbert-base-uncased"):
99      #perform inference and get similarity results through BERTscore
100     ##[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
       generate next speaker sentence)]
101
102     #SET UP PROMPTS
103     prompt1=sys_prompt[0]  #list of subtext prompts
104     prompt2=sys_prompt[1] # voting prompt
105     prompt3=sys_prompt[2] # generate next speaker sentence
106
107     #generating subtext thoughts
108     #want to keep the references size, so generate along the other
       column and then reshape the array
109     gen_text_thoughts=[]
110     for i in range (0,len(prompt1)):
111         myinput=[]
112         for j in range (0,len(_references)):
113             myinput.append([{"role": "system", "content": prompt1[i
       ]}, {"role": "user", "content": _references[j]}] )
114         gen_text1=generateInference(myinput, model, tokenizer,
       modelType)#subtext thought
115         gen_text1 = [re.sub("Thought:", "",i,3) for idx, i in
       enumerate(gen_text1)]
116         gen_text_thoughts.append(gen_text1)
117
118     # Transpose the list of lists
119     gen_text_thoughts = list(map(list, zip(*gen_text_thoughts)))
       #output should be a list of lists (x is reference, y is different
```

```python
        subtext prompts)
120     print("Gen thoughts shape:")
121     print(len(gen_text_thoughts))
122     print(len(gen_text_thoughts[0]))
123     #voting
124
125     for i in range(len(gen_text_thoughts)):#this should be length
    references
126         #create the list of sutexts to choose from
127         thoughtList=[]
128         thoughtPrompt = "\n".join([f"{idx+1}. {thought}" for idx,
    thought in enumerate(subtexts)])
129
130         myinput=[]
131         for i in range (0,len(_references)):
132             myinput.append([{"role": "system", "content": prompt2},
    {"role": "user", "content": "Conversation History: " +
    _references[i] + "Thought options: " + thoughtPrompt  }  ])
133
134      gen_text2=generateInference(myinput, model, tokenizer, modelType
    ,_temperature=0.0)#result of voting (should just be a single
    number)
135
136     #need to select which one was voted for
137     voted_thought=[]
138     numErrors=0
139     for i in range (len(gen_text_thoughts)):
140         parsedString=None
141         match = re.search(r'Output:\s*(\d+)', gen_text2[i])
142         if match:
143             parsedString = match.group(1)
144         else:
```

```
145            match = re.search(r'\d+', gen_text2[i])
146          if match:
147              parsedString = match.group(0)
148          else:
149              print("Error: no number found in the output: " +
     gen_text2[i])
150              parsedString = "1"  # default to "1" if no match is
     found
151              numErrors+=1
152
153        if int(parsedString) > len(gen_text_thoughts[i]):
154            print("Error: number found is greater than the number of
      thoughts. The gen_text was: " + gen_text2[i])
155            parsedString = "1"
156        if int(parsedString) < 1:
157            print("Error: number found is less than 1, which is not
     a valid option. The gen_text was: " + gen_text2[i])
158            parsedString = "1"
159        voted_thought.append(int(parsedString) - 1)
160
161    print("Number of errors in voting in this batch: " + str(
     numErrors))
162    #need to generate the next speaker sentence  based on thought
     like before
163
164    myinput=[]
165    for i in range (0,len(_references)):
166        myinput.append([{"role": "system", "content": prompt3}, {"
     role": "user", "content": _references[i]  + "Thought: " +
     gen_text_thoughts[i][voted_thought[i]]  }  ])
167
168    gen_text3=generateInference(myinput, model, tokenizer, modelType
```

95

```
      )#final output (generated speaker sentence(s))
169    gen_text3 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
       enumerate(gen_text3)]

170

171    results = bertscore.compute(predictions=gen_text3, references=
       _labels, model_type=bertscore_model)

172

173    return results['f1'], gen_text_thoughts, gen_text2, gen_text3,
       voted_thought

174

175 def getSimilarityResultsVoteScoring(_references,_labels, sys_prompt,
    subtexts, model, tokenizer, bertscore, modelType="Llama8B",
    bertscore_model="distilbert-base-uncased"):
176    #perform inference and get similarity results through BERTscore
177    ##[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
       generate next speaker sentence)]

178

179    #SET UP PROMPTS
180    prompt1=sys_prompt[0]  #list of subtext prompts
181    prompt2=sys_prompt[1] # scoring prompt
182    prompt3=sys_prompt[2] # generate next speaker sentence

183

184    #generating subtext thoughts and then scoring them
185    #want to keep the references size, so generate along the other
       column and then reshape the array
186    gen_text_thoughts=[]
187    for i in range (0,len(prompt1)):
188        myinput=[]
189        for j in range (0,len(_references)):
190            myinput.append([{"role": "system", "content": prompt1[i
       ]}, {"role": "user", "content": _references[j]}] )
191        gen_text1=generateInference(myinput, model, tokenizer,
```

```
            modelType)#subtext thought
192              gen_text1 = [re.sub("Thought:", "",i,3) for idx, i in
            enumerate(gen_text1)]
193              gen_text_thoughts.append(gen_text1)
194          gen_text_thoughts = list(map(list, zip(*gen_text_thoughts)))   #
            Transpose the list of lists, output should be a list of lists (x
            is reference, y is different subtext prompts)
195          print("Gen thoughts shape:")
196          print(len(gen_text_thoughts))
197          print(len(gen_text_thoughts[0]))
198
199          #voting - instead of inputting the whole thought options,
            evaluate all the thought options separately with a probability/
            score, parse the score and then choose the vote with the highest
            score
200
201          #generating scores for each subtext
202          gen_text_scores=[]
203          for i in range (0,len(prompt1)):
204              myinput=[]
205              for j in range (0,len(_references)):
206                  myinput.append([{"role": "system", "content": prompt2},
            {"role": "user", "content": "Conversation History: " +
            _references[j] + "Subtext: " + subtexts[i] }] )
207              gen_text2=generateInference(myinput, model, tokenizer,
            modelType,_temperature=0.0)
208              gen_text_scores.append(gen_text2)
209
210          gen_text_scores = list(map(list, zip(*gen_text_scores)))       #
            Transpose the list of lists, output should be a list of lists (x
            is reference, y is different subtext prompts)
211          print("Gen scores shape:")
```

97

```python
212    print(len(gen_text_scores))
213    print(len(gen_text_scores[0]))
214
215    #parsing the list of scores to get the best thought for each
       reference
216    all_subtext_scores=[]#contains all the scores for all references
       for each subtext
217    best_score_values=[]
218    best_scored_thoughts=[0]*len(gen_text_scores)
219    best_score_indexes=[]
220
221    numErrors=0
222    for i in range(len(gen_text_scores)):#this should be length
       references
223        subtext_scores=[]
224
225        for j in range(len(gen_text_scores[i])):#this should be
       length subtexts
226            parsedString=None
227            match=None
228            match = re.search(r'Output:\s*(\d+)', gen_text_scores[i
       ][j])
229            if match:
230                parsedString = match.group(1)
231            else:
232                match = re.search(r'\d+', gen_text_scores[i][j])
233                if match:
234                    parsedString = match.group(0)
235                else:
236                    print("Error: no number found in the output: " +
       gen_text_scores[i][j])
237                    parsedString = "1"  # default to "1" if no match
```

98

```
 is found
238                numErrors+=1
239
240        subtext_scores.append(int(parsedString))
241
242
243     #get the max of the subtext scores
244     max_score=max(subtext_scores)
245     max_score_index=subtext_scores.index(max_score)
246
247     all_subtext_scores.append(subtext_scores)
248     best_score_values.append(max_score)
249     best_scored_thoughts[i]=gen_text_scores[i][max_score_index]
250     best_score_indexes.append(max_score_index)
251
252  print("Number of errors in parsing scores in this batch: " + str
     (numErrors))
253
254
255  #need to generate the next speaker sentence on the thought that
     is evaluated to be the best one
256  myinput=[]
257  for i in range (0,len(_references)):
258      myinput.append([{"role": "system", "content": prompt3}, {"
     role": "user", "content": "Conversation History: " + _references[
     i] + "Thought: " + best_scored_thoughts[i] }] )
259
260  gen_text3=generateInference(myinput, model, tokenizer, modelType
     )#final output (generated speaker sentence(s))
261  gen_text3 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
     enumerate(gen_text3)]
262
```

```python
263     print(len(gen_text3))
264     results = bertscore.compute(predictions=gen_text3, references=
        _labels, model_type=bertscore_model)
265
266     return results['f1'], gen_text_thoughts, gen_text_scores,
        gen_text3, best_score_indexes, all_subtext_scores,
        best_score_values
267
268 def getSimilarityResultsVoteChoiceRigged(_references,_labels,
        sys_prompt, subtexts, model, tokenizer, bertscore, modelType="
        Llama8B", bertscore_model="distilbert-base-uncased"):
269     #perform inference and get similarity results through BERTscore
270     ##[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
        generate next speaker sentence)]
271
272     #SET UP PROMPTS
273     prompt1=sys_prompt[0]   #list of generation of thought prompts
        inc. subtexts
274     prompt2=sys_prompt[1] # voting prompt (rigged)
275     prompt3=sys_prompt[2] # generate next speaker sentence
276
277     for i in range(len(_references)):#this should be length
        references
278         #create the list of subtexts to choose from
279         thoughtList=[]
280         thoughtPrompt = "\n".join([f"{idx+1}. {thought}" for idx,
        thought in enumerate(subtexts)])
281
282         myinput=[]
283         for i in range (0,len(_references)):
284             myinput.append([{"role": "system", "content": prompt2},
        {"role": "user", "content": "Conversation History: " +
```

```
    _references[i] + "Last Speaker: " + _labels[i] + "Thought options
    : " + thoughtPrompt  }  ])

285

286  gen_text2=generateInference(myinput, model, tokenizer, modelType
    ,_temperature=0.0)#result of voting (should just be a single
    number)

287

288  #need to select which one was voted for
289  voted_thought=[]
290  numErrors=0
291  for i in range (len(_references)):
292      parsedString=None
293      match = re.search(r'Output:\s*(\d+)', gen_text2[i])
294      if match:
295          parsedString = match.group(1)
296      else:
297          match = re.search(r'\d+', gen_text2[i])
298          if match:
299              parsedString = match.group(0)
300          else:
301              print("Error: no number found in the output: " +
    gen_text2[i])
302              parsedString = "1"  # default to "1" if no match is
    found
303              numErrors+=1
304
305      if int(parsedString) > len(prompt1):
306          print("Error: number found is greater than the number of
     thoughts. The gen_text was: " + gen_text2[i])
307          parsedString = "1"
308      if int(parsedString) < 1:
309          print("Error: number found is less than 1, which is not
```

```python
                    a valid option. The gen_text was: " + gen_text2[i])
310                 parsedString = "1"
311             voted_thought.append(int(parsedString) - 1)
312
313      print("Number of errors in voting in this batch: " + str(
         numErrors))
314      #need to generate the next speaker sentence  based on thought
         like before
315
316      # Generate thoughts based on voted_thought directly in batches
317      gen_text_thoughts = [["no thoughts"] * len(sys_prompt[0]) for _
         in range(len(_references))]
318      batch_inputs = []
319      for i in range(len(_references)):
320          voted_thought_index = voted_thought[i]
321          batch_inputs.append([{"role": "system", "content":
         sys_prompt[0][voted_thought_index]},
322                              {"role": "user", "content": _references
         [i]}])
323
324      # Perform batch inference
325      batch_gen_texts = generateInference(batch_inputs, model,
         tokenizer, modelType)
326
327      # Update gen_text_thoughts with the generated results
328      for i, gen_text in enumerate(batch_gen_texts):
329          #print(gen_text)
330          voted_thought_index = voted_thought[i]
331          gen_text_thoughts[i][voted_thought_index] = re.sub("Thought:
         ", "", gen_text, 3)
332      print("Generated thoughts based on voted_thought.")
333
```

```python
334
335       myinput=[]
336       for i in range (0,len(_references)):
337           myinput.append([{"role": "system", "content": prompt3}, {"
          role": "user", "content": _references[i] + "Thought: " +
          gen_text_thoughts[i][voted_thought[i]]  }  ])
338
339       gen_text3=generateInference(myinput, model, tokenizer, modelType
          )#final output (generated speaker sentence(s))
340       gen_text3 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
          enumerate(gen_text3)]
341
342       results = bertscore.compute(predictions=gen_text3, references=
          _labels, model_type=bertscore_model)
343
344       return results['f1'], gen_text_thoughts, gen_text2, gen_text3,
          voted_thought
345
346
347  def getSimilarityResultsVoteScoringRigged(_references,_labels,
          sys_prompt,subtexts, model, tokenizer, bertscore, modelType="
          Llama8B", bertscore_model="distilbert-base-uncased"):
348       #perform inference and get similarity results through BERTscore
349       ##[[prompt1 (generate subtext)*10], prompt2 (voting), prompt3 (
          generate next speaker sentence)]
350
351       #SET UP PROMPTS
352       prompt1=sys_prompt[0]  #list of subtext prompts
353       prompt2=sys_prompt[1] # scoring prompt
354       prompt3=sys_prompt[2] # generate next speaker sentence
355
356
```

```
357     #voting - instead of inputting the whole thought options,
        evaluate all the thought options separately with a probability/
        score, parse the score and then choose the vote with the highest
        score

358

359     #generating scores for each subtext
360     gen_text_scores=[]
361     for i in range (0,len(prompt1)):
362         myinput=[]
363         for j in range (0,len(_references)):
364             myinput.append([{"role": "system", "content": prompt2},
        {"role": "user", "content": "Conversation History: " +
        _references[j] + ". Last Speaker: " + _labels[j] + ". Subtext
        Thought: " + subtexts[i] }] )
365         gen_text2=generateInference(myinput, model, tokenizer,
        modelType,_temperature=0.0)
366         gen_text_scores.append(gen_text2)

367

368     gen_text_scores = list(map(list, zip(*gen_text_scores)))        #
        Transpose the list of lists, output should be a list of lists (x
        is reference, y is different subtext prompts)
369     print("Gen scores shape:")
370     print(len(gen_text_scores))
371     print(len(gen_text_scores[0]))

372

373     #parsing the list of scores to get the best thought for each
        reference
374     all_subtext_scores=[]#contains all the scores for all references
        for each subtext
375     best_score_values=[]
376     best_scored_thoughts=[0]*len(gen_text_scores)
377     best_score_indexes=[]
```

```python
numErrors=0
for i in range(len(gen_text_scores)):#this should be length
references
    subtext_scores=[]


    for j in range(len(gen_text_scores[i])):#this should be
length subtexts
        parsedString=None
        match=None
        match = re.search(r'Output:\s*(\d+)', gen_text_scores[i
][j])
        if match:
            parsedString = match.group(1)
        else:
            match = re.search(r'\d+', gen_text_scores[i][j])
            if match:
                parsedString = match.group(0)
            else:
                print("Error: no number found in the output: " +
gen_text_scores[i][j])
                parsedString = "1"  # default to "1" if no match
is found
                numErrors+=1


        subtext_scores.append(int(parsedString))


    #get the max of the subtext scores
    max_score=max(subtext_scores)
    max_score_index=subtext_scores.index(max_score)
```

105

```
405        all_subtext_scores.append(subtext_scores)
406        best_score_values.append(max_score)
407        best_scored_thoughts[i]=gen_text_scores[i][max_score_index]
408        best_score_indexes.append(max_score_index)
409
410    print("Number of errors in parsing scores in this batch: " + str
       (numErrors))
411
412    # Generate thoughts based on voted_thought directly
413     # Generate thoughts based on voted_thought directly in batches
414    gen_text_thoughts = [["no thoughts"] * len(sys_prompt[0]) for _
       in range(len(_references))]
415    batch_inputs = []
416    for i in range(len(_references)):
417        voted_thought_index = int(best_score_indexes[i])  # Ensure
       it's an integer
418        batch_inputs.append([{"role": "system", "content":
       sys_prompt[0][voted_thought_index]},
419                             {"role": "user", "content": _references
       [i]}])
420
421    # Perform batch inference
422    batch_gen_texts = generateInference(batch_inputs, model,
       tokenizer, modelType)
423
424    # Update gen_text_thoughts with the generated results
425    for i, gen_text in enumerate(batch_gen_texts):
426        #print(gen_text)
427        voted_thought_index = best_score_indexes[i]
428        gen_text_thoughts[i][voted_thought_index] = re.sub("Thought:
       ", "", gen_text, 3)
429    print("Generated thoughts based on voted_thought.")
```

```python
430
431
432
433    #need to generate the next speaker sentence on the thought that
      is evaluated to be the best one
434    myinput=[]
435    for i in range (0,len(_references)):
436        myinput.append([{"role": "system", "content": prompt3}, {"
      role": "user", "content": "Conversation History: " + _references[
      i] + "Thought: " + gen_text_thoughts[i][best_score_indexes[i]] }]
      )
437
438    gen_text3=generateInference(myinput, model, tokenizer, modelType
      )#final output (generated speaker sentence(s))
439    gen_text3 = [re.sub("Speaker [0-9]:", "",i,3) for idx, i in
      enumerate(gen_text3)]
440
441    print(len(gen_text3))
442    results = bertscore.compute(predictions=gen_text3, references=
      _labels, model_type=bertscore_model)
443
444    return results['f1'], gen_text_thoughts, gen_text_scores,
      gen_text3, best_score_indexes, all_subtext_scores,
      best_score_values
445
446
447
448 #want this to be the ame for all prompting types
449 def performAnalysisOnWholeData(references,labels, sys_prompt,
      prompting_type, model_name, model, tokenizer, bertscore,
      shuffled_indices,subtexts=None,_temperature=1.0,
      random_thought_list=None):
```

```
450    batchSize=20
451    batchNum=math.ceil(len(references)/batchSize)
452

453

454    results_dict = {
455        "reference": references,
456        "label": labels,
457    }
458

459    if prompting_type != "multi_parallel":
460        results_dict["gen_speech"]= []
461        results_dict["bertscore"]= []
462

463    if prompting_type == "voting_score" or prompting_type == "
       voting_choice" or prompting_type == "voting_choice_rigged" or
       prompting_type == "voting_score_rigged":
464        results_dict["gens_thoughts"] = []
465        results_dict["gens_vote"] = []
466        results_dict["voted_thoughts_parsed"] = []
467        results_dict["voted_thoughts_true"] = []
468        results_dict["gen_for_voted_thought"] = []
469        results_dict["subtext_for_voted_thoughts"] = []
470

471    for i in range(batchNum):
472        start=batchSize*i
473        end = len(references) if i == batchNum - 1 else batchSize *
       (i + 1)
474        if prompting_type == "single":
475            batch_scores, batch_gens_speech = (
476                getSimilarityResults(references[start:end], labels[
       start:end], sys_prompt, model, tokenizer, bertscore, model_name,
       _temperature=_temperature)
```

```
477                    )
478            if prompting_type == "multi_parallel":
479                for p in range(len(sys_prompt[0])):
480                    print(f"Running for subtext {p}")
481                    batch_prompt=[sys_prompt[0][p],sys_prompt[1]]

482
483                    if random_thought_list is None:
484                        batch_scores, batch_gens_thought,
    batch_gens_speech = (
485                            getSimilarityResultsMultiStep(references[
    start:end], labels[start:end], batch_prompt, model, tokenizer,
    bertscore, model_name)
486                        )
487                    else:
488                        batch_scores, batch_gens_thought,
    batch_gens_speech = (
489                            getSimilarityResultsMultiStepRandomThought(
    references[start:end], labels[start:end], batch_prompt, model,
    tokenizer, bertscore, model_name, random_thought_list=
    random_thought_list)
490                        )
491                    results_dict.setdefault(f"subtext{p}", []).extend([
    subtexts[p]]*(end-start))#may need to expand this by batch size
492                    results_dict.setdefault(f"gen_thought{p}", []).
    extend(batch_gens_thought)
493                    results_dict.setdefault(f"gen_speech{p}", []).extend
    (batch_gens_speech)
494                    results_dict.setdefault(f"bertscore{p}", []).extend(
    batch_scores)

495
496            elif prompting_type == "voting_score":
497                batch_scores, batch_gens_thought, batch_gens_vote,
```

```python
    batch_gens_speech, batch_voted_thought, batch_all_subtext_scores,
     batch_best_score_values = (
498                 getSimilarityResultsVoteScoring(references[start:end
    ], labels[start:end], sys_prompt, subtexts, model, tokenizer,
    bertscore, model_name)
499             )
500         elif prompting_type == "voting_choice":
501             batch_scores, batch_gens_thought, batch_gens_vote,
    batch_gens_speech, batch_voted_thought = (
502                 getSimilarityResultsVoteChoice(references[start:end
    ], labels[start:end], sys_prompt, model, tokenizer, bertscore,
    model_name)
503             )
504         elif prompting_type == "voting_choice_rigged":#
505             batch_scores, batch_gens_thought, batch_gens_vote,
    batch_gens_speech, batch_voted_thought = (
506                 getSimilarityResultsVoteChoiceRigged(references[
    start:end], labels[start:end], sys_prompt, subtexts, model,
    tokenizer, bertscore, model_name)
507             )
508         elif prompting_type == "voting_score_rigged":
509             batch_scores, batch_gens_thought, batch_gens_vote,
    batch_gens_speech, batch_voted_thought, batch_all_subtext_scores,
     batch_best_score_values = (
510                 getSimilarityResultsVoteScoringRigged(references[
    start:end], labels[start:end], sys_prompt, subtexts, model,
    tokenizer, bertscore, model_name)
511             )
512
513         if prompting_type == "voting_choice" or prompting_type == "
    voting_score" or prompting_type == "voting_choice_rigged" or
    prompting_type == "voting_score_rigged":
```

```python
514             voted_thoughts_strings = [batch_gens_thought[i][
        voted_thought] for i, voted_thought in enumerate(
        batch_voted_thought)]
515             voted_thoughts_subtext = [subtexts[i] for i in
        batch_voted_thought]
516
517             results_dict.setdefault("gens_thoughts", []).extend(
        batch_gens_thought)
518             results_dict.setdefault("gens_vote", []).extend(
        batch_gens_vote)#for scoring choice this is the list of generated
         scores
519             if prompting_type=="voting_score" or prompting_type=="
        voting_score_rigged":
520                 results_dict.setdefault("gens_vote_unshuffled", []).
        extend(unshuffleArray(batch_gens_vote, shuffled_indices))
521             results_dict.setdefault("voted_thoughts_parsed", []).
        extend(batch_voted_thought)
522
523             # Undo the shuffling for subtexts, gens_thoughts,
        voted_thoughts
524             reverse_mapping =shuffled_indices                #map from
        order in shuffled array to order in original array
525             # Apply the reverse mapping to restore the original
        order
526             test_array=[shuffled_indices for _ in range(len(
        batch_gens_thought))]
527             unshuffled_subtexts=unshuffleArray(test_array,
        shuffled_indices)
528             unshuffled_gens_thoughts=unshuffleArray(
        batch_gens_thought, shuffled_indices)
529             voted_thoughts_true = [reverse_mapping[
        batch_voted_thought[i]] for i in range(len(batch_voted_thought))]
```

111

```python
            results_dict.setdefault("voted_thoughts_true", []).
    extend(voted_thoughts_true)
            results_dict.setdefault("gen_for_voted_thought", []).
    extend(voted_thoughts_strings)
            results_dict.setdefault("subtext_for_voted_thoughts",
    []).extend(voted_thoughts_subtext)
            results_dict.setdefault("gens_thoughts_unshuffled", []).
    extend(unshuffled_gens_thoughts)
            results_dict.setdefault("unshuffled_subtexts_indices",
    []).extend(unshuffled_subtexts)
            if prompting_type == "voting_score" or prompting_type ==
     "voting_score_rigged":#the part that's unique to voting_scores
                results_dict.setdefault("
    all_subtext_scores_unshuffled", []).extend(unshuffleArray(
    batch_all_subtext_scores, shuffled_indices))
                results_dict.setdefault("best_score_values", []).
    extend(batch_best_score_values)
        if prompting_type != "multi_parallel":
            results_dict["bertscore"].extend(batch_scores)
            results_dict["gen_speech"].extend(batch_gens_speech)


    for key, value in results_dict.items():
        print(f"{key}: {len(value)}")

    if "results" in results_dict:
        results = results_dict["results"]
        print("Mean result: " + str(np.mean(results)))
        print("Standard deviation results: " + str(np.std(results)))
        print("Number of samples: " + str(len(results)))
```

```python
553        print("Std deviation of mean: " + str(np.std(results) / np.
    sqrt(len(results))))

554

555    return results_dict

556

557

558 def unshuffleArray(gens_thoughts, shuffled_indices):
559    temp_array1=[]
560    for i in range(len(gens_thoughts)):
561        temp_array2=[]
562        for j in range(len(shuffled_indices)):
563            index_in_shuffled = shuffled_indices.index(j)
564            temp_array2.append(gens_thoughts[i][index_in_shuffled])
565        temp_array1.append(temp_array2)
566    gens_thoughts=temp_array1
567    return gens_thoughts

568

569

570 def main():
571    parser = argparse.ArgumentParser(description="Run analysis with
    specified model and data.")
572    parser.add_argument('--model_name', type=str, required=True,
    help='Name of the LM to use')
573    parser.add_argument('--data_path', type=str, required=True, help
    ='Path to the data file')
574    parser.add_argument('--prompt_path', type=str, required=True,
    help='Path to the prompt file')
575    parser.add_argument('--num_shuffles', type=int, required=True,
    help='Number of random shuffles to do (first set is always not
    shuffled)')
576    parser.add_argument('--prompting_type', type=str, required=True,
    choices=["single","multi_parallel","voting_choice","voting_score
```

```python
                   ","voting_choice_rigged","voting_score_rigged"], help='Prompting
                   type: standard or voting')
577    parser.add_argument('--temperature', type=float, required=False,
                   help='Temperature of generation sampling', default=1.0)
578    parser.add_argument('--random_subtext', type=str, required=False
                   , help='Whether a set of random subtexts are used', default="
                   False")
579    parser.add_argument('--random_thought', type=str, required=False
                   , help='Whether a set of random thoughts are used', default="
                   False")
580
581    args = parser.parse_args()
582
583    model_name = args.model_name
584    data_path = Path(args.data_path)
585    prompt_path = Path(args.prompt_path)
586    num_shuffles = args.num_shuffles
587    random_shuffle = num_shuffles > 1 #if num_shuffles is 1, don't
                   shuffle
588    prompting_type = args.prompting_type
589    _temperature = args.temperature
590    random_subtext = str2bool(args.random_subtext)
591    random_thought = str2bool(args.random_thought)
592
593
594    bertscore = load("bertscore")
595    #set up device to GPU if available
596    device = torch.device("cuda" if torch.cuda.is_available() else "
                   cpu")
597    print(f"Using device: {device}")
598
599
```

```python
     #LOADING AND SETTING UP THE MODEL
     #set up model id based on model type


     if (model_name=="Llama1B"):
         model_id = "meta-llama/Llama-3.2-1B-Instruct"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.bfloat16, device_map="auto").to(device)
     elif (model_name=="Llama8B"):
         model_id = "meta-llama/Llama-3.1-8B-Instruct"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.bfloat16, device_map="auto").to(device)
     elif (model_name=="Llama70B"):
         model_id ="hugging-quants/Meta-Llama-3.1-70B-Instruct-GPTQ-
     INT4"
         #model_id = "unsloth/Meta-Llama-3.1-70B-Instruct-bnb-4bit"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.float16).to(device)
         # model = AutoModelForCausalLM.from_pretrained(model_id,
     device_map="auto", cache_dir='/scratch/gpfs/ij9216/cached_models
     ').to(device)
         # model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.float16, device_map="auto", cache_dir='/scratch
     /gpfs/ij9216/cached_models')


     print("TESTING AA1")
     tokenizer = AutoTokenizer.from_pretrained(model_id, padding_side
     = "left")
     tokenizer.pad_token_id = tokenizer.eos_token_id


     print("TESTING AA2")
     #print model size
     module_sizes = accelerate.utils.compute_module_sizes(model)
```

```
623    print(f"The model size is {module_sizes['']} * 1e-9} GB")

624

625    if torch.cuda.is_available():
626        gpu_id = 0   # or whichever GPU you're using
627        total_memory = torch.cuda.get_device_properties(gpu_id).
    total_memory
628        print(f"Total GPU memory: {total_memory / 1024**3:.2f} GB")

629

630    with open(prompt_path, 'r') as file:
631        promptsDict = json.load(file)

632

633    with open(data_path, 'r') as file:
634        data = json.load(file)
635    #get the first X datapoints
636    numSamplesTotal = len(data["references"])

637

638    references=data["references"][0:numSamplesTotal]
639    labels=data["labels"][0:numSamplesTotal]
640    print("dataset " + str(data_path) + " loaded")

641

642

643    subtexts=None
644    if prompting_type== "single":
645        prompts=promptsDict["prompt_singleStep"]
646        subtexts=None
647    elif prompting_type== "multi_parallel":
648        if not random_subtext:
649            prompts=promptsDict["prompts_subtext_multiStep"]
650            subtexts=promptsDict["subtexts"]
651        else:
652            prompts=promptsDict["prompts_subtext_multiStep_random"]
653            subtexts = promptsDict["randomSubtextList"]
```

```python
        elif prompting_type== "voting_choice":
            prompts = promptsDict["prompts_subtext_voting1"]
            subtexts = promptsDict["subtexts"]
        elif prompting_type== "voting_score":
            prompts = promptsDict["prompts_subtext_voting2"]
            subtexts=promptsDict["subtexts"]
        elif prompting_type== "voting_choice_rigged":
            prompts = promptsDict["prompts_subtext_voting3"]
            subtexts=promptsDict["subtexts"]
        elif prompting_type== "voting_score_rigged":
            prompts = promptsDict["prompts_subtext_voting4"]
            subtexts=promptsDict["subtexts"]
    if random_thought:
        random_thought_list = pd.read_csv("data/
    RandomGens20250409_Llama8B.csv")["GeneratedText"].tolist()
    else:
        random_thought_list = None

    print(f"Using model: {model_name}")
    print(f"Data path: {data_path}")
    print(f"Prompt path: {prompt_path}")
    print(f"Prompting type: {prompting_type}")

    all_res_dfs = []

    for shuffle_id in range(1, num_shuffles + 1):
        original_indices = list(range(len(prompts[0])))
        shuffled_indices = original_indices[:]
        if random_shuffle:
            if shuffle_id != 1:  # don't shuffle the first time
                random.shuffle(shuffled_indices)
            if prompting_type == "single":
```

```
685                    pass
686              elif prompting_type == "multi_parallel":
687                  prompts[0] = [prompts[0][i] for i in
      shuffled_indices]
688              elif prompting_type == "voting_score":
689                  prompts[0] = [prompts[0][i] for i in
      shuffled_indices]
690              elif prompting_type == "voting_choice":
691                  prompts[0] = [prompts[0][i] for i in
      shuffled_indices]
692              elif prompting_type == "voting_choice_rigged":
693                  prompts[0] = [prompts[0][i] for i in
      shuffled_indices]
694
695              if subtexts is not None:
696                  subtexts = [subtexts[i] for i in shuffled_indices]
697
698          print(prompts)
699
700          today_date = datetime.now().strftime("%Y%m%d")
701          start_time = datetime.now()
702
703          results_dict = performAnalysisOnWholeData(references,labels,
       prompts, prompting_type,model_name, model, tokenizer, bertscore,
      shuffled_indices,subtexts,_temperature=_temperature,
      random_thought_list=random_thought_list)
704
705          res_df=pd.DataFrame.from_dict(results_dict)
706
707
708          #almost exactly the same, just need to get reference/label
      from csv instead of json file like before
```

118

```
709        #and then also add another column which says if voted
    subtext index matches real subtext index, (i.e. whether
    generation was correct
710
711        #also want to do it with multi parallel realistically
712
713        #can probably handle almost everything within the
    performAnalysisOnWholeData function, only the input needs to
    change asw
714
715        res_df.insert(0, 'sample_id', range(1, numSamplesTotal + 1))
716        res_df.insert(1, 'shuffle_id', shuffle_id)
717
718        res_df['shuffle_indices'] = [shuffled_indices] *
    numSamplesTotal
719        res_df['subtexts'] = [subtexts] * numSamplesTotal
720
721        all_res_dfs.append(res_df)
722
723     final_res_df = pd.concat(all_res_dfs, ignore_index=True)
724     print(final_res_df.head())
725
726     # Calculate the duration
727     duration = datetime.now() - start_time
728     print(f"Time taken: {duration}")
729
730     # save the df to a csv
731     # filename: file_path _ prompt_file _ prompt index _ model name
    _ date
732     res_f_name = f"noramlrand{today_date}_{data_path.stem}_{
    prompt_path.stem}_{model_name}_{prompting_type}_shuffle{
    num_shuffles}_temp{_temperature}_random{str(random_subtext)}
```

119

```
        _randomThought{str(random_thought)}.csv"

733

734     # Check and create the directory
735     if not os.path.exists("results"):
736         os.makedirs("results")
737     final_res_df.to_csv(f"results/{res_f_name}")

738

739     print(f"Results saved to results/{res_f_name}")

740

741

742 if __name__ == "__main__":
743     main()
```

```
1

2     import torch
3 import time, subprocess, argparse, json, re, math
4 import accelerate.utils
5 from transformers import AutoModelForCausalLM, AutoTokenizer,
     AutoModel
6 from evaluate import load

7

8 #default values from generation_config.json   "temperature": 0.6, "
     top_p": 0.9,

9

10 def generateInference(inputPrompts, model, tokenizer, modelType="
     Llama8B",debugTime=False, device='cuda',_temperature=0.6):

11

12     #tokenizing inputs
13     texts = tokenizer.apply_chat_template(inputPrompts,
     add_generation_prompt=True, tokenize=False)
14     inputs = tokenizer(texts, padding=True, return_tensors="pt").to(
     device)

15
```

120

```
16      temp_texts=tokenizer.batch_decode(inputs["input_ids"],
    skip_special_tokens=True)#to see how tokenizer structured input
17
18   if debugTime:
19        #generate inferred next sentence
20        print("Starting generation")
21        start_time = time.time()
22
23   if (modelType=="Llama1B"):
24        terminators = [
25            tokenizer.eos_token_id,
26            tokenizer.convert_tokens_to_ids("<|eot_id|>")
27        ]
28        if _temperature==0.0:
29            gen_tokens = model.generate(
30                **inputs,
31                max_new_tokens=50,
32                pad_token_id=tokenizer.eos_token_id,eos_token_id=
    terminators,
33                do_sample=False,
34            )
35        else:
36            gen_tokens = model.generate(
37                **inputs,
38                max_new_tokens=50,
39                pad_token_id=tokenizer.eos_token_id,eos_token_id=
    terminators,
40                temperature=_temperature,
41            )
42   else:
43        if _temperature==0.0:
44            gen_tokens = model.generate(
```

```
45                    **inputs,
46                    max_new_tokens=50,
47                    pad_token_id=tokenizer.eos_token_id,
48                    do_sample=False,
49                )
50          else:
51              gen_tokens = model.generate(
52                    **inputs,
53                    max_new_tokens=50,
54                    pad_token_id=tokenizer.eos_token_id,
55                    temperature=_temperature,
56                )
57
58      if debugTime:
59          print(f"Generation finished, Time: {time.time()-start_time}"
    )
60
61      #decode generated text and remove the prompt
62      gen_text = tokenizer.batch_decode(gen_tokens,
    skip_special_tokens=True)
63      gen_text = [i[len(temp_texts[idx]):] for idx, i in enumerate(
    gen_text)]
64      return gen_text
```

```
1 import torch
2 import time, subprocess, argparse, json, re, math
3 import accelerate.utils
4 import numpy as np
5 import pandas as pd
6 from transformers import AutoModelForCausalLM, AutoTokenizer,
    AutoModel
7 from evaluate import load
8 from pathlib import Path
```

```python
9  from datetime import datetime
10 import os
11 import math
12 import random
13 from inference import generateInference
14 from distutils.util import strtobool
15
16
17
18 def generateRandomText(_references, sys_prompt, model, tokenizer,
       modelType="Llama8B", bertscore_model="distilbert-base-uncased",
       _temperature=1.0):
19      myinput=[]
20      for i in range (0,len(_references)):
21          myinput.append([{"role": "system", "content": sys_prompt}, {
       "role": "user", "content": _references[i]}] )
22       #getting next sentence and outputting similarity results
23       gen_text=generateInference(myinput, model, tokenizer, modelType,
       _temperature=_temperature)
24       print("Generated text: with tempeature " + str(_temperature))
25
26       return gen_text
27
28
29 def main():
30      parser = argparse.ArgumentParser(description="Run analysis with
       specified model and data.")
31      parser.add_argument('--model_name', type=str, required=True,
       help='Name of the LM to use')
32      args = parser.parse_args()
33
34      model_name = args.model_name
```

```python
     #set up device to GPU if available
     device = torch.device("cuda" if torch.cuda.is_available() else "
     cpu")
     print(f"Using device: {device}")
     #LOADING AND SETTING UP THE MODEL
     #set up model id based on model type

     if (model_name=="Llama1B"):
         model_id = "meta-llama/Llama-3.2-1B-Instruct"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.bfloat16, device_map="auto").to(device)
     elif (model_name=="Llama8B"):
         model_id = "meta-llama/Llama-3.1-8B-Instruct"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.bfloat16, device_map="auto").to(device)
     elif (model_name=="Llama70B"):
         model_id ="hugging-quants/Meta-Llama-3.1-70B-Instruct-GPTQ-
     INT4"
         #model_id = "unsloth/Meta-Llama-3.1-70B-Instruct-bnb-4bit"
         model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.float16).to(device)
         # model = AutoModelForCausalLM.from_pretrained(model_id,
     device_map="auto", cache_dir='/scratch/gpfs/ij9216/cached_models
     ').to(device)
         # model = AutoModelForCausalLM.from_pretrained(model_id,
     torch_dtype=torch.float16, device_map="auto", cache_dir='/scratch
     /gpfs/ij9216/cached_models')

     print("TESTING AA1")
     tokenizer = AutoTokenizer.from_pretrained(model_id, padding_side
     = "left")
     tokenizer.pad_token_id = tokenizer.eos_token_id
```

```python
57
58     print("TESTING AA2")
59     #print model size
60     module_sizes = accelerate.utils.compute_module_sizes(model)
61     print(f"The model size is {module_sizes['']} * 1e-9} GB")
62
63     if torch.cuda.is_available():
64         gpu_id = 0  # or whichever GPU you're using
65         total_memory = torch.cuda.get_device_properties(gpu_id).
    total_memory
66         print(f"Total GPU memory: {total_memory / 1024**3:.2f} GB")
67
68
69     numGens=8000
70     batchSize=50
71     references=["Generate some random text" for _ in range(numGens)]
72
73     numBatches=math.ceil(len(references)/batchSize)
74
75     allGeneratedText=[]
76     for i in range(numBatches):
77         startIndex=i*batchSize
78         endIndex=min((i+1)*batchSize,len(references))
79         batchReferences=references[startIndex:endIndex]
80         #generate random text
81         generatedText=generateRandomText(batchReferences, "Generate
    some random text", model, tokenizer, modelType="Llama8B",
    bertscore_model="distilbert-base-uncased",_temperature=2.0)
82         allGeneratedText.extend(generatedText)
83
84     results_dict = {"GeneratedText": allGeneratedText}
85     final_res_df = pd.DataFrame.from_dict(results_dict)
```

```python
86
87
88      today_date = datetime.now().strftime("%Y%m%d")
89

90

91      # save the df to a csv
92      # filename: file_path _ prompt_file _ prompt index _ model name
      _ date
93      res_f_name = f"RandomGens{today_date}_{model_name}.csv"
94

95      # Check and create the directory
96      if not os.path.exists("results"):
97          os.makedirs("results")
98      final_res_df.to_csv(f"results/{res_f_name}")
99

100     print(f"Results saved to results/{res_f_name}")
101

102

103 if __name__ == "__main__":
104     main()
```

## A.5 Language model hyperparameters

Hyperparameters reterived from "generation$_c$$onfig.json$". $It is the same for LlaMa 3.1 1B, LlaMa 3.2 8$

```json
1 {
2   "bos_token_id": 128000,
3   "do_sample": true,
4   "eos_token_id": [
5     128001,
6     128008,
7     128009
8   ],
```

126

```
9    "temperature": 0.6,
10   "top_p": 0.9,
11   "transformers_version": "4.46.2"
12 }
```

## A.6 Banded Ridge Regression Code

dataframe_setup.py

```
1  import pandas as pd
2
3
4  #load transcript and extract sentences
5  # Extract sentences/sentence indices and add them to the transcript
        DataFrame.
6  puncTranscript_path = "monkey_util-whisperx_transcript.csv" # Change
        as needed
7
8  df = pd.read_csv(puncTranscript_path)
9  df.dropna(subset=['start'], inplace=True)
10 df.sort_values("start", inplace=True)
11 df.word.tail(10)
12
13 df.index = list(range(len(df.word)))# fix word index numbering
14 sentenceEnds = df.index[df.word.str.contains('[.?!]')]
15 df.loc[sentenceEnds]
16 sentenceBegs = sentenceEnds + 1
17 sentenceBegs = sentenceBegs.insert(0, 0)
18 sentenceBegs = sentenceBegs.delete(-1)
19
20 sentenceIdx = [0] * len(df.word);
21 sentence = [''] * len(df.word);
```

```
22 for i, j, k in zip(sentenceBegs, sentenceEnds, range(len(
      sentenceBegs))):
23   fullSentence = (df.word[i:(j+1)]).str.cat(sep=' ')#concatenate
      words in sentence
24
25   sentenceIdx[i:(j+1)] = [k] * len(range(i, (j+1)))#set sentence
      index for each word in sentence to be the same
26   sentence[i:(j+1)] = [fullSentence] * len(range(i, (j+1)))
27
28 sentenceIdx = pd.DataFrame({'sentence_idx': sentenceIdx})
29 sentence = pd.DataFrame({'sentence': sentence})
30 df = df.join(sentenceIdx)
31 df = df.join(sentence)
32 df.insert(0, "word_idx", df.index.values)
33
34 df.to_csv('transcript_inc_sentence.csv', index=False)
```

create_embeddings.py

```
1 import h5py
2 import torch
3 import numpy as np
4 import pandas as pd
5
6 from himalaya.backend import set_backend, get_backend
7
8 from accelerate import Accelerator, find_executable_batch_size
9 from transformers import AutoModelForCausalLM, AutoTokenizer
10 from sentence_transformers import SentenceTransformer
11 import gensim.downloader
12 import re
13
14 if torch.cuda.is_available():
15     set_backend("torch_cuda")
```

```python
16        print("Using cuda!")
17
18
19 #load transcript and extract sentences
20 # Extract sentences/sentence indices and add them to the transcript
        DataFrame.
21 transcript_path = "transcript_inc_sentence.csv" # Change as needed
22 df = pd.read_csv(transcript_path)
23
24 # Download 300-dimensional word2vec embeddings
25 model_name = 'word2vec-google-news-300'
26 n_features = 300
27
28 model = gensim.downloader.load(model_name)
29
30
31 #creating static word embeddings
32 transcript_w2v = df.copy()
33
34 # Convert words to lowercase
35 transcript_w2v['word'] = transcript_w2v.word.str.lower()
36 # Filter words
37 transcript_w2v['filtered_word'] = transcript_w2v.word.apply(lambda
        word: re.sub(r"[^\w\s]", "", word))
38
39 # Function to extract embeddings if available
40 def get_vector(word):
41      #word = ''.join(char for char in word if char.isalnum() and char
        not in ["'", ",", ".","?","!","-"])
42      if word in model.key_to_index:
43          return model.get_vector(word, norm=True).astype(np.float32)
44      return np.nan
```

```python
45
46 # Extract embedding for each word
47 transcript_w2v['embedding'] = transcript_w2v.filtered_word.apply(
       get_vector)  # word2vec embeddings are 300-dimensional
48 # Insert zero vectors for entries with no embeddings
49 embedding_length = n_features  # Length of the embeddings (300 for
       word2vec)
50 transcript_w2v['embedding'] = transcript_w2v['embedding'].apply(
51     lambda x: np.zeros(embedding_length, dtype=np.float32) if np.any
       (np.isnan(x)) else x
52 )
53
54
55 # Print out words not found in vocabulary
56 print(f'{(transcript_w2v.embedding.isna()).sum()} words not found:')
57 print(transcript_w2v.filtered_word[transcript_w2v.embedding.isna()].
       value_counts())
58
59
60 #creating contextual word embeddings
61 #model to tokenize the data
62 modelname = "gpt2"
63 context_len = 32
64 device = torch.device("cpu")
65 if torch.cuda.is_available():
66     device = torch.device("cuda", 0)
67     print("Using cuda!")
68
69 model_sentence = SentenceTransformer("intfloat/e5-mistral-7b-
       instruct")
70 # In case you want to reduce the maximum sequence length:
71 model_sentence.max_seq_length = 1000
```

```
72
73  # Load model for tokenization
74  tokenizer = AutoTokenizer.from_pretrained(modelname)
75
76  #explode dataframe using tokens
77  df["hftoken"] = df.word.apply(lambda x: tokenizer.tokenize(" " + x))
78
79  df = df.explode("hftoken", ignore_index=True)
80  df["token_id"] = df.hftoken.apply(tokenizer.convert_tokens_to_ids)
81
82  df.head(10)
83
84  #model to get token embeddings
85  print("Loading model...")
86  model = AutoModelForCausalLM.from_pretrained(modelname)
87
88  print(
89      f"Model : {modelname}"
90      f"\nLayers: {model.config.num_hidden_layers}"
91      f"\nEmbDim: {model.config.hidden_size}"
92      f"\nConfig: {model.config}"
93  )
94  model = model.eval()
95  model = model.to(device)
96
97  # setting up data for getting embeddings (using previous 32 tokens
        as context)
98  # could change this to use e.g. last sentence as context
99  token_ids = df.token_id.tolist()
100
101 print(f"Token id length: {len(token_ids)}")
102
```

131

```python
103 fill_value = 0
104 if tokenizer.pad_token_id is not None:
105     fill_value = tokenizer.pad_token_id
106
107 #data ends up being an array of length token_ids containing arrays
        of length context_len + 1
108 data = torch.full((len(token_ids), context_len + 1), fill_value,
        dtype=torch.long)
109 for i in range(len(token_ids)):
110     example_tokens = token_ids[max(0, i - context_len) : i + 1]
111     data[i, -len(example_tokens) :] = torch.tensor(example_tokens)
112
113
114 print(f"Data has a shape of: {data.shape}") #  torch.Size([5749,
        33])
115
116
117 #extract embeddings of all of the tokens
118 #use accekerator to make extracting features more efficient
119 accelerator = Accelerator()
120 @find_executable_batch_size(starting_batch_size=32)
121 def inference_loop(batch_size=32):
122     # nonlocal accelerator  # Ensure they can be used in our context
123     accelerator.free_memory()  # Free all lingering references
124
125     data_dl = torch.utils.data.DataLoader(
126         data, batch_size=batch_size, shuffle=False
127         )
128
129     top_guesses = []
130     ranks = []
131     true_probs = []
```

```python
132      entropies = []
133      embeddings = []
134
135      #batch size is number of entries in each batch (i.e 32,32... (
         something less than 32 at end))
136
137      with torch.no_grad():
138          i=0
139          for batch in data_dl:
140              # Get output from model
141              output = model(batch.to(device), output_hidden_states=
         True)
142              logits = output.logits
143              states = output.hidden_states
144
145              true_ids = batch[:, -1]
146              brange = list(range(len(true_ids)))
147              logits_order = logits[:, -2, :].argsort(descending=True)
148              batch_top_guesses = logits_order[:, 0]
149              batch_ranks = torch.eq(logits_order, true_ids.reshape
         (-1, 1).to(device)).nonzero()[:, 1]
150              batch_probs = torch.softmax(logits[:, -2, :], dim=-1)
151              batch_true_probs = batch_probs[brange, true_ids]
152              batch_entropy = torch.distributions.Categorical(probs=
         batch_probs).entropy()
153
154              #hidden states is a list of tensors where each tensor
         has shape (batch_size, sequence length, hidden_size)
155              #hidden size is dimensitonality of the model's hidden
         representations
156
157              #extracts the hidden state for the last token of each
```

```
          sequence in the batch, giving a list of 13 tensors of shape (
          batch size, hidden size) i.e (32, 768)
158           batch_embeddings = [state[:, -1, :].numpy(force=True)
          for state in states ]
159
160
161           top_guesses.append(batch_top_guesses.numpy(force=True))
162           ranks.append(batch_ranks.numpy(force=True))
163           true_probs.append(batch_true_probs.numpy(force=True))
164           entropies.append(batch_entropy.numpy(force=True))
165           embeddings.append(batch_embeddings)
166           i+=1
167       print(f"Number of batches is {i}")
168       return top_guesses, ranks, true_probs, entropies, embeddings
169
170 top_guesses, ranks, true_probs, entropies, embeddings =
        inference_loop()
171
172 #embeddings is a list of 180 lists, each list has 13 tensors of
        shape (32, 768)
173 #GPT 2 has 13 layers of embeddings, each word embedding is 768
        dimensions long
174 print(f"There are {len(embeddings[0])} layers of embeddings")
175 print(f"Each word embedding is {embeddings[0][0].shape[1]}
        dimensions long")
176
177
178 embeddingsA=np.hstack(embeddings)#this should be len 180, containing
         a bunch of ndarrays of shape (32,768)
179 print(embeddingsA.shape)
180 embeddingsB=embeddingsA[12,:,:]
181 print(embeddingsB.shape)# Shape: (num tokens, embedding size)
```

```python
182
183 print("List shape")
184 print(len(embeddingsB.tolist()))
185 print(len(embeddingsB.tolist()[0]))
186
187 df["rank"] = np.concatenate(ranks)
188 df["true_prob"] = np.concatenate(true_probs)
189 df["top_pred"] = np.concatenate(top_guesses)
190 df["entropy"] = np.concatenate(entropies)
191 df["tkn_embedding"] = embeddingsB.tolist()
192 df.head(10)
193
194
195 #average embeddings over tokens in each word, and add to dataframe
196 aligned_embeddings = []
197 for _, group in df.groupby("word_idx"): # group by word index
198     indices = group.index.to_numpy()
199     average_emb = embeddingsB[indices].mean(0) # average features
200     aligned_embeddings.append(average_emb)
201 aligned_embeddings = np.stack(aligned_embeddings)
202 print(f"LLM embeddings matrix has shape: {aligned_embeddings.shape}"
    )
203
204 print(len(aligned_embeddings.tolist()))
205
206 #get list of sentences in order
207 # Get the indices of the first occurrence of each unique sentence
208 unique_indices = df["sentence_idx"].drop_duplicates(keep="first").
    index.tolist()
209 # Access each sentence in order
210 sentences = df.loc[unique_indices, "sentence"].tolist()
211
```

```python
212 #create new list of embeddings for each sentence
213 sentence_embeddings=[]
214 sentence_embeddings = model_sentence.encode(sentences)
215 dataLength=df.shape[0]
216
217 # explode the sentence list to line up with words(for each word add
        its sentence embedding to array)
218 unique_word_indices = df["word_idx"].drop_duplicates(keep="first").
        index.tolist()
219 expanded_sentence_embeddings=[]
220 for idx in unique_word_indices:
221     curSentence=df["sentence_idx"].iloc[idx]
222     expanded_sentence_embeddings.append(sentence_embeddings[
        curSentence])
223
224 print(len(expanded_sentence_embeddings))
225
226
227 df_sentence_embeddings = pd.DataFrame({'sentence_embedding':
        expanded_sentence_embeddings})
228 df_word_embeddings = pd.DataFrame({'word_embedding':
        aligned_embeddings.tolist()})
229 df_static_word_embeddings =pd.DataFrame({'static_word_embedding':
        transcript_w2v.embedding.tolist()})
230
231 df_word_embeddings = df_word_embeddings.join(df_sentence_embeddings)
232 df_word_embeddings = df_word_embeddings.join(
        df_static_word_embeddings)
233
234
235
236 print(df_word_embeddings.head())
```

```
237  df_word_embeddings.to_hdf('embeddings.h5', key='df', mode='w')
```

joint_encoding.py

```
1   import mne
2   import h5py
3   import torch
4   import numpy as np
5   import pandas as pd
6   import matplotlib.pyplot as plt
7   from nilearn.plotting import plot_markers
8
9   from himalaya.backend import set_backend, get_backend
10  from himalaya.ridge import RidgeCV
11  from himalaya.scoring import correlation_score
12  from himalaya.scoring import correlation_score_split
13
14  from sklearn.model_selection import KFold
15  from sklearn.pipeline import make_pipeline
16  from sklearn.preprocessing import StandardScaler
17
18  from himalaya.kernel_ridge import Kernelizer, ColumnKernelizer,
        MultipleKernelRidgeCV
19  import joblib
20
21  if torch.cuda.is_available():
22      set_backend("torch_cuda")
23      print("Using cuda!")
24
25  #span of the epochs
26  tMin=-6.0
27  tMax=6.0
28
29  #sampling rate
```

```
30 sampling_rate = 512
31 step_size=16
32
33 #load transcript and load embeddings
34 transcript_path = "transcript_inc_sentence.csv"
35 df = pd.read_csv(transcript_path)
36
37 df_embeddings= pd.read_hdf('embeddings.h5', 'df')
38 #print(df_embeddings.head())
39 aligned_embeddings = np.vstack(df_embeddings["word_embedding"].
       values)
40 sentence_embeddings = np.vstack(df_embeddings["sentence_embedding"].
       values)
41 static_word_embeddings = np.vstack(df_embeddings["
       static_word_embedding"].values)
42
43
44 print(aligned_embeddings.shape)
45 print(sentence_embeddings.shape)
46 print(static_word_embeddings.shape)
47
48 #construct separate dataframe containins words with their start and
       end times
49 df_word = df.groupby("word_idx").agg(dict(word="first", start="first
       ", end="last"))
50
51 #load the preprocessed brain data
52 file_path = "sub-03_task-podcast_desc-highgamma_ieeg.fif"
53 raw = mne.io.read_raw_fif(file_path, verbose=False)
54 picks = mne.pick_channels_regexp(raw.ch_names, "LG[AB]*")
55 raw = raw.pick(picks)
56
```

```python
57 print(raw)
58
59 #map seconds to samples by multiplying to sampling rate
60 events = np.zeros((len(df_word), 3), dtype=int)
61 events[:, 0] = (df_word.start * raw.info['sfreq']).astype(int)
62 print(events.shape)
63
64 #create epochs and downsample to 32 Hz
65 epochs = mne.Epochs(raw,events,tmin=tMin,tmax=tMax,baseline=None,
      proj=False,event_id=None,preload=True,event_repeated="merge")
66 print(f"Epochs object has a shape of: {epochs._data.shape}")
67 epochs = epochs.resample(sfreq=32, npad='auto', method='fft', window
      ='hamming')
68 print(f"Epochs object has a shape of: {epochs._data.shape}")
69
70
71 #create Y matrices
72 #reshape target matrix Y by horizontally stacking electodes and lags
       along the second dimension
73 epochs_data = epochs.get_data(copy=True)
74 epochs_data = epochs_data.reshape(len(epochs), -1)
75 print(f"ECoG data matrix shape: {epochs_data.shape}")
76
77 #align features with target matrix/ECoG data (some bad epochs may be
       dropped)
78 selected_df = df_word.iloc[epochs.selection]
79 averaged_embeddings = aligned_embeddings[epochs.selection]
80 averaged_sentence_embeddings=sentence_embeddings[epochs.selection]
81
82
83 #Change the float precision to float32 for all data to take
      advantage of the GPU memory and computational speed.
```

```python
84  X_word_ctx = averaged_embeddings
85  X_word_static = static_word_embeddings[epochs.selection]
86  X_sentence=averaged_sentence_embeddings
87  X_word_used=X_word_ctx
88
89  X_joint = np.hstack([X_word_used, X_sentence]) # Horizontal-stack
        both embeddings to create joint model
90  Y = epochs_data
91
92  print(f"Joint predictor matrix shape: {X_joint.shape}")
93  print(f"X_word_ctx shape: {X_word_ctx.shape}")
94  print(f"X_word_static shape: {X_word_static.shape}")
95  print(f"X_sentence shape: {X_sentence.shape}")
96  print(f"Y shape: {Y.shape}")
97
98  #change float precision to float32 for all data to take advantage of
        GPU memory/computational speed
99  if "torch" in get_backend().__name__:
100     X_joint = X_joint.astype(np.float32)
101     X_word_ctx=X_word_ctx.astype(np.float32)
102     X_word_static=X_word_static.astype(np.float32)
103     Y = Y.astype(np.float32)
104
105 #do the banded ridge regression
106
107 # Cross validation is done with an outer part split into 2 and inner
        part split into 5
108 inner_cv = KFold(n_splits=5)
109
110 # Make pipeline with kernelizer for each feature space
111 column_pipeline = make_pipeline(
112     StandardScaler(with_mean=True, with_std=True),
```

```python
113    Kernelizer(kernel="linear"),
114 )
115
116 #set up slices
117 width_w = X_word_used.shape[1]
118 width_s = X_sentence.shape[1]
119
120 slice_w = slice(0, width_w)
121 slice_s = slice(width_w, width_w + width_s)
122 print(f"Word slice: {slice_w}")
123 print(f"Sentence slice: {slice_s}")
124
125 # Compile joint column kernelizer
126 column_kernelizer = ColumnKernelizer(
127     [('word', column_pipeline, slice_w),
128      ('sentence', column_pipeline, slice_s)])
129
130
131 # Ridge regression with alpha grid and nested CV
132 solver = 'random_search'
133 n_iter = 20
134 alphas = np.logspace(1, 10, 10)
135 solver_params = dict(n_iter=n_iter, alphas=alphas)
136
137 # Banded ridge regression with column kernelizer
138 banded_ridge = MultipleKernelRidgeCV(kernels="precomputed", solver=
        solver, solver_params=solver_params, cv=inner_cv)
139
140 # Chain transfroms and estimator into pipeline
141 pipeline = make_pipeline(column_kernelizer, banded_ridge)
142
143 def train_joint_encoding(X, Y,_split=False,_epochs_shape=None):
```

```python
corrs = [] # empty array to store correlation results
kfold = KFold(2, shuffle=False) # outer 2-fold cross-validation
setup
for train_index, test_index in kfold.split(X): # loop through
folds

    # Split train and test datasets
    X1_train, X1_test = X[train_index], X[test_index]
    Y_train, Y_test = Y[train_index], Y[test_index]

    # Standardize Y
    scaler = StandardScaler()
    Y_train = scaler.fit_transform(Y_train)
    Y_test = scaler.transform(Y_test)

    pipeline.fit(X1_train, Y_train) # Fit pipeline with
transforms and ridge estimator

    if _split:
        Y_preds = pipeline.predict(X1_test, split=True) # Use
trained model to predict on test set
        corr = correlation_score_split(Y_test, Y_preds) #
Compute correlation score
    else:
        Y_preds = pipeline.predict(X1_test)
        corr = correlation_score(Y_test, Y_preds).reshape(
epochs_shape) # Compute correlation score

    if "torch" in get_backend().__name__: # if using gpu,
transform tensor back to numpy
        corr = corr.numpy(force=True)
```

```python
169
170         corrs.append(corr) # append fold correlation results to
         final results
171     return np.stack(corrs)
172
173
174 alphas = np.logspace(1, 10, 10) # specify alpha values
175 inner_cv_single = KFold(n_splits=5, shuffle=False) # inner 5-fold
         cross-validation setup
176 single_model = make_pipeline(
177     StandardScaler(), RidgeCV(alphas, fit_intercept=True, cv=
         inner_cv_single) # pipeline
178 )
179
180
181 def train_single_encoding(X, Y,_epochs_shape):
182
183     corrs = [] # empty array to store correlation results
184     kfold = KFold(2, shuffle=False) # outer 2-fold cross-validation
         setup
185     for train_index, test_index in kfold.split(X): # loop through
         folds
186
187         # Split train and test datasets
188         X1_train, X1_test = X[train_index], X[test_index]
189         Y_train, Y_test = Y[train_index], Y[test_index]
190
191         # Standardize Y
192         scaler = StandardScaler()
193         Y_train = scaler.fit_transform(Y_train)
194         Y_test = scaler.transform(Y_test)
195
```

```python
196         single_model.fit(X1_train, Y_train) # Fit pipeline with
    transforms and ridge estimator
197         Y_preds = single_model.predict(X1_test) # Use trained model
    to predict on test set
198         corr = correlation_score(Y_test, Y_preds).reshape(
    _epochs_shape) # Compute correlation score
199
200         if "torch" in get_backend().__name__: # if using gpu,
    transform tensor back to numpy
201             corr = corr.numpy(force=True)
202
203         corrs.append(corr) # append fold correlation results to
    final results
204     return np.stack(corrs)
205
206
207 epochs_shape = epochs._data.shape[1:] # number of electrodes *
    number of lags
208 corrs_embedding = train_joint_encoding(X_joint, Y, _split=False,
    _epochs_shape=epochs_shape)
209 print(f"Encoding performance correlating matrix shape: {
    corrs_embedding.shape}")
210 # Encoding performance correlating matrix shape: (2, 127, 128)
211
212 corrs_embedding_split = train_joint_encoding(X_joint, Y, True)
213 corrs_embedding_split = corrs_embedding_split.reshape(2, 2, *
    epochs_shape)
214 print(f"Encoding performance correlating matrix shape: {
    corrs_embedding_split.shape}")
215 #Encoding performance correlating matrix shape: (2, 2, 127, 128)
216 corrs_embedding_word_ctx = train_single_encoding(X_word_ctx, Y,
    _epochs_shape=epochs_shape)  # predictions with just contextual
```

```
        embeddings
217 print(f"Encoding performance correlating matrix shape: {
        corrs_embedding_word_ctx.shape}")
218 # Encoding performance correlating matrix shape: (2, 127, 128)
219 corrs_embedding_word_static = train_single_encoding(X_word_static, Y
        , _epochs_shape=epochs_shape)  # predictions with just contextual
         embeddings
220 print(f"Encoding performance correlating matrix shape: {
        corrs_embedding_word_ctx.shape}")
221
222 #127 electrodes, 128 lags
223
224 joblib.dump(corrs_embedding, 'corrs_embedding_joint_ctx_sentence.pkl
        ')
225 joblib.dump(corrs_embedding_split, '
        corrs_embedding_split_ctx_sentence.pkl')
226 joblib.dump(corrs_embedding_word_ctx, 'corrs_embedding_word_ctx.pkl'
        )
227 joblib.dump(corrs_embedding_word_static, '
        corrs_embedding_word_static.pkl')
228
229 #Note:
230 #     - The dimensions of the following arrays are:
231 #         - corrs_embedding: (2, number_of_electrodes, number_of_lags
        )
232 #         - corrs_embedding_split: (2, 2, number_of_electrodes,
        number_of_lags)
233 #         - corrs_embedding_word_ctx: (2, number_of_electrodes,
        number_of_lags)
234
235
236 print("Done!")
```

# Appendix B

# Engineering Standards

This independent project adheres to the following engineering and industrial standards:

- **Programming Languages**: Python,

- **Open Source Software**: NumPy, PyTorch, HuggingFace Transformers Library

**Standards Compliance**

All tools, programming languages, and methodologies employed in the development of this thesis conform to well-established academic and industrial standards. This adherence ensures that both the theoretical insights and the associated codebases are robust, scalable, and easily adaptable to a variety of practical applications in the engineering domain. Additionally, the thesis consistently utilizes the International System of Units (SI units), ensuring standardization and precision in the presentation of technical data.

# Bibliography

[1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," June 2017.

[2] F. Xu, Q. Hao, Z. Zong, J. Wang, Y. Zhang, J. Wang, X. Lan, J. Gong, T. Ouyang, F. Meng, C. Shao, Y. Yan, Q. Yang, Y. Song, S. Ren, X. Hu, Y. Li, J. Feng, C. Gao, and Y. Li, "Towards large reasoning models: A survey of reinforced reasoning with large language models," 2025.

[3] A. Goldstein, Z. Zada, E. Buchnik, M. Schain, A. Price, B. Aubrey, S. A. Nastase, A. Feder, D. Emanuel, A. Cohen, A. Jansen, H. Gazula, G. Choe, A. Rao, C. Kim, C. Casto, L. Fanda, W. Doyle, D. Friedman, P. Dugan, L. Melloni, R. Reichart, S. Devore, A. Flinker, L. Hasenfratz, O. Levy, A. Hassidim, M. Brenner, Y. Matias, K. A. Norman, O. Devinsky, and U. Hasson, "Shared computational principles for language processing in humans and deep language models," *Nat. Neurosci.*, vol. 25, pp. 369–380, Mar. 2022.

[4] Z. Wu, L. Qiu, A. Ross, E. Akyürek, B. Chen, B. Wang, N. Kim, J. Andreas, and Y. Kim, "Reasoning or reciting? exploring the capabilities and limitations of language models through counterfactual tasks," 2023.

[5] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro,

and Y. Zhang, "Sparks of artificial general intelligence: Early experiments with GPT-4," 2023.

[6] J. Gauthier, J. Hu, E. Wilcox, P. Qian, and R. Levy, "SyntaxGym: An online platform for targeted evaluation of language models," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, (Stroudsburg, PA, USA), Association for Computational Linguistics, 2020.

[7] A. Warstadt and S. R. Bowman, "What artificial neural networks can tell us about human language acquisition," 2022.

[8] K. Mahowald, A. A. Ivanova, I. A. Blank, N. Kanwisher, J. B. Tenenbaum, and E. Fedorenko, "Dissociating language and thought in large language models," 2023.

[9] A. M. Paunov, I. A. Blank, O. Jouravlev, Z. Mineroff, J. Gallée, and E. Fedorenko, "Differential tracking of linguistic vs. mental state content in naturalistic stimuli by language and theory of mind (ToM) brain networks," *Neurobiol. Lang. (Camb.)*, vol. 3, pp. 413–440, July 2022.

[10] E. Fedorenko, A. A. Ivanova, and T. I. Regev, "The language network as a natural kind within the broader landscape of the human brain," *Nat. Rev. Neurosci.*, vol. 25, pp. 289–312, May 2024.

[11] E. Fedorenko and R. Varley, "Language and thought are not the same thing: evidence from neuroimaging and neurological patients," *Ann. N. Y. Acad. Sci.*, vol. 1369, pp. 132–153, Apr. 2016.

[12] G. Granato, A. M. Borghi, and G. Baldassarre, "A computational model of language functions in flexible goal-directed behaviour," *Sci. Rep.*, vol. 10, p. 21623, Dec. 2020.

[13] T. Ruffman, L. Slade, and E. Crowe, "The relation between children's and mothers' mental state language and theory-of-mind understanding," *Child Dev.*, vol. 73, pp. 734–751, May 2002.

[14] O. et al., "OpenAI o1 system card," Dec. 2024.

[15] D. et al., "DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning," 2025.

[16] Q. Team, "Qwq: Reflect deeply on the boundaries of the unknown," Nov 2024.

[17] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," 2022.

[18] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are Zero-Shot reasoners," 2022.

[19] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, "Improving factuality and reasoning in language models through multiagent debate," 2023.

[20] S. Chen, B. Li, and D. Niu, "Boosting of thoughts: Trial-and-error problem solving with large language models," 2024.

[21] N. Shinn, F. Cassano, E. Berman, A. Gopinath, K. Narasimhan, and S. Yao, "Reflexion: Language agents with verbal reinforcement learning," 2023.

[22] L. Yang, Z. Yu, T. Zhang, S. Cao, M. Xu, W. Zhang, J. E. Gonzalez, and B. Cui, "Buffer of thoughts: Thought-augmented reasoning with large language models," 2024.

[23] S. Parashar, B. Olson, S. Khurana, E. Li, H. Ling, J. Caverlee, and S. Ji, "Inference-time computations for LLM reasoning and planning: A benchmark and insights," 2025.

[24] S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan, "Tree of thoughts: Deliberate problem solving with large language models," 2023.

[25] L. Chen, J. Q. Davis, B. Hanin, P. Bailis, I. Stoica, M. Zaharia, and J. Zou, "Are more LLM calls all you need? towards scaling laws of compound inference systems," 2024.

[26] M. K. Ho, R. Saxe, and F. Cushman, "Planning with theory of mind," *Trends Cogn. Sci.*, vol. 26, pp. 959–971, Nov. 2022.

[27] T. Ullman, "Large language models fail on trivial alterations to Theory-of-Mind tasks," 2023.

[28] N. Shapira, M. Levy, S. H. Alavi, X. Zhou, Y. Choi, Y. Goldberg, M. Sap, and V. Shwartz, "Clever hans or neural theory of mind? stress testing social reasoning in large language models," in *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)* (Y. Graham and M. Purver, eds.), (St. Julian's, Malta), pp. 2257–2273, Association for Computational Linguistics, Mar. 2024.

[29] H. Kim, M. Sclar, X. Zhou, R. Bras, G. Kim, Y. Choi, and M. Sap, "FAN-ToM: A benchmark for stress-testing machine theory of mind in interactions," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing* (H. Bouamor, J. Pino, and K. Bali, eds.), (Singapore), pp. 14397–14413, Association for Computational Linguistics, Dec. 2023.

[30] K. Gandhi, J.-P. Fränken, T. Gerstenberg, and N. Goodman, "Understanding social reasoning in language models with language models," in *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

[31] C. Jin, Y. Wu, J. Cao, J. Xiang, Y.-L. Kuo, Z. Hu, T. Ullman, A. Torralba, J. Tenenbaum, and T. Shu, "MMToM-QA: Multimodal theory of mind question answering," in *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (L.-W. Ku, A. Martins, and V. Srikumar, eds.), (Bangkok, Thailand), pp. 16077–16102, Association for Computational Linguistics, Aug. 2024.

[32] Y. Wu, Y. He, Y. Jia, R. Mihalcea, Y. Chen, and N. Deng, "Hi-ToM: A benchmark for evaluating higher-order theory of mind reasoning in large language models," in *Findings of the Association for Computational Linguistics: EMNLP 2023* (H. Bouamor, J. Pino, and K. Bali, eds.), (Singapore), pp. 10691–10706, Association for Computational Linguistics, Dec. 2023.

[33] M. Sclar, S. Kumar, P. West, A. Suhr, Y. Choi, and Y. Tsvetkov, "Minding language models' (lack of) theory of mind: A plug-and-play multi-character belief tracker," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (A. Rogers, J. Boyd-Graber, and N. Okazaki, eds.), (Toronto, Canada), pp. 13960–13980, Association for Computational Linguistics, July 2023.

[34] Y. Gu, O. Tafjord, H. Kim, J. Moore, R. L. Bras, P. Clark, and Y. Choi, "SimpleToM: Exposing the gap between explicit ToM inference and implicit ToM application in LLMs," 2024.

[35] C. Jung, D. Kim, J. Jin, J. Kim, Y. Seonwoo, Y. Choi, A. Oh, and H. Kim, "Perceptions to beliefs: Exploring precursory inferences for theory of mind in large language models," 2024.

[36] J. Wu, Z. Chen, J. Deng, S. Sabour, H. Meng, and M. Huang, "COKE: A cognitive knowledge graph for machine theory of mind," May 2023.

151

[37] H. Kim, M. Sclar, T. Zhi-Xuan, L. Ying, S. Levine, Y. Liu, J. B. Tenenbaum, and Y. Choi, "Hypothesis-driven theory-of-mind reasoning for large language models," Feb. 2025.

[38] A. Afrasiyabi, E. Busch, R. Singh, D. Bhaskar, L. Caplette, N. Turk-Browne, and S. Krishnaswamy, "Looking through the mind's eye via multimodal encoder-decoder networks," 2024.

[39] G. Tuckute, N. Kanwisher, and E. Fedorenko, "Language in brains, minds, and machines," *Annu. Rev. Neurosci.*, vol. 47, pp. 277–301, Aug. 2024.

[40] C. Caucheteux and J.-R. King, "Brains and algorithms partially converge in natural language processing," *Commun. Biol.*, vol. 5, p. 134, Feb. 2022.

[41] A. Goldstein, Z. Zada, E. Buchnik, M. Schain, A. Price, B. Aubrey, S. A. Nastase, A. Feder, D. Emanuel, A. Cohen, A. Jansen, H. Gazula, G. Choe, A. Rao, C. Kim, C. Casto, L. Fanda, W. Doyle, D. Friedman, P. Dugan, L. Melloni, R. Reichart, S. Devore, A. Flinker, L. Hasenfratz, O. Levy, A. Hassidim, M. Brenner, Y. Matias, K. A. Norman, O. Devinsky, and U. Hasson, "Shared computational principles for language processing in humans and deep language models," *Nat. Neurosci.*, vol. 25, pp. 369–380, Mar. 2022.

[42] A. Goldstein, E. Ham, S. A. Nastase, Z. Zada, A. Dabush, B. Bobbi Aubrey, M. Schain, H. Gazula, A. Feder, W. Doyle, S. Devore, P. Dugan, D. Friedman, M. Brenner, A. Hassidim, O. Devinsky, A. Flinker, O. Levy, and U. Hasson, "Correspondence between the layered structure of deep language models and temporal structure of natural language processing in the human brain." July 2022.

[43] G. Tuckute, A. Sathe, S. Srikant, M. Taliaferro, M. Wang, M. Schrimpf, K. Kay, and E. Fedorenko, "Driving and suppressing the human language network using large language models," *Nat. Hum. Behav.*, vol. 8, pp. 544–561, Mar. 2024.

[44] F. Pereira, B. Lou, B. Pritchett, S. Ritter, S. J. Gershman, N. Kanwisher, M. Botvinick, and E. Fedorenko, "Toward a universal decoder of linguistic meaning from brain activation," *Nat. Commun.*, vol. 9, Mar. 2018.

[45] C. Kauf, G. Tuckute, R. Levy, J. Andreas, and E. Fedorenko, "Lexical-semantic content, not syntactic structure, is the main contributor to ANN-brain similarity of fMRI responses in the language network," *Neurobiol. Lang. (Camb.)*, vol. 5, pp. 7–42, Apr. 2024.

[46] D. Jurafsky and J. H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models.* 2024. Online manuscript released August 20, 2024.

[47] E. Smith, O. Hsu, R. Qian, S. Roller, Y.-L. Boureau, and J. Weston, "Human evaluation of conversations is an open problem: comparing the sensitivity of various methods for evaluating dialogue agents," in *Proceedings of the 4th Workshop on NLP for Conversational AI*, (Stroudsburg, PA, USA), Association for Computational Linguistics, 2022.

[48] C.-W. Liu, R. Lowe, I. V. Serban, M. Noseworthy, L. Charlin, and J. Pineau, "How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation," 2016.

[49] A. B. Sai, A. K. Mohankumar, and M. M. Khapra, "A survey of evaluation metrics used for NLG systems," 2020.

[50] S. Mehri, J. Choi, L. F. D'Haro, J. Deriu, M. Eskenazi, M. Gasic, K. Georgila, D. Hakkani-Tur, Z. Li, V. Rieser, S. Shaikh, D. Traum, Y.-T. Yeh, Z. Yu,

Y. Zhang, and C. Zhang, "Report from the NSF future directions workshop on automatic evaluation of dialog: Research directions and challenges," 2022.

[51] Y.-T. Yeh, M. Eskenazi, and S. Mehri, "A comprehensive assessment of dialog evaluation metrics," 2021.

[52] J. I. Choi, M. Collins, E. Agichtein, O. Rokhlenko, and S. Malmasi, "Combining multiple metrics for evaluating retrieval-augmented conversations," in *Proceedings of the Third Workshop on Bridging Human–Computer Interaction and Natural Language Processing*, (Stroudsburg, PA, USA), Association for Computational Linguistics, 2024.

[53] D. Khashabi, G. Stanovsky, J. Bragg, N. Lourie, J. Kasai, Y. Choi, N. A. Smith, and D. S. Weld, "GENIE: Toward reproducible and standardized human evaluation for text generation," 2021.

[54] T. Hashimoto, H. Zhang, and P. Liang, "Unifying human and statistical evaluation for natural language generation," in *Proceedings of the 2019 Conference of the North*, (Stroudsburg, PA, USA), Association for Computational Linguistics, 2019.

[55] P. Gupta, S. Mehri, T. Zhao, A. Pavel, M. Eskenazi, and J. Bigham, "Investigating evaluation of open-domain dialogue systems with human generated multiple references," in *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*, (Stroudsburg, PA, USA), Association for Computational Linguistics, 2019.

[56] R. Lowe, M. Noseworthy, I. V. Serban, N. Angelard-Gontier, Y. Bengio, and J. Pineau, "Towards an automatic Turing test: Learning to evaluate dialogue responses," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (R. Barzilay and M.-Y.

Kan, eds.), (Vancouver, Canada), pp. 1116–1126, Association for Computational Linguistics, July 2017.

[57] A. B. Sai, M. D. Gupta, M. M. Khapra, and M. Srinivasan, "Re-evaluating ADEM: A deeper look at scoring dialogue responses," *Proc. Conf. AAAI Artif. Intell.*, vol. 33, pp. 6220–6227, July 2019.

[58] C. Callison-Burch, M. Osborne, and P. Koehn, "Re-evaluating the role of Bleu in machine translation research," in *11th Conference of the European Chapter of the Association for Computational Linguistics* (D. McCarthy and S. Wintner, eds.), (Trento, Italy), pp. 249–256, Association for Computational Linguistics, Apr. 2006.

[59] A. Smith, C. Hardmeier, and J. Tiedemann, "Climbing mont BLEU: The strange world of reachable high-BLEU translations," in *Proceedings of the 19th Annual Conference of the European Association for Machine Translation*, pp. 269–281, 2016.

[60] E. Sulem, O. Abend, and A. Rappoport, "BLEU is not suitable for the evaluation of text simplification," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, eds.), (Brussels, Belgium), pp. 738–744, Association for Computational Linguistics, Oct.-Nov. 2018.

[61] J. Deriu, A. Rodrigo, A. Otegi, G. Echegoyen, S. Rosset, E. Agirre, and M. Cieliebak, "Survey on evaluation methods for dialogue systems," *arXiv [cs.CL]*, 2019.

[62] F. Jelinek, R. L. Mercer, L. R. Bahl, and J. K. Baker, "Perplexity—a measure of the difficulty of speech recognition tasks," *J. Acoust. Soc. Am.*, vol. 62, pp. S63–S63, Dec. 1977.

[63] S. F. Chen, D. Beeferman, and R. Rosenfeld, "Evaluation metrics for language models," 1998.

[64] X. Gao, Y. Zhang, M. Galley, C. Brockett, and B. Dolan, "Dialogue response ranking training with large-scale human feedback data," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (B. Webber, T. Cohn, Y. He, and Y. Liu, eds.), (Online), pp. 386–395, Association for Computational Linguistics, Nov. 2020.

[65] A. K. Vijayakumar, M. Cogswell, R. R. Selvaraju, Q. Sun, S. Lee, D. Crandall, and D. Batra, "Diverse beam search: Decoding diverse solutions from neural sequence models," 2016.

[66] J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan, "A diversity-promoting objective function for neural conversation models," in *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (K. Knight, A. Nenkova, and O. Rambow, eds.), (San Diego, California), pp. 110–119, Association for Computational Linguistics, June 2016.

[67] Y. Zhang, M. Galley, J. Gao, Z. Gan, X. Li, C. Brockett, and B. Dolan, "Generating informative and diverse conversational responses via adversarial information maximization," in *Advances in Neural Information Processing Systems* (S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds.), vol. 31, Curran Associates, Inc., 2018.

[68] W. Xu, S. Pakhomov, P. Heagerty, E. Horvitz, E. R. Bradley, A. Campbell, J. D. Woolley, A. Cohen, D. Ben-Zeev, and T. Cohen, "Perplexity and proximity: Large language model perplexity complements semantic distance metrics for the detection of incoherent speech." 2025.

[69] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, "BERTScore: Evaluating text generation with BERT," 2019.

[70] G. Datta, N. Joshi, and K. Gupta, "Analysis of automatic evaluation metric on low-resourced language: Bertscore vs bleu score," in *Speech and Computer* (S. R. M. Prasanna, A. Karpov, K. Samudravijaya, and S. S. Agrawal, eds.), (Cham), pp. 155–162, Springer International Publishing, 2022.

[71] S. Xu, W. Xie, L. Zhao, and P. He, "Chain of draft: Thinking faster by writing less," 2025.

[72] C. Beaudoin, É. Leblanc, C. Gagner, and M. H. Beauchamp, "Systematic review and inventory of theory of mind measures for young children," *Front. Psychol.*, vol. 10, p. 2905, 2019.

[73] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Sravankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru, B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell, C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz, D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hupkes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán, F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cucurell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra, I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, and J. a. Shah, "The llama 3 herd of models,"

[74] R. Stureborg, D. Alikaniotis, and Y. Suhara, "Large language models are inconsistent and biased evaluators," May 2024.

[75] J. E. Eicher and R. F. Irgolič, "Reducing selection bias in large language models," 2024.

[76] M. Sap, H. Rashkin, D. Chen, R. LeBras, and Y. Choi, "SocialIQA: Commonsense reasoning about social interactions," Apr. 2019.

[77] T. Lanham, A. Chen, A. Radhakrishnan, B. Steiner, C. Denison, D. Hernandez, D. Li, E. Durmus, E. Hubinger, J. Kernion, K. Lukošiūtė, K. Nguyen, N. Cheng, N. Joseph, N. Schiefer, O. Rausch, R. Larson, S. McCandlish, S. Kundu, S. Kadavath, S. Yang, T. Henighan, T. Maxwell, T. Telleen-Lawton, T. Hume, Z. Hatfield-Dodds, J. Kaplan, J. Brauner, S. R. Bowman, and E. Perez, "Measuring faithfulness in Chain-of-Thought reasoning," 2023.

[78] M. Turpin, J. Michael, E. Perez, and S. R. Bowman, "Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting," 2023.

[79] S. H. Tanneru, D. Ley, C. Agarwal, and H. Lakkaraju, "On the hardness of faithful Chain-of-Thought reasoning in large language models," 2024.

[80] C. Agarwal, S. H. Tanneru, and H. Lakkaraju, "Faithfulness vs. plausibility: On the (un)reliability of explanations from large language models," 2024.

[81] A. Welivita, Y. Xie, and P. Pu, "A large-scale dataset for empathetic response generation," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, eds.), (Online and Punta Cana, Dominican Republic), pp. 1251–1264, Association for Computational Linguistics, Nov. 2021.

[82] M. J. Traxler and M. A. Gernsbacher, eds., *Handbook of Psycholinguistics*. San Diego, CA: Academic Press, 2 ed., Nov. 2006.

[83] M. K. Tanenhaus, M. J. Spivey-Knowlton, K. M. Eberhard, and J. C. Sedivy, "Integration of visual and linguistic information in spoken language comprehension," *Science*, vol. 268, pp. 1632–1634, June 1995.

[84] J. E. Hanna, M. K. Tanenhaus, and J. C. Trueswell, "The effects of common ground and perspective on domains of referential interpretation," *J. Mem. Lang.*, vol. 49, pp. 43–61, July 2003.

[85] C. G. Chambers, M. K. Tanenhaus, and J. S. Magnuson, "Actions and affordances in syntactic ambiguity resolution," *J. Exp. Psychol. Learn. Mem. Cogn.*, vol. 30, pp. 687–696, May 2004.

[86] S. E. Brennan, "Invited talk: Processes that shape conversation and their implications for computational linguistics," in *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, (Hong Kong), pp. 1–11, Association for Computational Linguistics, Oct. 2000.

[87] G. Aist, J. Allen, E. Campana, L. Galescu, C. A. G. Gallo, S. C. Stoness, M. Swift, and M. Tanenhaus, "Software architectures for incremental understanding of human speech," in *Interspeech 2006*, (ISCA), pp. aper 1869–Wed2FoP.5–0, ISCA, Sept. 2006.

[88] Z. Zada, S. A. Nastase, B. Aubrey, I. Jalon, S. Michelmann, H. Wang, L. Hasenfratz, W. Doyle, D. Friedman, P. Dugan, L. Melloni, S. Devore, A. Flinker, O. Devinsky, A. Goldstein, and U. Hasson, "The "podcast" ECoG dataset for modeling neural activity during natural language comprehension." Feb. 2025.

[89] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,

D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in python," *J. Mach. Learn. Res.*, vol. 12, p. 2825–2830, Nov. 2011.

[90] T. Dupré la Tour, M. Visconti di Oleggio Castello, and J. L. Gallant, "The voxelwise modeling framework: a tutorial introduction to fitting encoding models to fMRI data." Apr. 2024.

[91] T. Dupré la Tour, M. Eickenberg, A. O. Nunez-Elizalde, and J. L. Gallant, "Feature-space selection with banded ridge regression," *NeuroImage*, vol. 264, p. 119728, 2022.

[92] D. Barrett, *Supernormal Stimuli: How Primal Urges Overran Their Evolutionary Purpose.* New York: W. W. Norton & Company, 2010.

[93] C. Fernyhough and A. M. Borghi, "Inner speech as language process and cognitive tool," *Trends Cogn. Sci.*, vol. 27, pp. 1180–1193, Dec. 2023.