
Text2Decision: Decoding Latent Variables in Risky Decision Making from Think Aloud Text

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Understanding human thoughts can be difficult, as scientists usually rely on observ-
2 ing behaviors. The Think-Aloud protocol, where people talk about their thoughts
3 while making decisions, provides a more direct way to study thoughts. However,
4 past research on this topic has mostly been qualitative. Recent advancements in
5 artificial intelligence and natural language processing provide the potential for
6 more quantitative analysis of language data. This study introduces **Text2Decision**,
7 a model trained on task questions from a large-scale task collection, used to decode
8 decision tendencies in risky decision-making from Think-Aloud texts. We test
9 our model in both human and GPT-4 simulated Think-Aloud text data about risky
10 decision-making, which are out-of-distribution in the training. Our findings demon-
11 strate the model’s performance in capturing GPT-4 manipulated decision personas
12 and in unveiling heuristic decision tendencies from humans. **Text2Decision** demon-
13 strates its capability by training on basic task outlines and theoretical frameworks
14 and generalizing to unseen empirical Think-Aloud text data. This not only allows
15 decoding individual differences from these texts but also extends to analyzing
16 large-scale domain datasets. This study shed light on AI integration in cognitive
17 research for the AI4Science paradigm.

18 1 Introduction

19 Understanding human thoughts is one of the major goals of Cognitive Science, yet observing thoughts
20 is hard. Inspired by the behaviorist approach, modern Computational Cognitive Scientists have tried
21 to infer hidden thought processes by fitting computational models to behavioral data that is easy to
22 observe — such as button presses and response times. While this approach has had some notable
23 successes, it also suffers from several limitations, not least of which is that the models’ design is often
24 colored by the researcher’s own cognitive experiences and introspections[Wilson and Collins, 2019].

25 One of the more direct methodologies to access human thoughts is to simply ask people to speak them
26 aloud, via the Think-Aloud procedure [Simon and Ericsson, 1984]. However, due to the complexities
27 and intricacies of linguistic data, traditional analyses of Think Aloud data have been limited by
28 human coding capacities to be largely qualitative and relatively small scale. As an example of this
29 traditional approach, and of the type of task we will later use in this paper, Brandstätter and Gussmack
30 [2013] used the Think-Aloud procedure in risky decision-making tasks. In this study, the researchers
31 hand-coded people’s utterances according to which kind of decision-making strategy they might
32 be using and whether this was closer to a holistic strategy like Prospect Theory, where options in

33 a choice problem are assigned a single Expected Utility [Kahneman, 2011], or a heuristic strategy,
34 where the features of options are compared one by one [Gigerenzer and Gaissmaier, 2011].

35 While this study provided support for the heuristic view of decision-making, the hand-coding approach
36 was time-consuming, subjective, and hard to replicate, all of which challenge the deeper application
37 of the Think-Aloud method [Gu, 2014].

38 The contemporary landscape of artificial intelligence, notably in the realm of natural language
39 processing (NLP), has brought about transformative changes. Large language models (LLMs) are
40 now equipped to process, understand, and even reason with language data, thereby offering the
41 potential to bridge the gap between qualitative nuances and quantitative rigor [Zhao et al., 2023].
42 With these advances, we have the opportunity to process Think-Aloud text data with advanced
43 quantitative techniques, decoding latent interpretable variables directly from Think-Aloud text,
44 without the intervention of human coding.

45 In our present study, we introduce the **Text2Decision** neural network model. This model is designed
46 to decode both overt human behaviors and the subtle, underlying variables at play during risky
47 decision-making tasks, all from Think-Aloud text narratives.

48 2 Risky decision-making task

49 To help ground our exposition of the **Text2Decision** model, we first introduce the cognitive task that
50 is the focus of this paper. As our first foray into the Think-Aloud paradigm, we sought to replicate
51 the decision-making study of Brandstätter and Gussmack [2013]. This study combines 18 decision
52 problems from Kahneman and Tversky’s classic studies [Kahneman and Tversky, 1979] with the
53 Think-Aloud procedure. For each decision problem, participants choose between two gambles, each
54 offering varying outcome-probability combinations (e.g. \$10 for sure or \$20 with 50% chance). We
55 had both humans and artificial agents complete the task.

Experiments

Human Experiment

Please try to speak your thoughts aloud. When you reach a decision, register your choice by clicking or tapping the option you prefer.

Win \$1000 with 25% chance Lose \$0 with 75% chance Win \$240 with 100% chance

RECORDING

Example think-aloud text:

Okay, win 1000 dollars with 25 percent chance lose 100, 100, lose 0 dollars with 75 percent chance win 240 dollars with 100 percent chance well obviously, I am going with the 240 dollars because I know and I can guarantee that I am going to get those 240 dollars

GPT-4 Experiment

H You are a human participant in a think-aloud experiment. And you are {type, description}. Now please read the question and try to make a decision between two gambles. Please indicate your thoughts about making a decision, and explicitly indicate your decision with a standard format "I will choose Option A/B". Which offer do you prefer? Option A: win 1000 dollars with 25% chance and lose 0 dollars with 75% chance; Option B: win 240 dollars with 100% chance.


 To make this decision, I'll calculate the expected value of each option. The expected value is the sum of the probabilities of each outcome multiplied by the value of each outcome. For Option A, the expected value is (1000.0 dollars * 25.0 %) + (0.0 dollars * 75.0 %), which equals 250 dollars. For Option B, ..., Option A has a higher expected value than Option B. So, I will choose Option A

Figure 1: Risky Think-Aloud decision-making task run on humans and GPT-4

56 2.1 Human experiment

57 76 undergraduate students recruited from the University of Arizona Psychology Subject Pool completed the task online recording both their choices and verbal utterances of their thoughts. Transcripts

59 were initially generated using OpenAI’s Whisper [OpenAI, 2023b] and subsequently verified by
60 human research assistants. The study was approved by the University of Arizona IRB.

61 2.2 GPT-4 experiment

62 To acquire a Think-Aloud and behavioral dataset from GPT-4, we instruct GPT-4 to Think-Aloud
63 and decide as if it were a human participant with the exact same trials in the human experiment.
64 Simulating the same sample size (N=76), we also employed GPT-4 and attributed it to one of five
65 decision-making personas: ‘Risk-Averse’, ‘Risk-Seeking’, ‘Rational Decision-Maker’, ‘Probability-
66 Weighted Decision-Maker’, or ‘Outcome-Focused Decision-Maker’. The specifics of these personas
67 and instruction prompts are detailed in the Appendix Section A. By manipulating the prompts,
68 GPT-4 was instructed to generate varied styles of Think-Aloud responses. For decisions, GPT-4 was
69 guided to use a consistent format, “I will choose A/B”, enabling us to extract choices using regular
70 expressions.

71 3 Text2Decision model

72 To decode interpretable latent decision variables from Think-Aloud data, we need a model that
73 bridges between semantic and theory-driven domains. It is essential that this model be robust and
74 versatile, able to capture varied latent variables across distinct data patterns. While many models
75 require substantial empirical datasets, our **Text2Decision** framework is designed to train efficiently
76 on easily generated synthetic datasets and seek interpretations on empirical Think-Aloud datasets.

77 3.1 Model

78 The **Text2Decision** model is a fully connected neural network that translates text into a compact
79 set of pre-defined decision variables. The goal is that applying the model to Think-Aloud data will
80 allow us to extract the decision variables used by the agent (human or GPT-4) to make its decision. In
81 the risky decision-making task, these decision variables include key features of the gambles that are
82 thought to drive decisions such as the expected values, probabilities, losses, and gains associated with
83 each gamble.

84 To train the network to perform this transformation from text to decision variables, we made use of a
85 large collection of 14,568 risky decision problems used in Peterson et al. [2021]. Our basic approach
86 was to input the text of the question and train the network to predict the vector of decision variables
87 we computed from the decision problem.

88 More concretely, inputs to the network took the form of text embeddings produced using OpenAI’s
89 *text_embedding_ada_002* model, based on the mere question descriptions [OpenAI, 2023a]. For
90 example, for the question used in Figure 2, the input text would be ‘1000 dollars with 50.0% chance,
91 0 dollars with 50.0% chance.’.

92 The decision variables to be predicted by the network included both heuristic and normative measures
93 that have been hypothesized to drive human decision-making Kahneman [2011]. Heuristic decision
94 variables included **maximum gain, minimum gain, maximum loss, minimum loss, maximum
95 plus median gain, probability of maximum gain, probability of minimum gain, probability
96 of maximum loss, probability of minimum loss, probability of maximum plus median gain**
97 [Brandstätter and Gussmack, 2013]. Normative decision variables included **Expected Utility** and
98 **Entropy**. Thus, for the example decision used in Figure 2, the decision variable vector would be
99 [1000, 0, 0, 0, 1000, 0.5, 0.5, 0, 0, 1, 500, 1].

100 The model’s objective is to learn the correlation between text and decision variables, decoding
101 pertinent Think-Aloud text into interpretable decision variables, which could be further investigated
102 in clustering individual differences or uncovering algorithmic information in decision-making.

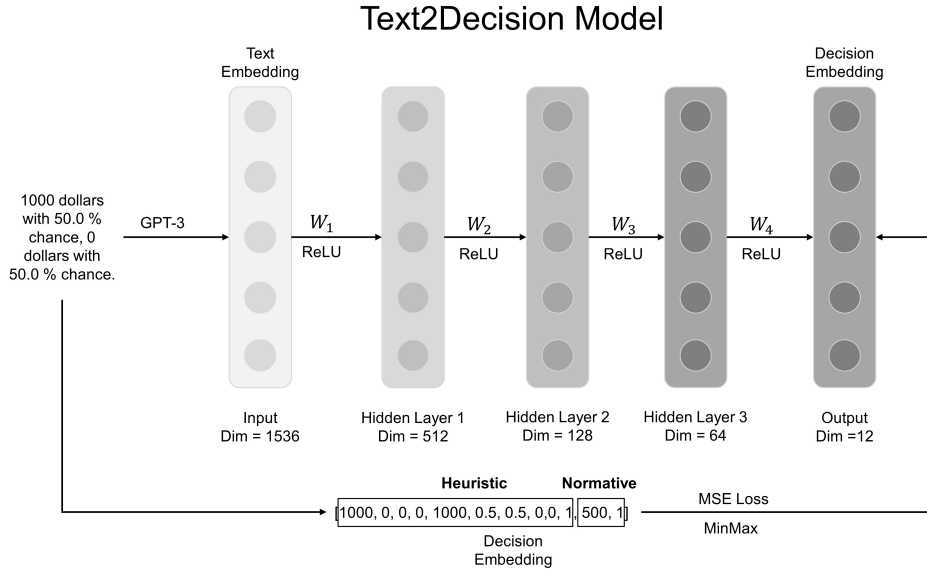


Figure 2: Text2Decision model structure and training illustration

103 3.2 Training and validation

104 For each question in the collection, we generated corresponding text and decision variables. The
 105 dataset is partitioned into 80%/10%/10% segments for training, validation, and testing, respectively.
 106 Our loss function is the Mean Squared Error (MSE), defined for a 12-dimensional decision variable
 107 as:

$$\text{MSE} = \frac{1}{12} \sum_{d=1}^{12} \sum_{i=1}^N (y_{i,d} - \hat{y}_{i,d})^2$$

108 Considering the varied scales across the decision variable dimensions (e.g., probabilities between 0
 109 and 1 versus outcomes from -3000 to 3000), we employed min-max normalization for each dimension
 110 prior to training, which ensures a balanced training impact. Using the Adam optimizer with batch
 111 gradient descent, we minimize the MSE across 200 epochs at a learning rate of 0.001. Both training
 112 and validation losses are monitored per dimension, verifying thorough model transformations (Figure
 113 3). Validation outcomes underscore the model’s proficiency in computing decision variables from
 114 question descriptions.

115 4 Results

116 Next, we asked whether transforming Think-Aloud text to decision variables via the **Text2Decision**
 117 network enabled us to better predict choices of both humans and GPT-4.

118 To predict choices from Think-Aloud texts, we devised four logistic regression models to assess if
 119 our **Text2Decision** model offers enhanced predictive performance relative to basic text embeddings
 120 from GPT-3.

121 In particular, we tested three different flavors of **Text2Decision**. In the first, ‘Text2Decision em-
 122 bedding,’ we used the decoded decision variables as inputs to the logistic model. In the second,
 123 ‘Text2Decision relative Euclidean distance,’ we compared the decoded decision variables from the
 124 Think-Aloud text to the decoded variables generated from each of the two options. The idea here
 125 is that the closest option to the text is more likely to be chosen. In the third flavor of the model,

Training & Validation

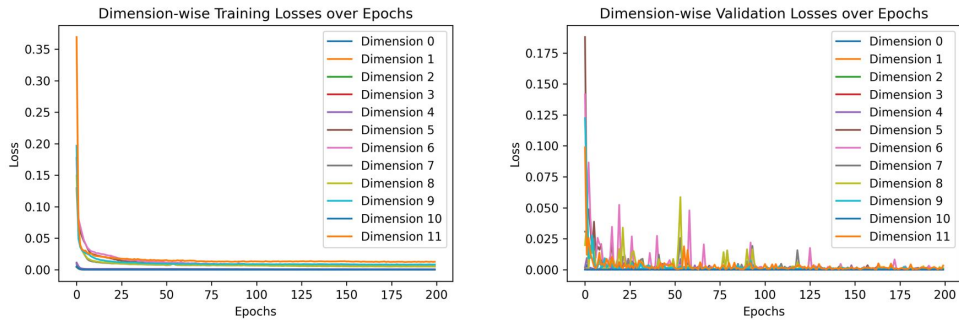


Figure 3: Text2Decision model training and validation loss

Table 1: Decision prediction performance of logistic regression models

Model Name	Mean Accuracy (\pm SE)	
	GPT-4	Human
Raw embedding (PCA to 12 Dimensions)	0.69 ± 0.01	0.55 ± 0.01
Text2Decision embedding	0.58 ± 0.01	0.62 ± 0.01
Text2Decision relative Euclidean distance ^a	0.59 ± 0.01	0.52 ± 0.01
Text2Decision multi-dimensional relative Euclidean distance ^a	0.65 ± 0.01	0.68 ± 0.01

^a The relative distance means the differences in distances from the Text2Decision model transformed text embedding to decision embeddings of two given options.

126 ‘Text2Decision multi-dimensional relative Euclidean distance,’ we adopted a similar idea of compar-
 127 ing the Euclidean distance of decoded decision variables from the Think-Aloud text to the decoded
 128 variables generated from each of the two options but leaving them separately computed for each
 129 dimension as a per regressor in the logistic model. Finally, as a baseline, we trained a model based
 130 on the raw text embeddings compressed to 12 dimensions by principal component analysis, ‘Raw
 131 embedding.’

132 All models we evaluated based on Leave-One-Out-Cross-Validation (LOOCV). Results of this
 133 analysis are shown in Table 1 for both the GPT-4 and human Think-Aloud data.

134 For the GPT-4 dataset, the baseline model (Raw embedding) is the most accurate in predicting agents’
 135 choices.¹ For the human dataset, the model leveraging multi-dimensional relative distances between
 136 transformed text embedding and the two provided options (Text2Decision multi-dimensional relative
 137 Euclidean distance) excels over its counterparts. This suggests that our Text2Decision model can
 138 effectively predict choices based on Think-Aloud text.

139 4.1 Decoding manipulated personas in GPT-4

140 A key application of our model is to decode latent interpretable variables. In the GPT-4 experiment,
 141 we introduced varied decision-maker types to assess the model’s capability in discerning individual
 142 differences in risky decision-making. Using **Text2Decision**, we transformed each text embedding into
 143 the decision variables and computed the variance for each dimension across the five decision-maker

¹In fact, we achieved near-perfect accuracy of 99% with GPT-4’s original output. This likely stems from GPT-4’s consistent phrasing, such as ‘I will choose A/B’, potentially being deterministically encoded into the text embedding. However, when we masked decisions by substituting A/B with X, performance significantly declined. We masked the decision information because we want to ensure the model learns the capacity of inference from the Think-Aloud reasoning process, but not keywords of decisions in a statistical pattern reflected in text embedding.

144 types. Given that all decision-makers underwent identical trials, we attribute variance differences
145 among them to individual disparities.

146 As depicted in Figure 4, the five decision-maker types display varied variances across certain
147 dimensions. Notably, the 'Probability-Weighted Decision Maker' type manifests the largest variances
148 in 'Expected Utility', 'maximum loss', and 'probability of minimum loss'. This aligns with its
149 descriptive prompt: 'Relies on explicit probabilities to estimate expected values and opts for choices
150 with the highest perceived value.' For a comprehensive description of all types, refer to the Appendix
151 (see Section A).

152 These findings indicate that the transformed text embeddings from Think-Aloud data can unveil
153 latent variables in risky decision-making tasks, aligning with the ground truth from our GPT-4
154 manipulations. Consequently, this method holds promise for hypothesis testing in experimental
155 settings.

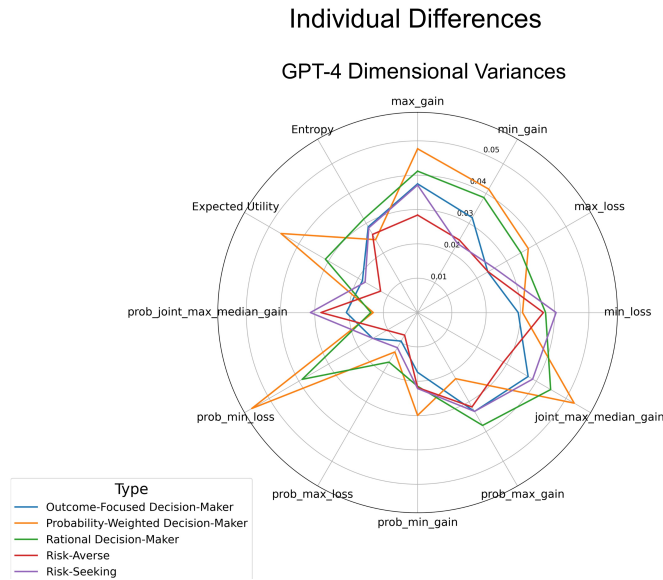


Figure 4: Decoding individual differences in GPT-4 data by computing variance for each dimension of Text2Decision transformed Think-Aloud embedding

156 4.2 Decoding individual differences in humans

157 In a similar vein, we sought to decode latent variables from human data. Without explicit labels for
158 categorization in this dataset, we employed K-means to cluster and segment individuals into five
159 categories. Variance calculations for each label's participants are illustrated in Figure 5.

160 Distinctive decision-making styles emerge for each participant cluster. For instance, Cluster 0
161 participants show the highest variance in 'probability of maximum loss' and 'probability of minimum
162 loss', suggesting a loss-focused approach. Conversely, Cluster 3 participants exhibit pronounced
163 variance in 'probability of max gain' and 'probability of maximum plus median gain', indicating a
164 gain-centric perspective with a probability over outcome emphasis. Interestingly, no clusters display
165 significant variance in the Expected Utility dimension. This hints at humans being more heuristic-
166 driven in their decision-making for this task, in contrast to the normative-driven GPT-4. The heuristic
167 nature of human risky decision-making revealed by the variance representations in Text2Decision
168 transformed embeddings aligns with the findings in behavioral experiments [Kahneman and Tversky,
169 1979] and Think-Aloud experiments with hand-coding [Brandstätter and Gussmack, 2013].

170 Integrating unsupervised learning methods with our **Text2Decision** model's output, we present a
171 framework that allows for interpretable insights into clusters, labels, or principal components.

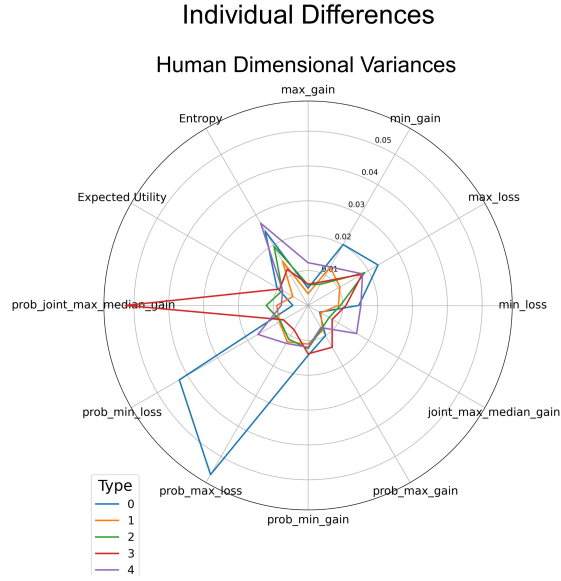


Figure 5: Decoding individual differences in human data by computing variance for each dimension of Text2Decison transformed Think-Aloud embedding

172 5 Discussion

173 In this research, we introduced the **Text2Decision** framework model designed for decoding latent
 174 variables from Think-Aloud texts. We assessed the model’s efficacy in behavioral prediction and
 175 discerning individual differences. Our findings indicate that the model adeptly extracts pertinent
 176 information from raw Think-Aloud text embeddings, enhancing behavioral predictions, particularly
 177 with human datasets. Moreover, the model demonstrated the capability to decode latent, interpretable
 178 variables.

179 Our framework seeks to bridge the gap between hypothesis-driven and data-driven methodologies in
 180 the Cognitive Sciences by integrating the Think-Aloud protocol with Large Language Models. His-
 181 torically, Cognitive Scientists have depended on experimental manipulations to validate hypotheses.
 182 With the advent of computational modeling, there’s been a shift towards quantitatively characterizing
 183 behaviors to provide generative explanations. Yet, a persistent challenge is navigating the trade-off
 184 between the interpretability inherent in hypothesis-driven approaches and the precision of data-driven
 185 ones. Unlike models anchored to specific hypotheses, our **Text2Decision** framework is trained on
 186 general task settings, ensuring a broader explanatory capacity on empirical data, accommodating
 187 the potential for both hypotheses testing and interpretable data-driven investigations. By mapping
 188 semantic spaces to interpretable decision spaces, we aim to decode latent variables, making them
 189 comprehensible and primed for further exploration.

190 Moving forward, our aim is to both deepen and broaden the **Text2Decision** framework. In terms
 191 of depth, we plan to integrate computational modeling, behavioral analysis, and neural recordings
 192 (e.g., fMRI or EEG) to facilitate more robust hypothesis testing and extract clearer, interpretable
 193 patterns[Schneider et al., 2023]. Moreover, we intend to explore diverse participant populations,
 194 considering factors such as race, gender, culture, age, and mental health, to better understand their
 195 Think-Aloud representations during risky decision-making. Broadening our scope, we aspire to adapt
 196 our framework to more intricate tasks, including learning, planning, and challenges like sorting,
 197 clustering, and compositionality. These tasks usually contain rich slow cognitive processes, whereas
 198 Think-Aloud texts may be more useful to decode complex strategies and algorithms. We are also
 199 keen to assess whether our strategy of training on basic task settings (or easily generated synthetic
 200 data) retains its efficacy in these diverse contexts.

201 In conclusion, our framework, in tandem with LLMs, heralds a promising avenue for deciphering
202 human thought processes via Think-Aloud methodologies.

203 **References**

- 204 Eduard Brandstätter and Manuela Gussmack. The cognitive processes underlying risky choice.
205 *Journal of Behavioral Decision Making*, 26(2):185–197, 2013.
- 206 Gerd Gigerenzer and Wolfgang Gaissmaier. Heuristic decision making. *Annual review of psychology*,
207 62:451–482, 2011.
- 208 Yongqi Gu. To code or not to code: Dilemmas in analysing think-aloud protocols in learning strategies
209 research. *System*, 43:74–81, 2014.
- 210 Daniel Kahneman. *Thinking, fast and slow*. macmillan, 2011.
- 211 Daniel Kahneman and Amos Tversky. Oprospect theory: An analysis of decision under risk, 1
212 *econometrica*. *March*, 47(2):2635291, 1979.
- 213 OpenAI. OpenAI embedding models, 2023a. URL [https://platform.openai.com/docs/
214 guides/embeddings/embedding-models](https://platform.openai.com/docs/guides/embeddings/embedding-models). Accessed: [September 24, 2023].
- 215 OpenAI. OpenAI whisper speech-to-text model, 2023b. URL [https://platform.openai.com/
216 docs/guides/speech-to-text](https://platform.openai.com/docs/guides/speech-to-text). Accessed: [January 15, 2023].
- 217 Joshua C Peterson, David D Bourgin, Mayank Agrawal, Daniel Reichman, and Thomas L Griffiths.
218 Using large-scale experiments and machine learning to discover theories of human decision-making.
219 *Science*, 372(6547):1209–1214, 2021.
- 220 Steffen Schneider, Jin Hwa Lee, and Mackenzie Weygandt Mathis. Learnable latent embeddings for
221 joint behavioural and neural analysis. *Nature*, pages 1–9, 2023.
- 222 Herbert Alexander Simon and K Anders Ericsson. Protocol analysis: Verbal reports as data. (*No
223 Title*), 1984.
- 224 Robert C Wilson and Anne GE Collins. Ten simple rules for the computational modeling of behavioral
225 data. *Elife*, 8:e49547, 2019.
- 226 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
227 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv
228 preprint arXiv:2303.18223*, 2023.

229 **A GPT-4 Experiment Prompts**

230 "You are a human participant in a Think-Aloud experiment. And you are {type, description}. Now
231 please read the question and try to make a decision between two gambles. Please indicate your
232 thoughts about making a decision, and explicitly indicate your decision with a standard format 'I will
233 choose Option A/B'. Which offer do you prefer? Option A: win 1000 dollars with 25% chance and
234 lose 0 dollars with 75% chance; Option B: win 240 dollars with 100% chance."

235 **A.1 Types of Decision Makers**

236 **1. Risk-Averse:**

237 Prefers options with predictable outcomes and minimal risk, even if potential rewards are
238 lower.

239 **2. Risk-Seeking:**

240 Drawn to high-reward options even if they come with significant risks.

241

3. Rational Decision-Maker:

242

Analyzes all available information and weighs pros and cons to maximize the outcome.

243

4. Probability-Weighted Decision-Maker:

244

Relies on explicit probabilities to estimate expected values and chooses options with the highest perceived value.

245

246

5. Outcome-Focused Decision-Maker:

247

Prioritizes potential outcomes, especially extreme values, and may avoid options with possible losses even if the expected value is positive.

248