

Linking Strategies to Think Aloud in A Stochastic Learning Task

Zhenlong Zhang¹¹, Hanbo Xie²¹, Travis Baker³, Megan Peters⁴, Robert C. Wilson²²

¹ Bloomberg School of Public Health, Johns Hopkins University, Baltimore, USA

² School of Psychology, Georgia Institute of Technology, Atlanta, USA

³ Center for Molecular and Behavioral Neuroscience, Rutgers State University, Newark, USA

⁴ Department of Cognitive Sciences, University of California, Irvine, USA

zzhan352@jh.edu, hanboxie1997@gatech.edu, teb81@newark.rutgers.edu, megan.peters@uci.edu, rwilson337@gatech.edu

Abstract

Understanding human thoughts is a key goal of cognitive science. Behavioral observations alone limit insight into cognition. The think-aloud protocol, where participants verbalize thoughts, offers a direct probe into reasoning but is underutilized due to challenges in subjectivity and scalability. Advancements in natural language processing (NLP) enable computational analysis of think-aloud data, yet little work explores its role in strategy learning. We test whether think-aloud reports reveal strategy use in a stochastic learning task where participants verbalized their strategies. Our results show diverse strategy usage, with a preference for persistent choices. Think-aloud analysis suggests participants rely on distinct meta-strategies to guide learning. Clustering and predictive modeling reveal strong alignment between choices and verbalized strategies. These findings highlight think-aloud as a scalable tool with NLP techniques for studying high-level cognition, shedding light on a promising paradigm for cognitive sciences.

Keywords: Think Aloud, Natural Language Processing, Strategies, Learning, Reasoning

Introduction

Understanding human thoughts is an ultimate goal of cognitive science. For a long time, cognitive scientists have relied on designing experiments and observing human behavior to infer underlying mental processes. Typically, scientists form hypotheses and test them either by manipulating variables in experiments or by proposing computational models to analyze behavioral data. However, this framework can be biased and may limit long-term understanding of the human mind, especially for high-level cognition, whose processes are often intractable from pure choice behavior. These indirect probes into cognitive processes can, to a large extent, be mitigated by the *Think Aloud* protocol—a traditional method in which participants verbalize their thoughts during an experiment (Simon & Ericsson, 1984). Despite its advantages, the qualitative nature of think-aloud data makes its analysis subjective, labor-intensive, and difficult to scale, which hinders its broader contribution to modern cognitive science and psychology research paradigms.

With advancements in Natural Language Processing (NLP) and Large Language Models (LLMs), there is now potential to revisit this traditional protocol using modern quantitative methods, enabling a more efficient and scalable approach to

handling think-aloud data. Indeed, recent research has leveraged neural network models to decode cognitive variables from think-aloud text embeddings in a risky decision-making task (Xie, Xiong, & Wilson, 2023). More advanced and diverse approaches include using LLMs directly as cognitive models to predict choice behavior (Xie, Xiong, & Wilson, 2024a) or translating think-aloud data into code-like symbolic cognitive models to predict behavior (Xie, Xiong, & Wilson, 2024b). Despite these advancements, little attention has been paid to the role of **strategy usage** in the learning process.

Learning is central to human development and adaptation to uncertain environments. Humans not only learn which actions yield the most benefits from the environment (*model-free learning*) but also infer the underlying rules governing the environment (*model-based learning*). These “*hypothesized*” rules, often referred to as “*strategies*” in the learning process, are difficult to extract from pure choice behavior alone. Past studies have developed various approaches to approximate these rules based on behavioral observations, such as **information-theoretic measures** of choices and behavior (Trepka et al., 2021), **hidden Markov models (HMMs)** (Guenouni & Speekenbrink, 2021), and **neural network models** (Rmus, Pan, Xia, & Collins, 2024). However, all these methods approximate human mental processes by making inferences solely from behavioral data, providing only indirect insight into underlying cognitive mechanisms.

Therefore, our study aims to investigate **strategy usage** through the think-aloud protocol in a strategic stochastic learning task (i.e., the *matching pennies game* (Barraclough, Conroy, & Lee, 2004)). The original study demonstrated that monkeys attempt to “learn” rules in a random rewarding task with their own biases. These biases drive spontaneous strategy hypotheses, testing, and shifts (recrafting), making the task an ideal paradigm for exploring diverse strategy usage. In our research, we seek to determine whether human strategy usage, identified through behavior data, maps onto verbalized think-aloud responses. By bridging the gap between subjective introspection and objective measurement, this approach ultimately offers a powerful framework for uncovering the nuanced cognitive mechanisms underlying learning and adaptation in uncertain environments.

¹Equal contribution.

²Corresponding author: rwilson337@gatech.edu

Methods

Experiment and Participants

To elicit diverse strategy usage, we employed a modified version of the matching-pennies game, inspired by a previous study on monkeys (Barraclough et al., 2004). In this task, participants see two options on each trial—a yellow star and a green star (see Figure 1A). Similar to classical bandit paradigms, participants must select one option, after which they receive feedback indicating whether they earned a reward (1 point) or not. Crucially, the two options have identical underlying probabilities of reward, making the task fundamentally stochastic and devoid of any true “optimal” strategy. To maintain engagement and discourage rapid convergence to random guessing, our cover story instructs participants that they are playing against a “highly intelligent computer agent,” and must “match” this agent’s choice to gain rewards. This narrative encourages participants to hypothesize—and revise—the complex strategies they believe the computer might be using.

Unlike more typical learning tasks, after every 10 trials, participants are prompted to describe the strategies they just used verbally and those they plan to use in the subsequent 10 trials. These verbal reports are recorded as audio and then transcribed automatically using OpenAI’s Whisper speech-to-text model (Radford et al., 2023).

We recruited $N = 68$ undergraduate students from the university to complete the task online. All participants provided informed consent before starting the study and were fully debriefed upon completion regarding the task’s true nature. The study protocol was approved by the university’s Institutional Review Board, ensuring that all ethical guidelines were followed.

Behavioral Metrics

To measure human participants’ behavior in the task, we use some simple and heuristic measurements, like the probability of stay and reward-conditioned probability of stay (aka. win-stay or lose-stay). These metrics will help us to know a basic behavior pattern as well as their reaction to reward feedback. So specifically, for the probability of stay, we define it as:

$$p(\text{stay}) = \frac{\sum_{t=2}^T \mathbb{I}(a_t = a_{t-1})}{T - 1}.$$

While win-stay and lose-stay consider the conditional probability of stay given the reward outcome is positive or negative:

$$p(\text{stay}|\text{win}) = \frac{\sum_{t=2}^T \mathbb{I}(a_t = a_{t-1} \wedge r_{t-1} = 1)}{\sum_{t=2}^T \mathbb{I}(r_{t-1} = 1)}.$$

And:

$$p(\text{stay}|\text{lose}) = \frac{\sum_{t=2}^T \mathbb{I}(a_t = a_{t-1} \wedge r_{t-1} = 0)}{\sum_{t=2}^T \mathbb{I}(r_{t-1} = 0)}.$$

Where:

- a_t : The choice made by the participant at time t .
- a_{t-1} : The choice made by the participant at time $t - 1$.
- r_{t-1} : The reward received at time $t - 1$ (1 for a reward, 0 for no reward).
- T : Total number of trials.

Computational Models of Choice Behaviors

To quantitatively characterize participants’ learning strategies, we proposed five simple heuristic models to fit behavior: (1) Random, (2) Win-Stay Lose-Switch (WSLS), (3) Reinforcement Learning (RL), (4) Choice Kernel (CK) and (5) RL-CK Model.

Random model simply estimates each participant’s choice proportion based on the observations.

Win-Stay Lose-Switch (WSLS) model estimates the trend that participants use a win-stay-lose-switch strategy, and attribute the rest of the actions to noises:

$$p(a_t | a_{t-1}, r_{t-1}; \epsilon) = \begin{cases} 1 - \frac{\epsilon}{2}, & \text{if } r_{t-1} = 1 \text{ and } a_t = a_{t-1}, \\ \frac{\epsilon}{2}, & \text{if } r_{t-1} = 1 \text{ and } a_t \neq a_{t-1}, \\ \frac{\epsilon}{2}, & \text{if } r_{t-1} = 0 \text{ and } a_t = a_{t-1}, \\ 1 - \frac{\epsilon}{2}, & \text{if } r_{t-1} = 0 \text{ and } a_t \neq a_{t-1}. \end{cases} \quad (1)$$

Here, $(1 - \epsilon)$ corresponds to a pure WSLS strategy, while ϵ captures occasional deviations, attributing them to noise or exploratory choices.

Reinforcement learning model considers the reward learning process in the task. Specifically, we use model-free learning, assuming participants simply bonding values of each option to the rewards, without explicit models about the task structure. We use Rescolar-Wagner (RW) model as such:

$$V_{t+1}(a_t) = V_t(a_t) + \alpha(r_t - V_t(a_t)), \quad (2)$$

Where:

- $V_{t+1}(a_t)$: The value of option a_t at time $t + 1$, which is updated based on the reward received at time t .
- $V_t(a_t)$: The value of option a_t at time t , representing the current estimate of the option’s value.
- $\alpha \in [0, 1]$: The learning rate, a parameter that determines how quickly the value of the option is updated based on the new reward.
- r_t : The reward received at time t , where r_t is typically 1 for a reward and 0 for no reward.

The probability of choosing option i on trial t is then given by a softmax function:

$$p(a_t = i) = \frac{\exp(\beta V_t(i))}{\sum_j \exp(\beta V_t(j))}. \quad (3)$$

Choice kernel model, however, only considers the internal repetition trend of one’s behaviors. Therefore, it only considers the last action that the participant has done.

$$CK_{t+1}^k = CK_t^k + \alpha_c (a_t^k - CK_t^k). \quad (4)$$

where $\alpha_c \in [0, 1]$ is a learning rate that adjusts the kernel toward 1 if option k was chosen, and toward 0 otherwise. The probability of choosing option k on trial t follows a softmax function like above (3).

Finally, the **RL-CK model** is a hybrid model. While considering both the updates for values (2) and kernel (4), the model combines both values in the decision-making process with an integrated softmax function:

$$p(a_t = i) = \frac{\exp(\beta [V_t(i) + wCK_t^i])}{\sum_j \exp(\beta [V_t(j) + wCK_t^j])}. \quad (5)$$

All the parameters in the five models are estimated by Maximum Likelihood Estimation (MLE) with the Python package ‘optimize’ on an individual basis. The dynamic of likelihood is post-hoc simulations from the estimated parameters.

Think-Aloud Data Analysis

For all transcriptions of think-aloud audios, we first preprocess the text data, including punctuations, removing uncommon symbols, and correcting misspellings and capitalization. To quantify the semantic meaning of those think-aloud texts, we used an embedding model to convert think-aloud texts into text embeddings. Text embeddings are high-dimensional vectors that represent relative semantic meaning in the vast training dataset. We used **text-embedding-ada-002**, which returns each think-aloud text as a 1536-dimensional vector. These vectors are then visualized in 2D space with t-SNE. We apply an elbow test to determine the number of k-mean and use K-mean clustering to cluster all think-aloud text embeddings into four clusters for post-block descriptions and pre-block planning, respectively. These clusters will help us understand underlying strategies that participants may reveal and their mappings to specific behavioral patterns.

We also deployed machine learning classifier models, including Random Forest (200 estimators, maximum depth of 10, balanced class weights), Support Vector Machines (SVM, RBF kernel, balanced class weights, probability estimation), and Logistic Regression (1000 iterations, balanced class weights) to see whether text embeddings are predictive for block-wise strategies (the best model for each participant at each block). To address class imbalance, we applied SMOTE (Synthetic Minority Over-sampling Technique) to resample the training data, ensuring a balanced distribution of strategies in the dataset. The training process is conducted

with 5-fold cross-validation to ensure the models are not over-fitting. The models were evaluated using accuracy as the primary metric, along with standard error calculated from the cross-validation results. Additionally, confusion matrices and other classification metrics were used to assess model performance.

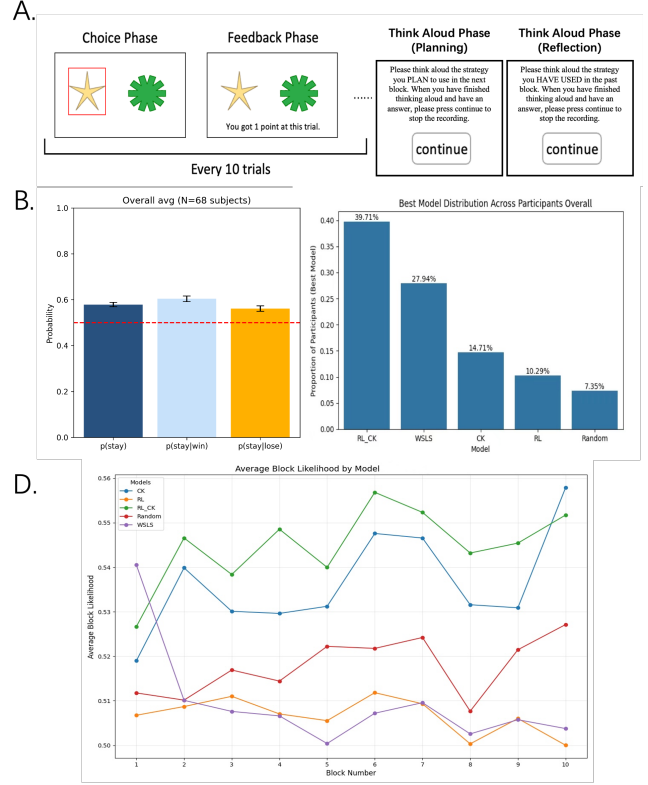


Figure 1: **A. Experimental Procedure.** On each trial, participants choose between two options to maximize rewards. After each choice, the outcome is displayed. Every ten trials, participants are asked to think aloud about past strategies and future plans. **B. Behavioral Tendency.** Participants tend to repeat their previous choice, especially after receiving a reward. **C. Computational Modeling Results.** Five models were fit to participants’ behavior; RL-CK models best explained the majority. **D. Temporal Dynamics.** Block-wise likelihoods show population-level shifts in strategy over time.

Results

Identifying Strategies from Choice Behaviors

To analyze the behavioral patterns of participants in this task, we first computed several basic behavioral metrics. As shown in Figure 1B, participants exhibited a probability of staying higher than chance levels, as well as a conditioned probability of staying after both winning and losing outcomes, all of which were significantly above chance ($p(\text{stay})$: $p < 0.001$, $p(\text{stay}|\text{win})$: $p < 0.001$, $p(\text{stay}|\text{lose})$: $p = 0.001$). This suggests that participants display strong persistent behaviors, regardless of whether they win or lose.

To further quantify participants' behaviors more precisely, we fit each participant's choice data to five proposed computational models (see Method 2.3, *Computational Models of Choice Behaviors*). These models range from naive statistical models to reward-based learning models. As shown in Figure 1C, most participants were best described by the RL-CK model, indicating that their behaviors are both reward-sensitive and persistent in relation to their previous choices. Notably, a relatively large proportion of participants followed heuristic **win-stay-lose-switch (WSLS)** strategies or exhibited simple persistence. The most classical reinforcement learning model accounted for only 10.29% of participants, indicating minimal reliance on this strategy in the task.

We further analyzed model performance by simulating the likelihood of participants' behaviors using fitted parameters, aggregated by blocks (every 10 trials). As shown in Figure 1D, during the first block, the WSLS model was the dominant strategy. However, as the task progressed, both CK and RL-CK models became more dominant, while WSLS rapidly declined in model performance. This trend suggests that strategy structures evolve from simpler to more complex ones, echoing the persistent choice behaviors observed in the basic metrics.

Think-Aloud Contents Reveal Diverse Strategies

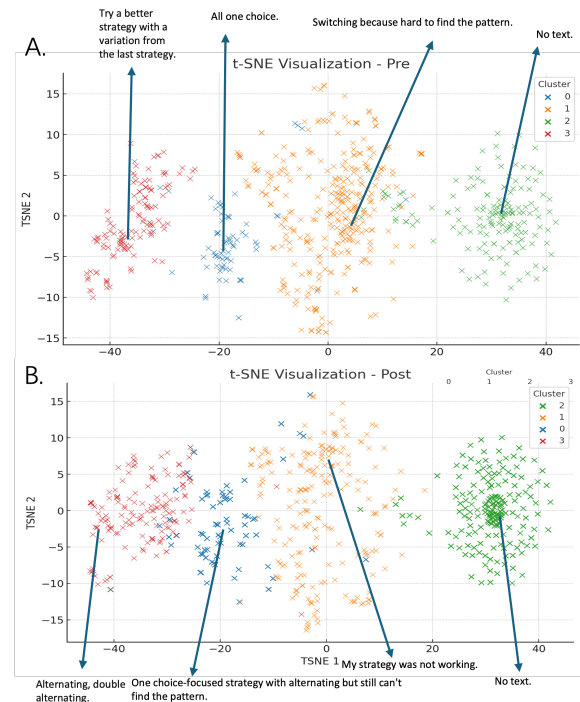


Figure 2: **Think-Aloud Cluster Analysis and Visualizations.** **A.** Think-aloud embedding clusters related to planning for the next blocks. **B.** Think-aloud embedding clusters related to strategy usage in previous blocks.

To quantitatively analyze Think-Aloud text data, we used an embedding model to convert each think-aloud response

into a high-dimensional vector (see Method 2.4, *Think-Aloud Data Analysis*). We then applied t-SNE to reduce these high-dimensional vectors into a 2D space, allowing us to visualize their distribution (Figure 2). To further interpret the semantic meaning of these embeddings, we applied K-means clustering to divide the embeddings into four clusters, each represented by a different color. We then visually inspected actual think-aloud responses in each cluster to conduct a preliminary qualitative analysis.

Pre-block Think-Aloud One prominent cluster in the pre-block data represents participants who frequently switched choices due to difficulty identifying a stable pattern (i.e., Cluster 1 in yellow in Figure 2). A key characteristic of this group is their uncertainty regarding the reward structure, as reflected in statements like:

“So, I feel like it’s more like, whether it’s my own bias or not, but I feel like they wouldn’t typically pick the same one in a row by chance, because that would be a 25% chance should they pick one of two options twice in a row.”

Participants in this cluster struggled to determine whether their choices were based on an actual pattern or just random fluctuations in rewards.

Another cluster in the pre-block data represents participants who committed to a single choice, potentially as an exploratory strategy (i.e., Cluster 0 in blue in Figure 2). This is exemplified by statements such as:

“I’m gonna just press all stars again.”

Rather than attempting to detect a pattern, these participants adopted a simple, fixed-choice strategy, selecting the same option repeatedly in an attempt to test whether it yielded consistent rewards.

A separate pre-block cluster captures participants who focused on an exploration strategy. These participants attempted to find a better strategy by making slight variations from their previous choices (i.e., Cluster 3 in red in Figure 2). One example is:

“I plan to use the same strategy that I use with switching off if it’s the same, but I’m going to start with green again.”

This cluster reflects participants who aimed to maximize rewards by revising their strategy.

Finally, a distinct pre-block cluster consists of participants who did not provide any verbal report, leading to a “No text” category.

Post-block Think-Aloud There are similar strategy clusters in post-block think-aloud responses as in pre-block think-aloud responses. One post-block cluster represents participants who focused on a single-choice strategy with alternating patterns but struggled to identify a clear structure. These

participants attempted systematic switching but remained uncertain about whether a discernible pattern existed. One participant expressed this difficulty:

“That time I actually did use the strategy I was talking about before. I just went with the stars and it happened to be literally all stars. So it kind of gave me a little confidence, but you know, in the previous ones, I also wasn’t, I didn’t do very well. So confidence is there, but it’s not going to last very long.”

This suggests that while some participants committed to a structured choice strategy, they continued to question its effectiveness due to inconsistent outcomes.

Another post-block cluster captures participants who explicitly recognized their strategy as ineffective and sought to adjust it. These individuals reflected on their approach and acknowledged the need for change, as exemplified by the following statement:

“Okay, the strategy has changed. Or the sequence has changed. I think I’m putting too much thought into this, but I’m pretty sure the sequence just changed on me.”

This cluster highlights participants who perceived shifts in the reward structure, leading them to rethink and modify their decision-making strategies.

A separate post-block cluster consists of participants who engaged in alternating and double-alternating strategies, actively attempting to refine their approach through systematic pattern detection. One participant described their evolving method:

“So I’ve tried alternating, double alternating, and now I’ve tried star, star, green, star, green, star... I feel like it worked in the beginning. Then halfway through, it might switch to double greens, then star, then double green, then star. So I’m going to try that one next.”

These participants demonstrated a deliberate effort to uncover a hidden rule, using increasingly complex alternation patterns in an attempt to anticipate future rewards.

For the last cluster, it remains the same as in the pre-block phase, in which no verbal response was recorded.

These patterns reveal a more complex model that participants may be attempting to build—**action sequences** (also known as **successor representations** (Momennejad et al., 2017)). This may explain why the CK model outperforms the pure RL model in capturing decision-making behavior.

Linking Behavioral Strategies to Introspective Descriptions

To investigate the relationship between participants’ learning strategies and their introspective descriptions, we analyzed how the best-fitting computational models were distributed across think-aloud clusters. Additionally, we assessed whether text embeddings derived from verbal reports

could predict participants’ decision-making strategies. For pre-block and post-block think-aloud clusters, we calculated the proportion of each model that best fits each block for each participant. The results show distinctive distributions in the best-fitting model across think-aloud clusters, whether pre-block or post-block (Figure 3A and B).

Specifically, Cluster 0 and Cluster 1 are dominated by blocks where the best-fitting model is the WSL model, while Cluster 2 is dominated by blocks best described by the CK model. Cluster 3 is primarily associated with blocks best captured by the RL-CK model. This distribution aligns with findings from the qualitative analysis of clusters: Cluster 0 reflects persistent choice behavior, while Cluster 1 represents switching due to failed outcomes, both of which correspond to the WSL model. Similarly, Cluster 3 captures participants attempting to refine their strategy, which aligns with the RL-CK model. Meanwhile, Cluster 2, which mainly contains blocks with no verbal responses, may indicate that participants either repeated or switched their choices without actively engaging in the task. Similar patterns were observed in the post-block strategic descriptions.

To further test whether participants’ verbal reports contained meaningful information for predicting their strategies across blocks, we trained classification models using think-aloud text embeddings as features. Using 5-fold cross-validation, we found that all three models performed above the chance level (20%), with Random Forest achieving the highest accuracy (Figure 3C). The classification accuracy was comparable between pre-block and post-block embeddings, suggesting that participants’ verbalized strategies before and after each block contained stable, meaningful patterns reflective of their actual decision-making models. Notably, post-block think-aloud text data was more predictive than pre-block think-aloud text data, likely due to variations between participants’ initial plans and their actual behavior.

These results suggest that introspective verbalizations—though subjective—encode structured information that aligns with computationally inferred learning strategies. By linking text-based insights to behavioral models, this approach offers a novel way to integrate qualitative and quantitative perspectives in cognitive science research.

Discussion

Understanding human thoughts is at the heart of cognitive science. In this study, we investigated human strategies in a stochastic learning task using think-aloud protocols. Behaviorally, we found that participants employed a diverse range of strategies, from heuristic win-stay-lose-switch (WSL) to reinforcement learning with choice kernels (RL-CK). By tracking model likelihoods over time, we observed a transition in strategies at the population level.

To analyze think-aloud data, we applied **visualization and cluster analysis** on text embeddings, allowing us to inspect different strategy patterns. Our findings suggest that clusters in the embedding space reflect distinct strategies, such

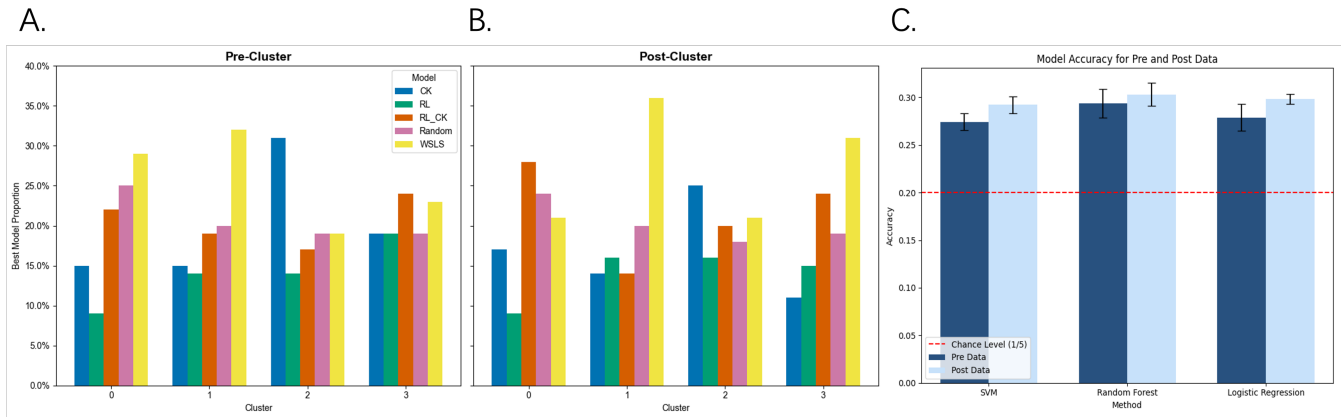


Figure 3: Linking Behavioral Strategies to Think-Aloud Descriptions. (A) Best-fitting model proportions across pre-block think-aloud clusters. (B) Best-fitting model proportions across post-block think-aloud clusters. (C) Model accuracy in predicting decision-making strategies using think-aloud embeddings.

as switching actions, persisting in a chosen action, or varying strategies over time. When combining behavioral model fitting with cluster analysis on think-aloud data, we found a strong alignment: for instance, one cluster was dominated by the WLS model, while another aligned with the RL-CK model. This suggests a meaningful link between participants’ verbalized thoughts and their actual decision-making strategies. Furthermore, by training machine learning models on think-aloud text embeddings, we demonstrated that these verbal descriptions can predict behavioral strategies with high accuracy. These results highlight the potential of think-aloud data as a rich source of information about human cognition in learning tasks.

While our findings confirm that think-aloud data can map onto behavioral strategies, it is important to acknowledge that these strategies are approximations based purely on observed behavior. A key next step is to explore how introspective descriptions might lead to even better approximations of both behavior and cognitive processes. One promising direction is leveraging **Large Language Models (LLMs)** to infer participants’ strategies in the form of programmatic code (i.e., program induction) (Xie et al., 2024b). This approach could generate block-wise models tailored to individuals rather than imposing a one-size-fits-all framework, potentially offering a more precise and hypothesis-free alignment with participants’ own descriptions.

Another important avenue for future research is understanding how strategies evolve over time. In our task, participants actively refined their strategies to maximize rewards, constructing and updating mental models of the task. This process of mental exploration and subsequent testing through real-world actions (Johnson-Laird, 1983) could provide deeper insights into meta-learning. Investigating how these evolving strategies manifest in both think-aloud data and behavior would further clarify the cognitive mechanisms underlying learning and adaptation.

Overall, our study presents a preliminary yet promising step toward bridging think-aloud text data with behavioral strategies in a learning task. By demonstrating how verbalized thoughts align with computational models of learning, we highlight the value of think-aloud protocols in studying higher-level cognition.

References

- Barracough, D. J., Conroy, M. L., & Lee, D. (2004). Prefrontal cortex and decision making in a mixed-strategy game. *Nature neuroscience*, 7(4), 404–410.
- Guenouni, I., & Speekenbrink, M. (2021). Transfer of learned opponent models in repeated games. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 43).
- Johnson-Laird, P. N. (1983). *Mental models: Towards a cognitive science of language, inference, and consciousness* (No. 6). Harvard University Press.
- Momennejad, I., Russek, E. M., Cheong, J. H., Botvinick, M. M., Daw, N. D., & Gershman, S. J. (2017). The successor representation in human reinforcement learning. *Nature human behaviour*, 1(9), 680–692.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2023). Robust speech recognition via large-scale weak supervision. In *International conference on machine learning* (pp. 28492–28518).
- Rmus, M., Pan, T.-F., Xia, L., & Collins, A. G. (2024). Artificial neural networks for model identification and parameter estimation in computational cognitive models. *PLOS Computational Biology*, 20(5), e1012119.
- Simon, H. A., & Ericsson, K. A. (1984). Protocol analysis: Verbal reports as data.
- Trepka, E., Spitman, M., Bari, B. A., Costa, V. D., Cohen, J. Y., & Soltani, A. (2021). Entropy-based metrics for predicting choice behavior based on local response to reward. *Nature communications*, 12(1), 6567.

- Xie, H., Xiong, H., & Wilson, R. (2024a). Evaluating predictive performance and learning efficiency of large language models with think aloud in risky decision making. In *Conference for computational cognitive neuroscience*.
- Xie, H., Xiong, H., & Wilson, R. (2024b). From strategic narratives to code-like cognitive models: An llm-based approach in a sorting task. In *First conference on language modeling*.
- Xie, H., Xiong, H., & Wilson, R. C. (2023). Text2decision: Decoding latent variables in risky decision making from think aloud text. In *Neurips 2023 ai for science workshop*.