# Automated scientific minimization of regret

**Marcel Binz**
Helmholtz Munich
Institute for Human-Centered AI
marcel.binz@helmholtz-munich.de

**Akshay K. Jagadish**
Helmholtz Munich
Institute for Human-Centered AI
akshay.jagadish@helmholtz-munich.de

**Milena Rmus**
Helmholtz Munich
Institute for Human-Centered AI
milena.rmus@helmholtz-munich.de

**Eric Schulz**
Helmholtz Munich
Institute for Human-Centered AI
eric.schulz@helmholtz-munich.de

## Abstract

We introduce *automated scientific minimization of regret* (ASMR) – a framework for automated computational cognitive science. Building on the principles of scientific regret minimization, ASMR leverages Centaur – a recently proposed foundation model of human cognition – to identify gaps in an interpretable cognitive model. These gaps are then addressed through automated revisions generated by a language-based reasoning model. We demonstrate the utility of this approach in a multi-attribute decision-making task, showing that ASMR discovers cognitive models that predict human behavior at noise ceiling while retaining interpretability. Taken together, our results highlight the potential of ASMR to automate core components of the cognitive modeling pipeline.

## Introduction

The combination of large-scale behavioral data sets and advances in machine learning has enabled the development of highly predictive models of human behavior [Binz et al., 2024, Peterson et al., 2021, Eckstein et al., 2024]. Yet, while these models excel at predicting behavior, they offer limited insight into the underlying cognitive mechanisms. A central challenge, therefore, is how to move beyond prediction and leverage these models to improve our understanding of human cognition.

A promising framework for achieving this goal is *scientific regret minimization* [Agrawal et al., 2020]. It takes a black-box predictive model and compares it against an interpretable cognitive model on a per-data-point basis. This comparison is used to identify data points that are, in principle, predictable – because they are correctly predicted by the black-box model – but are not yet captured by the cognitive model. Patterns in these data points are analyzed and incorporated back into the cognitive model through an iterative refinement process.

Two factors have prevented the broader adoption of this framework: (1) it requires large data sets to train a predictive model, and (2) identifying patterns in the resulting data points can be challenging. The first challenge is typically addressed by conducting a large-scale study within the experimental paradigm of interest and training a black-box model – typically some form of neural network – on the resulting data. This approach imposes substantial overhead, thereby limiting the scope of scientific regret minimization. The second challenge remains largely unresolved and is typically addressed through what might be called the "method of staring" where researchers manually inspect the identified data points until a recognizable pattern is recognized that can be incorporated into the cognitive model.
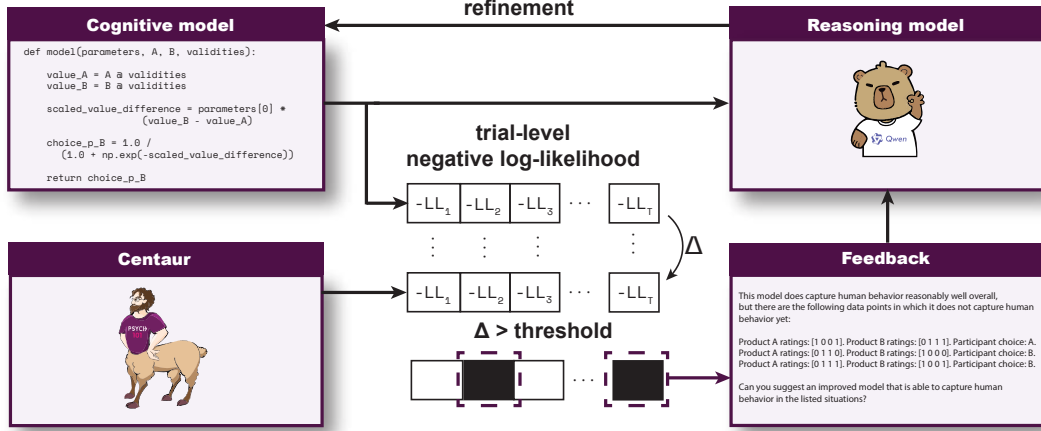
Figure 1: ASMR pipeline. Centaur is used as a guide to identify gaps in an interpretable cognitive model, resulting in a set of data points that are, in principle, predictable but are not currently accounted for. These points, together with the model's code and a brief instruction, are then provided to a language-based reasoning model. The reasoning model generates modifications to the cognitive model – a process which can be iterated multiple times.

The present paper introduces the idea of *automated scientific minimization of regret* (ASMR) – a framework that offers a solution to both of these issues. ASMR relies on Centaur – a foundation model of human cognition – as the predictive model [Binz et al., 2024]. Because Centaur was trained on a large collection of behavioral experiments, it can predict human behavior across domains without requiring additional data collection or task-specific training. To address the second issue, ASMR uses state-of-the-art reasoning models such as DeepSeek-R1 [DeepSeek-AI, 2025] or Qwen3 [Qwen, 2025]. Leveraging their language-based reasoning capabilities, these models can analyze failure modes and suggest new candidate models in an automated fashion and without the need for human intervention.

We present a case study illustrating the potential of ASMR. Without any human input, ASMR discovers cognitive models that match Centaur in predictive performance while retaining interpretability. Taken together, these results demonstrate the feasibility of fully automated scientific discovery [Musslick et al., 2024, Binz et al., 2025, Musslick et al., 2025, Castro et al., 2025, Rmus et al., 2025] guided by large-scale predictive models.

## Results

We apply ASMR to a multi-attribute decision-making paradigm [Hilbig and Moshagen, 2014], in which participants made repeated judgments between two fictitious products, each rated by four experts. Each expert provides a binary rating for both products, either approving or disapproving of them. Participants completed 96 trials in which they indicated which product they believed to be of higher quality. To guide their decision-making process, they were furthermore provided with the validity of each expert, defined as the probability that the product an expert approves is indeed the objectively better one on those occasions where the expert's two ratings differ. We focus our analysis on a subset of participants that are not part of Centaur's training data.

ASMR begins by extracting log-likelihoods of human responses for each participant and trial based on Centaur's predictions. In each iteration, it compares the resulting log-likelihoods to those of the cognitive model under evaluation. Free parameters of the cognitive model are fitted to the data using standard maximum likelihood estimation. Data points where the difference in log-likelihood exceeds a predefined threshold are identified and, along with the cognitive model code and a set of natural language instructions, passed to Qwen3-32B – a state-of-the-art reasoning model [Qwen, 2025]. Qwen3-32B then modifies the cognitive model with the goal of improving its predictions on the identified data points (see Figure 1). This process repeats for multiple iterations. In our experiments, five iterations were sufficient, though the optimal number may vary depending on task complexity.

**a** **b**



**c**

```python
def model(parameters, option_A, option_B):
    # define feature validity weights (importance of each feature)
    validities = np.array([parameters[0] * 0.9, 0.8, 0.7, 0.6])

    # compute weighted value for each option
    value_A = option_A @ validities
    value_B = option_B @ validities

    # compute scaled difference in value
    scale_value_difference = parameters[1] * (value_B - value_A)

    # apply logistic function to obtain choice probabilities
    choice_probability_B = 1.0 / (1.0 +
        np.exp(-scale_value_difference))

    return choice_probability_B
```
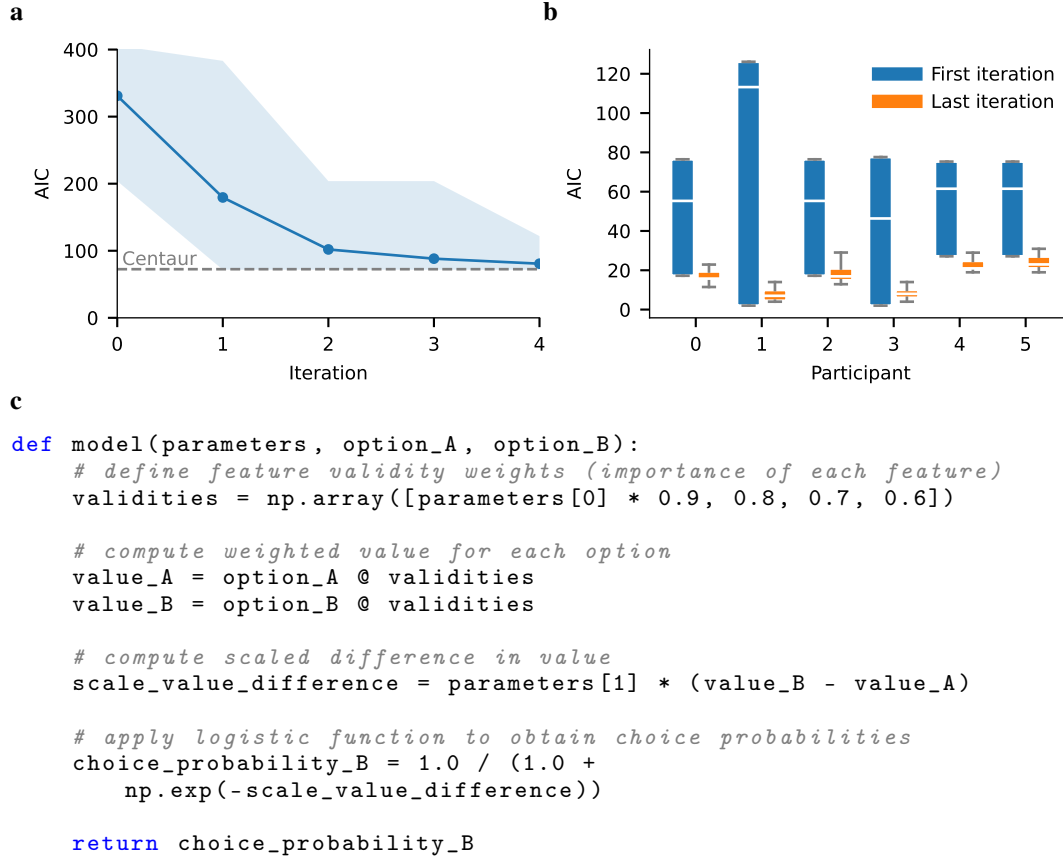
Figure 2: Summary of results. **a**, Improvement of aggregated AIC scores across iterations of ASMR. The solid line show the average AIC score, while the shaded areas represent the worst and best model at a given iteration. **b**, AIC scores from the models at the first and last iteration for each individual participant. **c**, Python code for one of the discovered models with the lowest AIC score.

We initialize the cognitive model using three strategies frequently employed in decision-making literature: the take-the-best heuristic, an equal weighting heuristic, and a weighted additive strategy [Gigerenzer et al., 2000, Binz et al., 2022]. Reported results are averaged over ten simulations per model class and final model performance is evaluated using the Akaike information criterion (AIC).

We found that ASMR rapidly improves the initial model. Within five iterations, the discovered models reach an average AIC score (M = 80.72, SD = 17.05) that approaches Centaur's goodness-of-fit ($AIC_{Centaur} = 72.5$; see Figure 2a). Notably, the best-performing model ($AIC_{ASMR} = 71.73$) even surpasses Centaur and matches the performance of prior approaches [Rmus et al., 2025]. Furthermore, ASMR consistently improves the AIC score of every participant, as shown in Figure 2b. Figure 2c presents one of the discovered models with the lowest AIC score. This example illustrates an adaptive upweighting of the highest-validity expert, enabling interpolation between take-the-best and weighted-additive strategies [Parpart et al., 2018] – a pattern that consistently emerges across multiple discovered models.

## Discussion

We have presented ASMR as a framework for the automated discovery of interpretable cognitive models. ASMR invokes Centaur – a foundation model of human cognition – as a reference model to reveal data points that are, in principle, predictable but are not captured by a given cognitive model. A state-of-the-art reasoning model is then prompted with this information and tasked with revising the cognitive model. When applied iteratively, this process enables the discovery of novel, interpretable cognitive theories.

While the results presented in this article serve as a first proof-of-concept, they open up several avenues for future exploration. A central question for future work concerns the nature of the feedback signal. We have shown that simply presenting negative data points is sufficient for discovering novel cognitive models. However, in more complex domains, it may be necessary to supplement this with additional forms of feedback, such as positive examples, or performance metrics. In our simulations, the reasoning model was furthermore only given the most recent version of the cognitive model for revision. Improvements in context window size may eventually lift this constraint and allow reasoning models to process and revise multiple proposals within a single prompt. Fully evaluating these alternatives necessitates a benchmark for automated scientific discovery in the cognitive sciences. Ideally, such a benchmark should span multiple experimental paradigms and cover a representative spectrum of cognitive modeling scenarios.

Ultimately, we expect that ASMR will enable the rapid generation of a wide range of cognitive models. To manage this growing set, we envision a searchable database that systematically links cognitive models to the experimental data they aim to explain. This database could be queried to identify models relevant to specific tasks or populations and would provide a foundation for large-scale meta-analyses across studies, paradigms, and modeling approaches.

## Methods

We initialize ASMR with the following cognitive models:

```python
# weighted addtive strategy
NUM_PARAMETERS = 1

def model(parameters, option_A, option_B):
    """
    Compute the probability of choosing Option B over Option A.

    Parameters
    ----------
    parameters : np.ndarray of shape (num_parameters,)
        Model parameters.

    option_A : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option A across trials.

    option_B : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option B across trials.

    Returns
    -------
    choice_probability_B : np.ndarray of shape (num_trials,)
        The predicted probability of choosing Option B on each trial.
    """

    # define feature validity weights (importance of each feature)
    validities = np.array([0.9, 0.8, 0.7, 0.6])

    # compute weighted value for each option
    value_A = option_A @ validities
    value_B = option_B @ validities

    # compute scaled difference in value
    scale_value_difference = parameters[0] * (value_B - value_A)

    # apply logistic function to obtain choice probabilities
    choice_probability_B = 1.0 / (1.0 +
        np.exp(-scale_value_difference))

    # clip probabilities to avoid numerical issues
    choice_probability_B = np.clip(choice_probability_B, 0.00001, 1 -
        0.00001)

    return choice_probability_B

# take-the-best heuristic
NUM_PARAMETERS = 1

def model(parameters, option_A, option_B):
    """
    Compute the probability of choosing Option B over Option A.

    Parameters
    ----------
    parameters : np.ndarray of shape (num_parameters,)
        Model parameters.

    option_A : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option A across trials.

    option_B : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option B across trials.
```

```
    Returns
    -------
    choice_probability_B : np.ndarray of shape (num_trials,)
        The predicted probability of choosing Option B on each trial.
    """

    # define feature validity weights (importance of each feature)
    validities = np.array([1.0, 0.5, 0.25, 0.125])

    # compute weighted value for each option
    value_A = option_A @ validities
    value_B = option_B @ validities

    # compute scaled difference in value
    scale_value_difference = parameters[0] * (value_B - value_A)

    # apply logistic function to obtain choice probabilities
    choice_probability_B = 1.0 / (1.0 +
        np.exp(-scale_value_difference))

    # clip probabilities to avoid numerical issues
    choice_probability_B = np.clip(choice_probability_B, 0.00001, 1 -
        0.00001)

    return choice_probability_B

# equal weighting heuristic
NUM_PARAMETERS = 1

def model(parameters, option_A, option_B):
    """
    Compute the probability of choosing Option B over Option A.

    Parameters
    ----------
    parameters : np.ndarray of shape (num_parameters,)
        Model parameters.

    option_A : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option A across trials.

    option_B : np.ndarray of shape (num_trials, num_features)
        Feature matrix for Option B across trials.

    Returns
    -------
    choice_probability_B : np.ndarray of shape (num_trials,)
        The predicted probability of choosing Option B on each trial.
    """

    # compute weighted value for each option
    value_A = option_A.sum(-1)
    value_B = option_B.sum(-1)

    # compute scaled difference in value
    scale_value_difference = parameters[0] * (value_B - value_A)

    # apply logistic function to obtain choice probabilities
    choice_probability_B = 1.0 / (1.0 +
        np.exp(-scale_value_difference))

    # clip probabilities to avoid numerical issues
    choice_probability_B = np.clip(choice_probability_B, 0.00001, 1 -
        0.00001)
```

```
        return choice_probability_B
```

In each iteration, the cognitive model is fitted to human responses on a per-subject basis using maximum likelihood estimation. This procedure is implemented using SCIPY's minimize function with the BFGS algorithm. We then subtract the cached negative log-likelihoods obtained from Centaur from those of the fitted cognitive model. Next, we submit all data points where this difference exceeds a threshold of $\Delta \geq 0.05$ – together with the model's code and a brief instruction – to Qwen3-32B [Qwen, 2025]. For this, we rely on the four-bit quantized version of the UNSLOTH package with sampling parameters set to the recommended values [Qwen, 2025]. We use the following prompt template:

```
I am studying human behavior in a multi-attribute decision-making experiment.
In this experiment, participants encounter a number of trials, in which they
have to choose between two options labelled A and B.
These options are fictitious products that are each characterized by four
features.
Each feature corresponds to a binary rating of an expert, either approving of
the product (1) or not (0).
The four experts are ordered based on their validity (taking values of 90%, 80%,
70%, and 60%), with the first feature corresponding to the ratings from the
highest validity expert.
In each trial, people have to predict which of the shown options is superior in
terms of quality based on the presented information.
I have the following computational model that is currently my best guess for how
people make decisions in this experiment:

[INSERT MODEL CODE]

This model does capture human behavior reasonably well overall, but there are
the following data points in which it does not capture human behavior yet:

[INSERT DATA POINTS]

Can you suggest an improved model that is able to capture human behavior in the
listed situations?

Please structure your answer as follows:
* Keep the structure of the function exactly the same.  * Do not change the
docstring.
* State the number of free parameters before the model function using the
NUM_PARAMETERS variable.
* Do not write any text besides that and do not elaborate any further.
```

Finally, the cognitive model is updated based on Qwen3's output, and the entire process is repeated. We ran ten simulations per model class, each with five iterations, and report results averaged across simulations unless otherwise noted. ASMR – in its current form – requires access to a single 80GB GPU (e.g., A100) but could be modified to run on fewer computational resources by using smaller models.

# References

Marcel Binz, Elif Akata, Matthias Bethge, Franziska Brändle, Fred Callaway, Julian Coda-Forno, Peter Dayan, Can Demircan, Maria K Eckstein, Noémi Éltető, et al. Centaur: a foundation model of human cognition. *arXiv preprint arXiv:2410.20268*, 2024.

Joshua C Peterson, David D Bourgin, Mayank Agrawal, Daniel Reichman, and Thomas L Griffiths. Using large-scale experiments and machine learning to discover theories of human decision-making. *Science*, 372(6547):1209–1214, 2021.

Maria K Eckstein, Christopher Summerfield, Nathaniel Daw, and Kevin J Miller. Hybrid neural-cognitive models reveal how memory shapes human reward learning, 2024.

Mayank Agrawal, Joshua C Peterson, and Thomas L Griffiths. Scaling up psychology via scientific regret minimization. *Proceedings of the National Academy of Sciences*, 117(16):8825–8835, 2020.

DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL `https://arxiv.org/abs/2501.12948`.

Qwen. Qwen3 technical report. `https://github.com/QwenLM/Qwen3/blob/main/Qwen3_Technical_Report.pdf`, 2025.

Sebastian Musslick, Younes Strittmatter, and Marina Dubova. Closed-loop scientific discovery in the behavioral sciences. *PsyArXiv*, 10, 2024.

Marcel Binz, Stephan Alaniz, Adina Roskies, Balazs Aczel, Carl T Bergstrom, Colin Allen, Daniel Schad, Dirk Wulff, Jevin D West, Qiong Zhang, et al. How should the advancement of large language models affect the practice of science? *Proceedings of the National Academy of Sciences*, 122(5):e2401227121, 2025.

Sebastian Musslick, Laura K Bartlett, Suyog H Chandramouli, Marina Dubova, Fernand Gobet, Thomas L Griffiths, Jessica Hullman, Ross D King, J Nathan Kutz, Christopher G Lucas, et al. Automating the practice of science: Opportunities, challenges, and implications. *Proceedings of the National Academy of Sciences*, 122(5):e2401238121, 2025.

Pablo Samuel Castro, Nenad Tomasev, Ankit Anand, Navodita Sharma, Rishika Mohanta, Aparna Dev, Kuba Perlin, Siddhant Jain, Kyle Levin, Noémi Éltető, et al. Discovering symbolic cognitive models from human and animal behavior. *bioRxiv*, pages 2025–02, 2025.

Milena Rmus, Akshay K Jagadish, Marvin Mathony, Tobias Ludwig, and Eric Schulz. Towards automation of cognitive modeling using large language models. *arXiv preprint arXiv:2502.00879*, 2025.

Benjamin E Hilbig and Morten Moshagen. Generalized outcome-based strategy classification: Comparing deterministic and probabilistic choice models. *Psychonomic bulletin & review*, 21:1431–1443, 2014.

Gerd Gigerenzer, Peter M Todd, the ABC Research Group, et al. *Simple heuristics that make us smart*. Oxford University Press, 2000.

Marcel Binz, Samuel J Gershman, Eric Schulz, and Dominik Endres. Heuristics from bounded meta-learned inference. *Psychological review*, 129(5):1042, 2022.

Paula Parpart, Matt Jones, and Bradley C Love. Heuristics as bayesian inference under extreme priors. *Cognitive psychology*, 102:127–144, 2018.