



# Human-like Category Learning by Injecting Ecological Priors from Large Language Models into Neural Networks

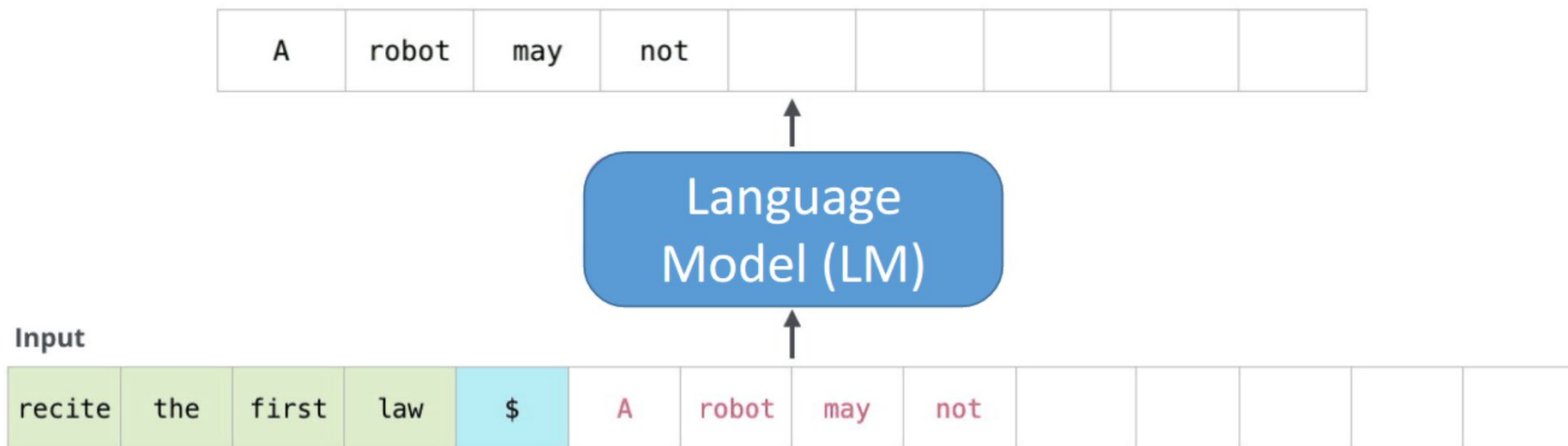
Akshay K Jagadish, Julian Coda-Forno, Mirko Thalmann,  
Eric Schulz\*, & Marcel Binz\*

# Ecological rationality



how do we define ecologically valid tasks?

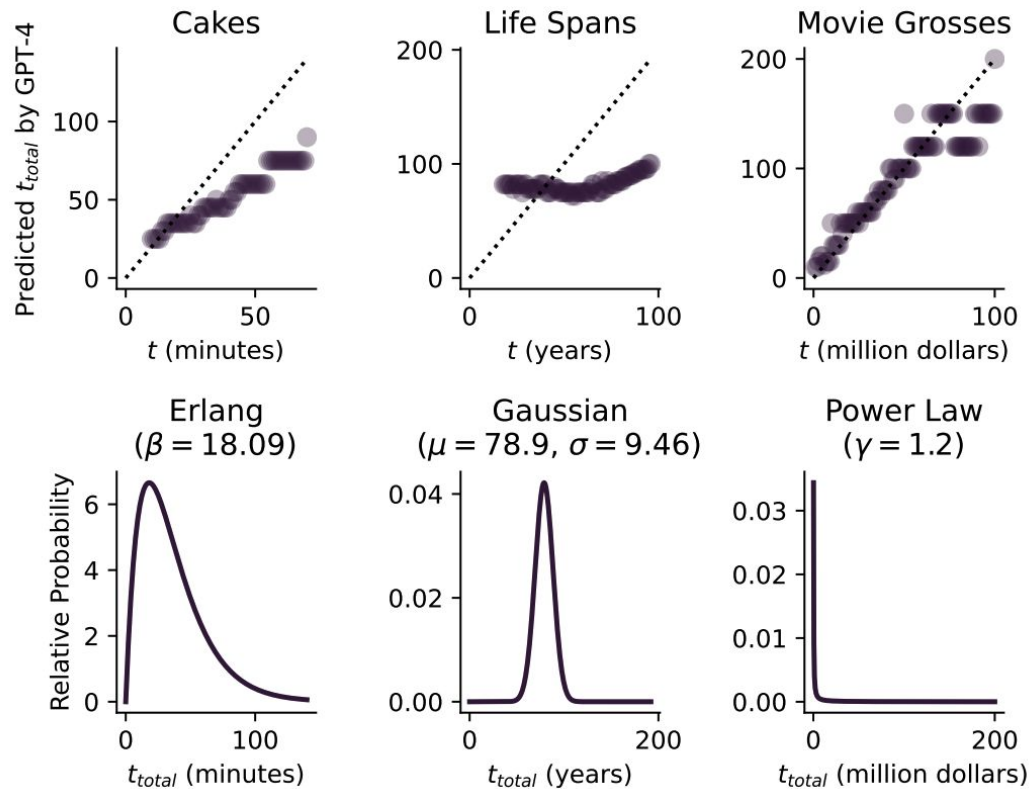
# Large Language Models (LLMs)



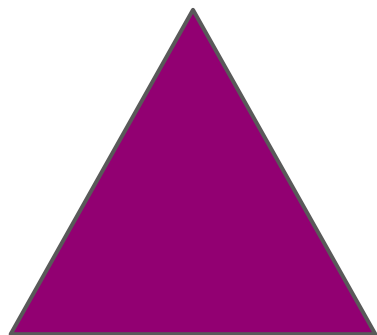
“LLMs [...] are powerful new cultural technologies, analogous to earlier technologies like writing, print, libraries, internet search and even language itself.”

Yiu, ..., Gopnik, 2023

# LLMs internalize everyday distributions similar to humans



Can LLMs be used to generate ecologically  
valid tasks?

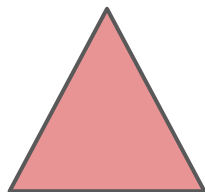




Choose A or B?

Choice: A

Correct Choice!

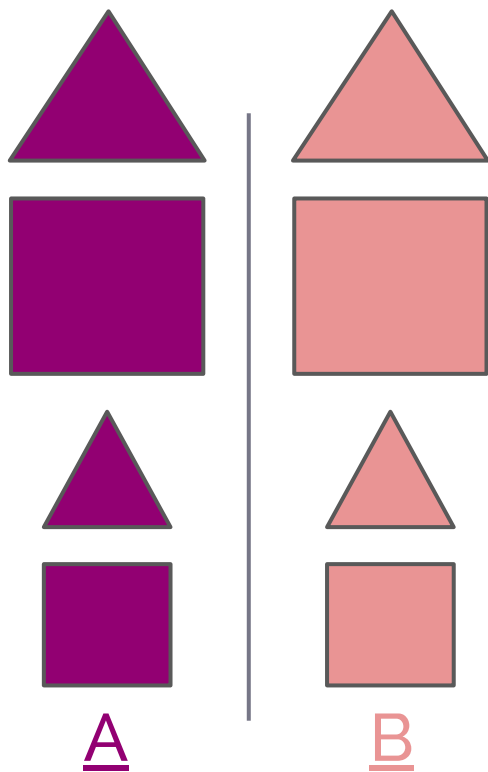


Choose A or B?

Choice: A

Wrong Choice!

# Category learning



Canonical task in psychology (Shepard et al. 1961)

Effects replicated successfully (Nosofsky et al. 1986, Medin et al. 1978, Badham et al. 2017, Devraj et al. 2022)

Many models of human categorisation exist (Smith et al. 1998, Anderson et al. 1991, Sandborn et al. 2010, Griffiths et al. 2017)



How can we generate category learning tasks using LLMs?

## Two step approach

Step 1: Synthesize feature names and category labels

Step 2: Generate data for category learning tasks

## Quick side note on LLMs considered

Achieved varying degrees of success with

- LLaMA
- GPT-3.5
- GPT-4

Results reported for **Claude-2**

- process up to 100k tokens
- instruction tuned
- good performance on preliminary tests
- temperature set to 1.

## Two step approach

Step 1: Synthesize feature names and category labels

Step 2: Generate data for category learning tasks

## Prompt

I am a psychologist who wants to run a category learning experiment. In a category learning experiment, there are many different {three}-dimensional stimuli, each of which belongs to one of two possible real-world categories.

Please generate names for {three} stimulus feature dimensions and {two} corresponding categories for {250} different category learning experiments:

- feature dimension 1, feature dimension 2, ..., feature dimension {3}, category label 1, category label 2

Feature 1	Feature 2	Feature 3	Category 1	Category 2
sodium	fat	protein	healthy	unhealthy
price	comfort	mileage	economy car	luxury car
rhythm	melody	harmony	jazz	classical

## Two step approach

Step 1: Synthesize feature names and category labels

Step 2: Generate data for category learning tasks

## Prompt

I am a psychologist who wants to run a category learning experiment. For a category learning experiment, I need a list of stimuli and their category labels. Each stimulus is characterized by {three} distinct features: {sodium}, {fat}, and {protein}. These features can take only numerical values. The category label can take the values {healthy} or {unhealthy} and should be predictable from the feature values of the stimulus.

Please generate a list of {400} stimuli with their feature values and their corresponding category labels using the following template for each row:

- feature value 1, feature value 2,..., feature value {3}, category label



## An example task

Sodium	Fat	Protein	Category
250	15	20	healthy
220	17	22	unhealthy
320	26	31	unhealthy
145	11	20	healthy

Is there an emerging theme in the LLM generated tasks?

A word cloud featuring a variety of nouns and adjectives. The words are arranged in a dense, overlapping manner. The colors used include shades of purple, blue, yellow, orange, red, and green. The sizes of the words vary, with 'classical' and 'trucks' being the largest. The words are: mammals, kitchen knives, vodka, electric, english, scotch, nonmetals, classical, dramas, electric guitars, suvs, cabernets, sauvignon, laptops, roses, conifers, ales, reptiles, trucks, lagers, comedy, lions, movies, tropical, insects, pinot noir, birds, sedans, europe, vegetables, diamonds, running shoes, stouts, metals, lattes, fish, ipas, sportscars, rye, rock, desktops, acoustic, chardonnay, jazz, pop, drama, tigers, fruits, comedies, acoustic guitars.

Do these tasks capture real-world statistics?

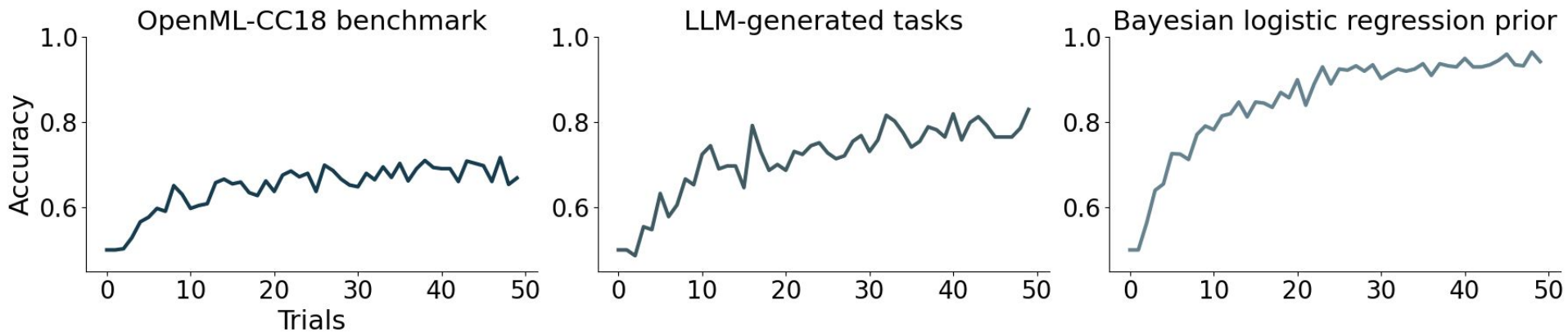
Comparison:

- OpenML-CC18 classification benchmark (Bischl et al. 2019)
- Bayesian logistic regression prior (Binz et al. 2022, Speekenbrink et al. 2008, 2010)
- Bayesian neural network prior (Müller et al. 2022, Levering et al. 2020)

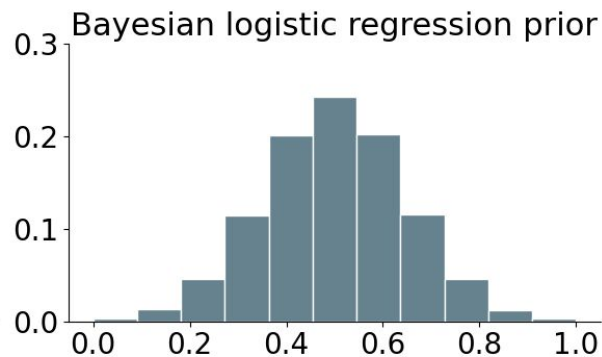
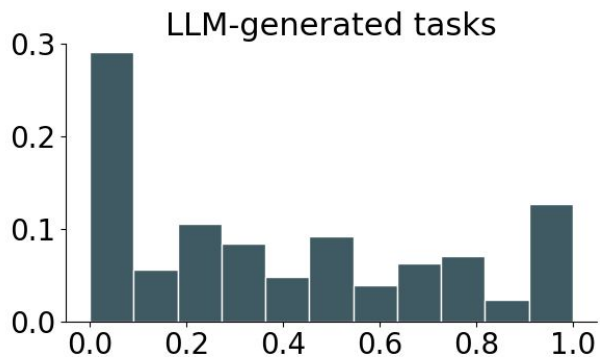
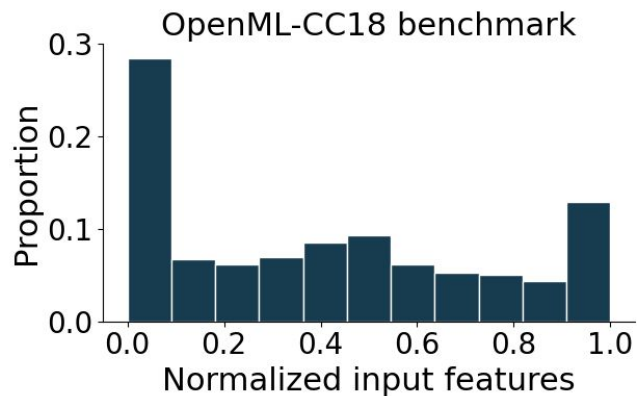
Statistics:

1. Classification performance
2. Normalized input features
3. Pair-wise input feature correlations
4. Sparsity
5. Linearity

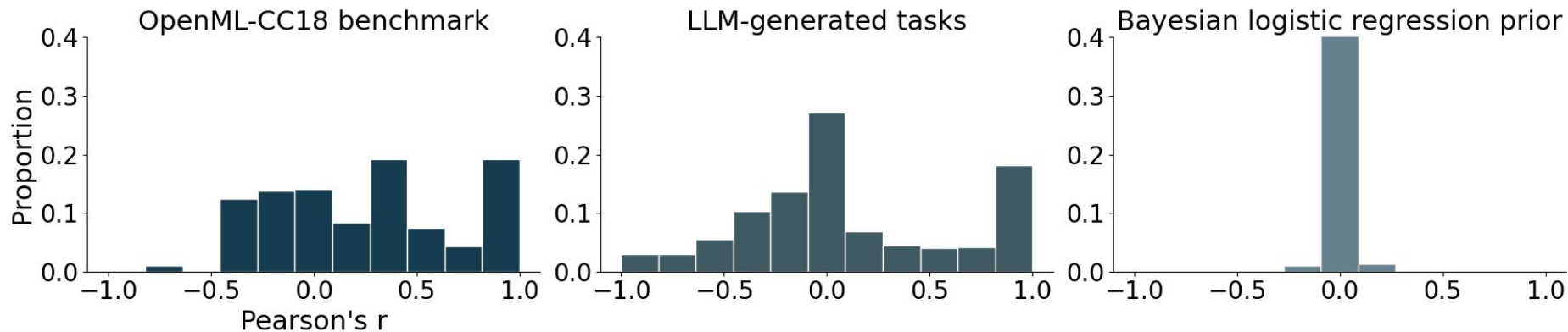
## Learning curves of LLMs match real-world data



## Distribution of normalized input features match real-world data

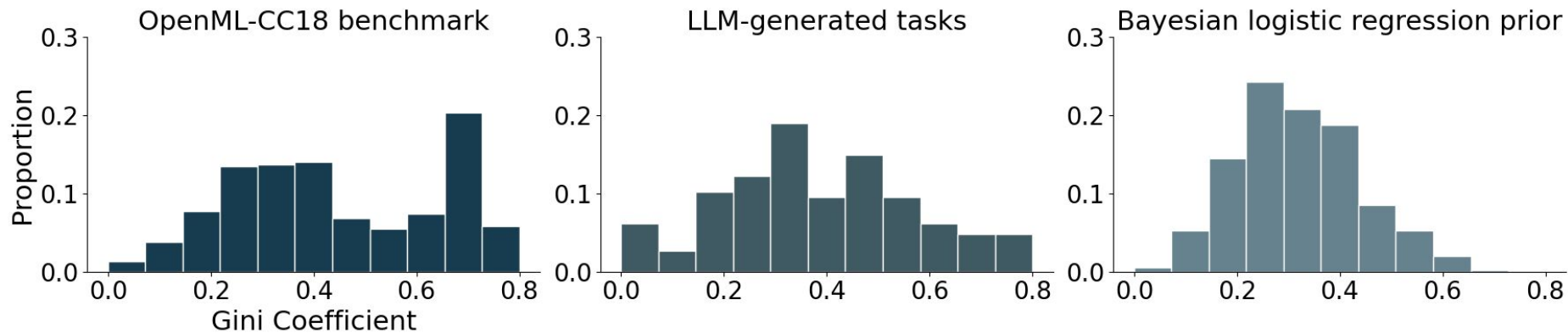


## Distribution of input feature correlations match real-world data

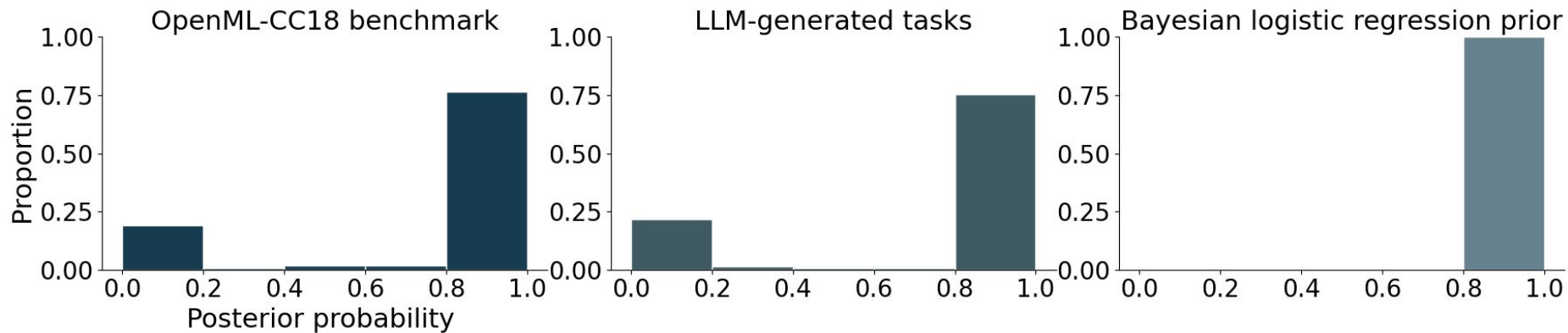




## Distribution of sparsity match real-world data



## Proportion of linear tasks match real-world data



1. Two step approach to generate category learning tasks from LLMs
2. Statistics of tasks generated by LLMs match real-world classification tasks

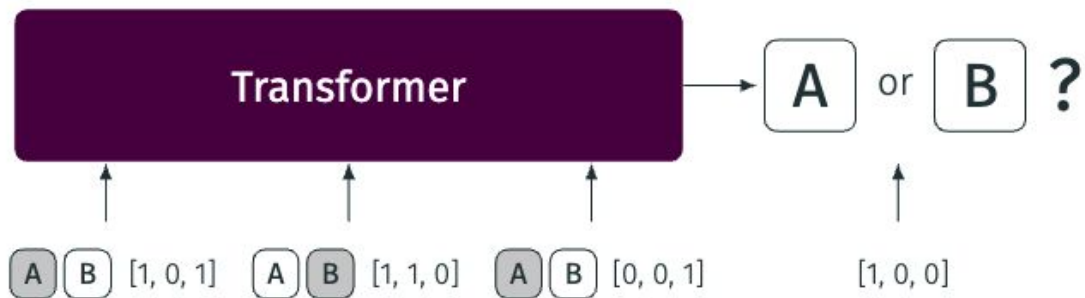
## Interim summary



Derving an ecologically rational model

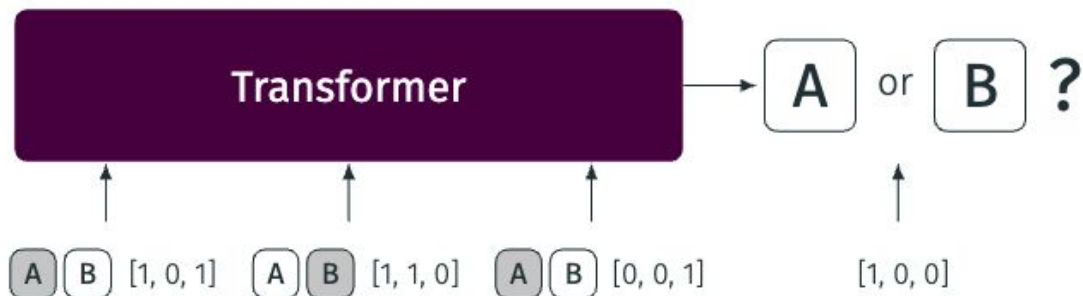
## In-weights learning (or meta-learning)

### In-context learning



## In-weights learning (or meta-learning)

### In-context learning



How much of human behavior does ERM  
explain?

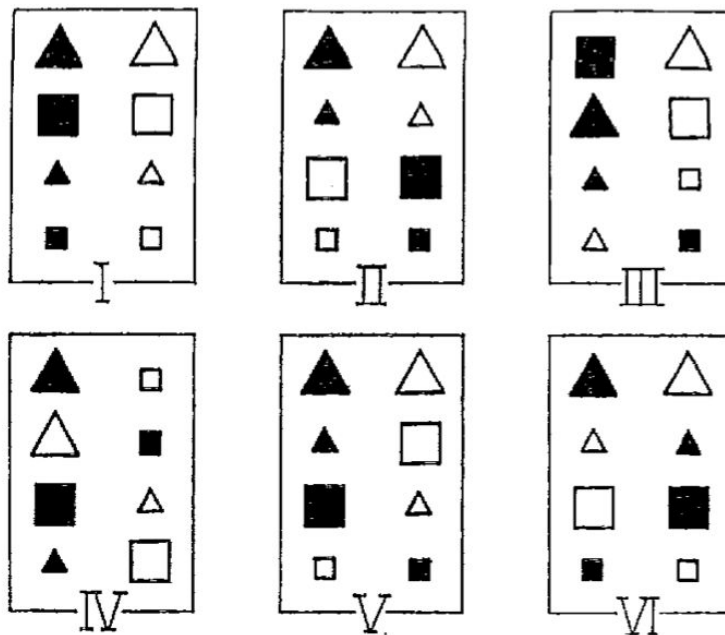


We looked at three different effects in human category learning:

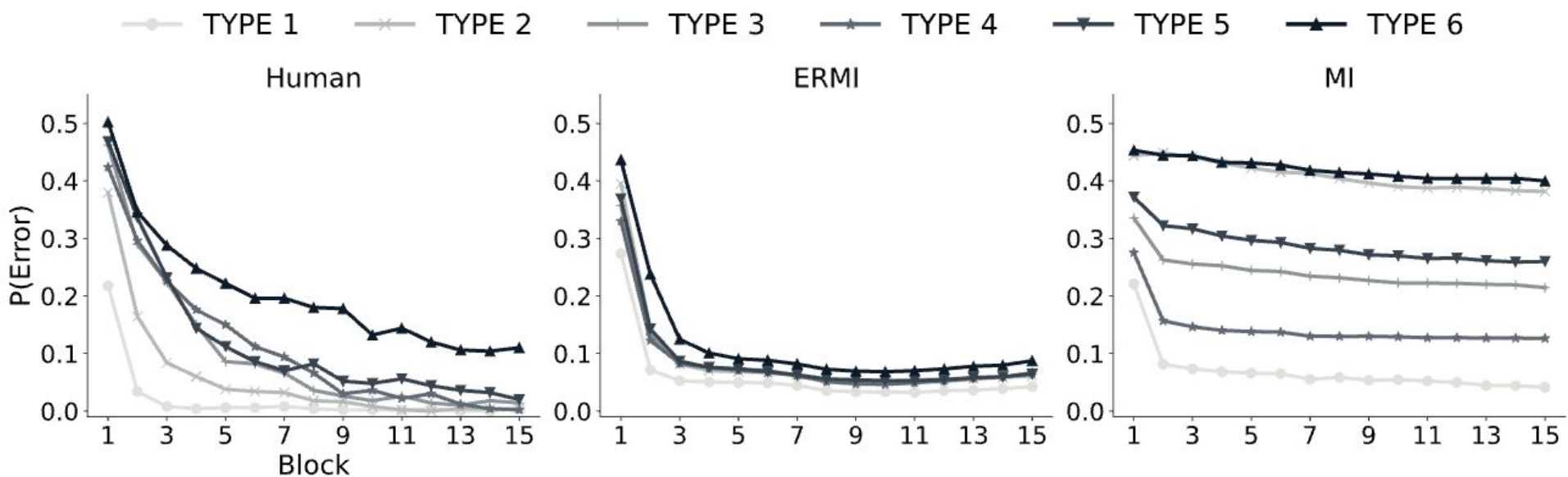
1. Learning difficulties
2. Learning strategy
3. Generalisation

Can ERMl explain human learning  
difficulties?

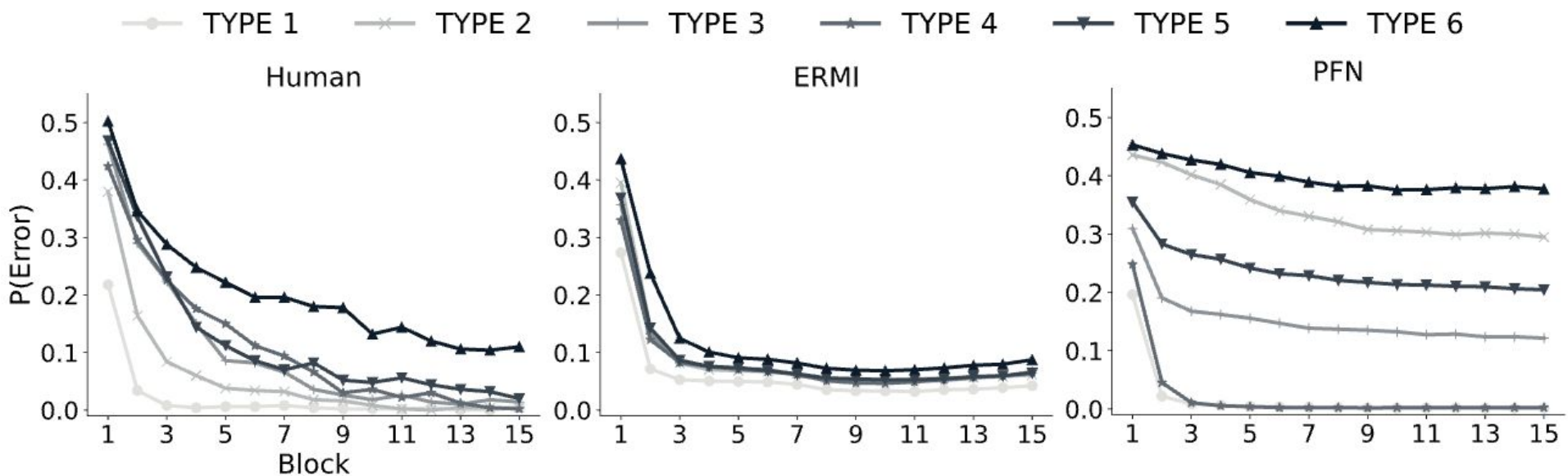
## Task difficulty increases from Type I to Type VI



## ERMI shows human-like learning difficulties



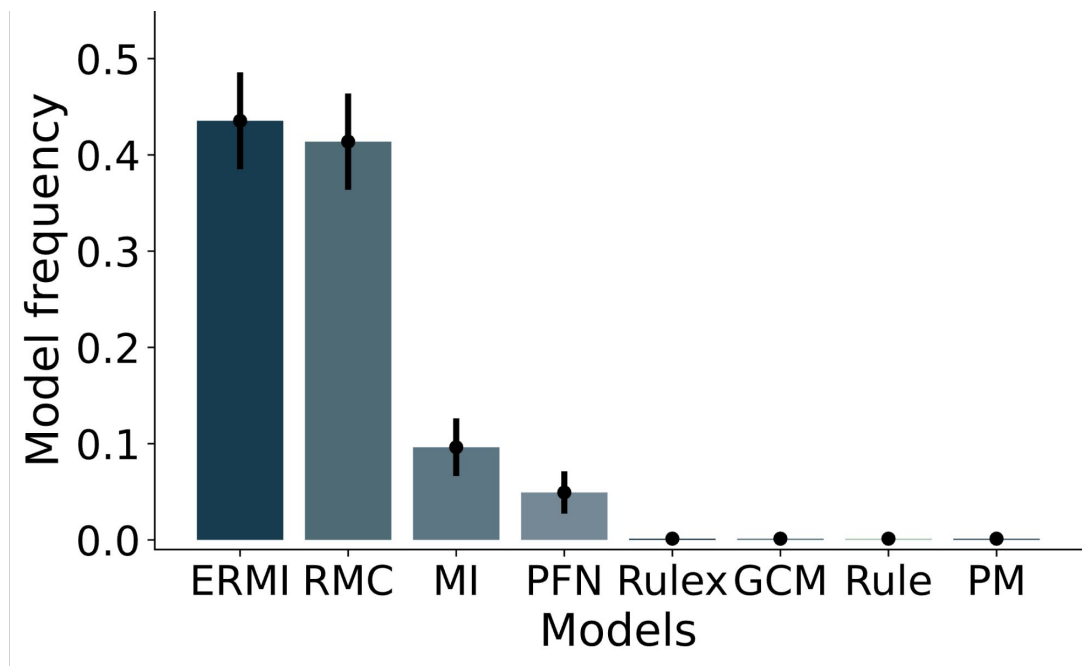
## ERMI shows human-like learning difficulties



## ERMI explains human data from Badham et al. 2017 better than cognitive models

1. **MI:** Meta-learned Inference on Bayesian logistic regression prior (Binz et al. 2022, Speekenbrink et al. 2008, 2010)
2. **PFN:** Meta-learned inference on Bayesian neural network prior (Müller et al. 2022, Levering et al. 2020)
3. **RMC:** Rational model of categorisation (Anderson et al. 1991)
4. **GCM:** Generalized context model (Nosofsky, 1986)
5. **PM:** Prototype Model (Homa and Cultice, 1984)
6. **Rule:** Rule-based model (Ashby and Townsend, 1986)
7. **Rulex:** Rule plus exception model (Nosofsky, 1994)

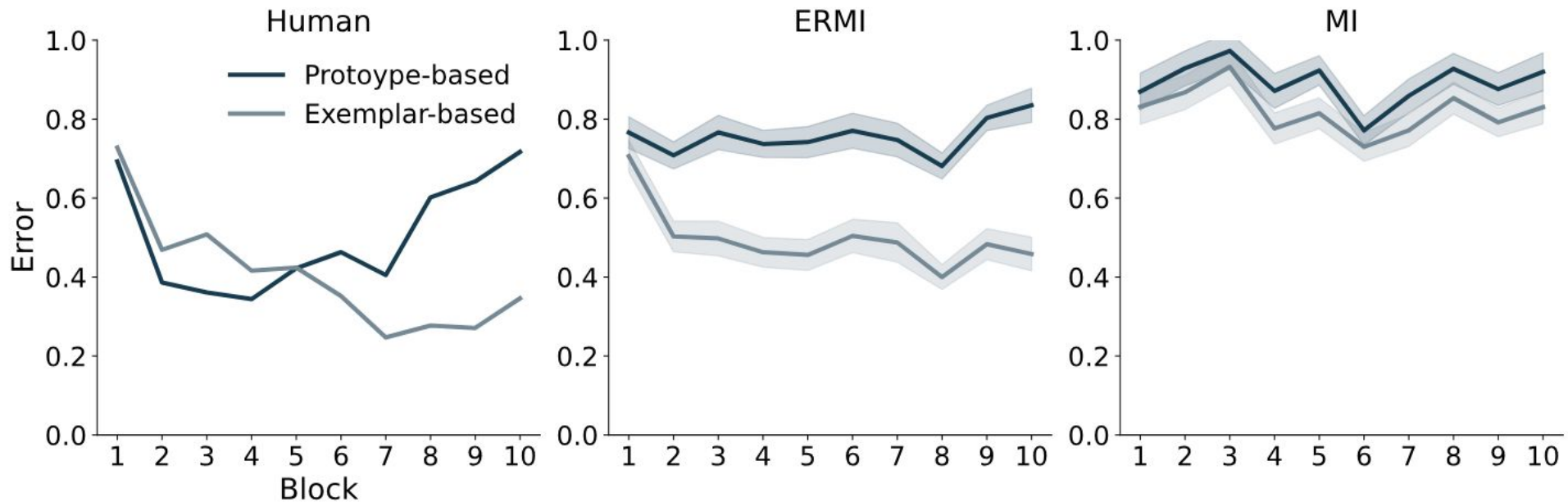
## ERMI explains human data from Badham et al. 2017 better than cognitive models



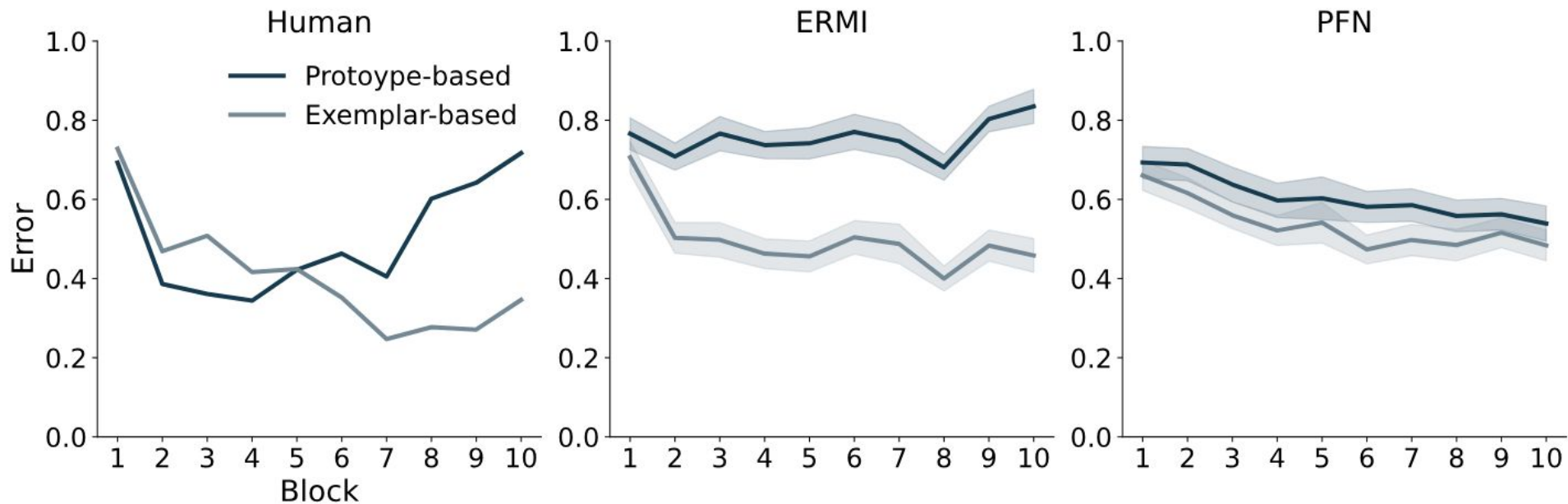
Does ERMI shift to exemplar-based learning strategy?



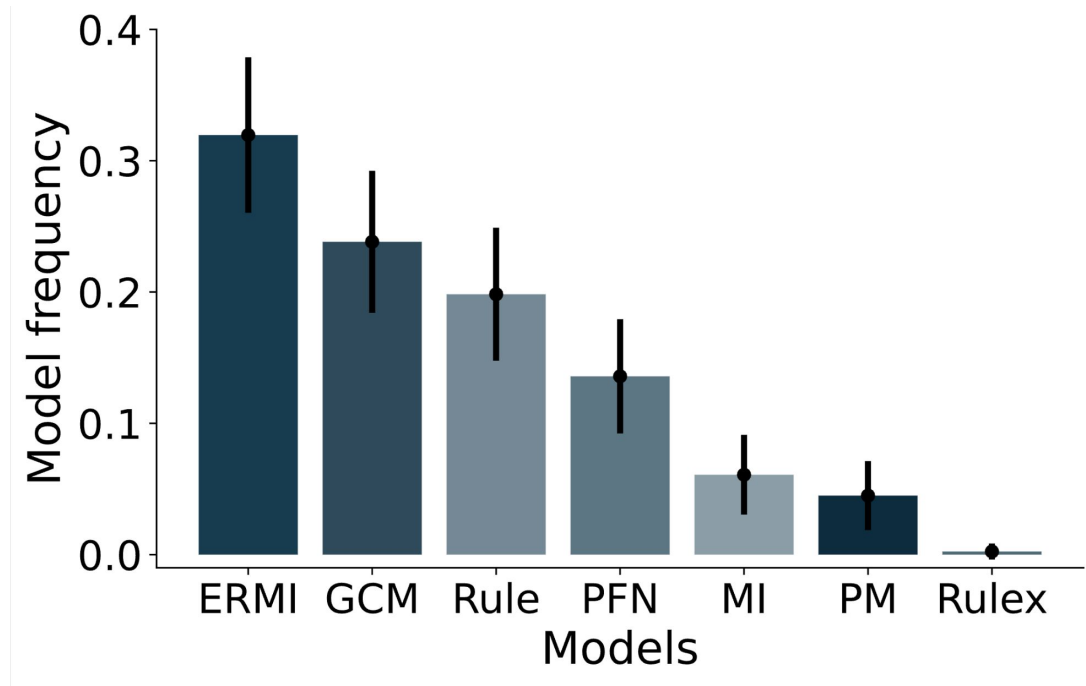
## ERMI also becomes more exemplar-based with learning



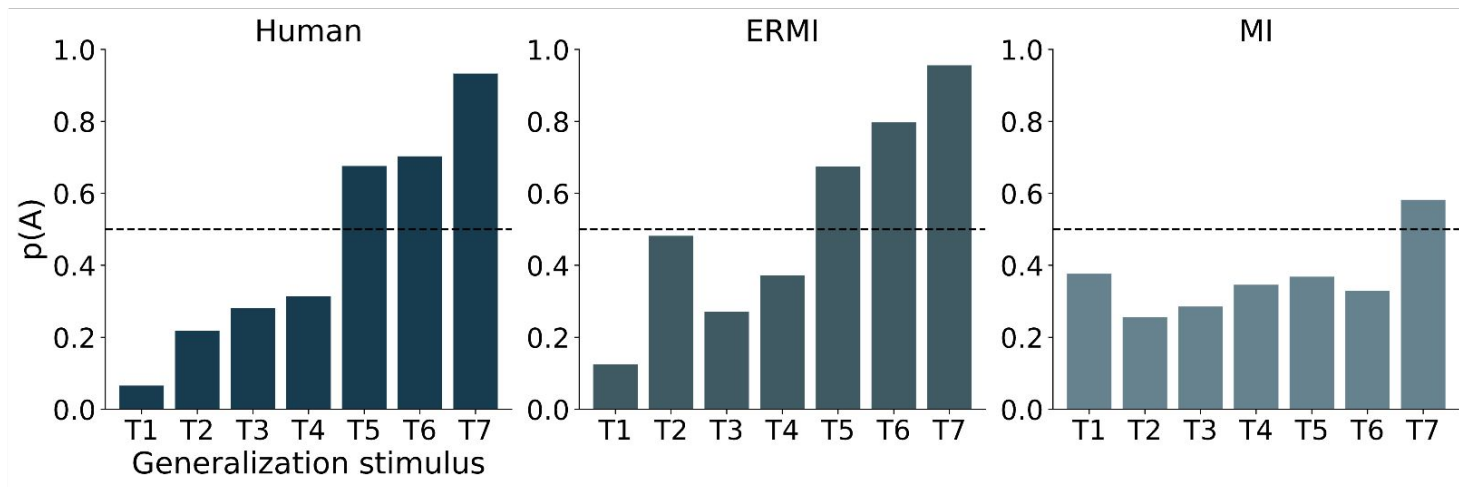
## ERMI also becomes more exemplar-based with learning



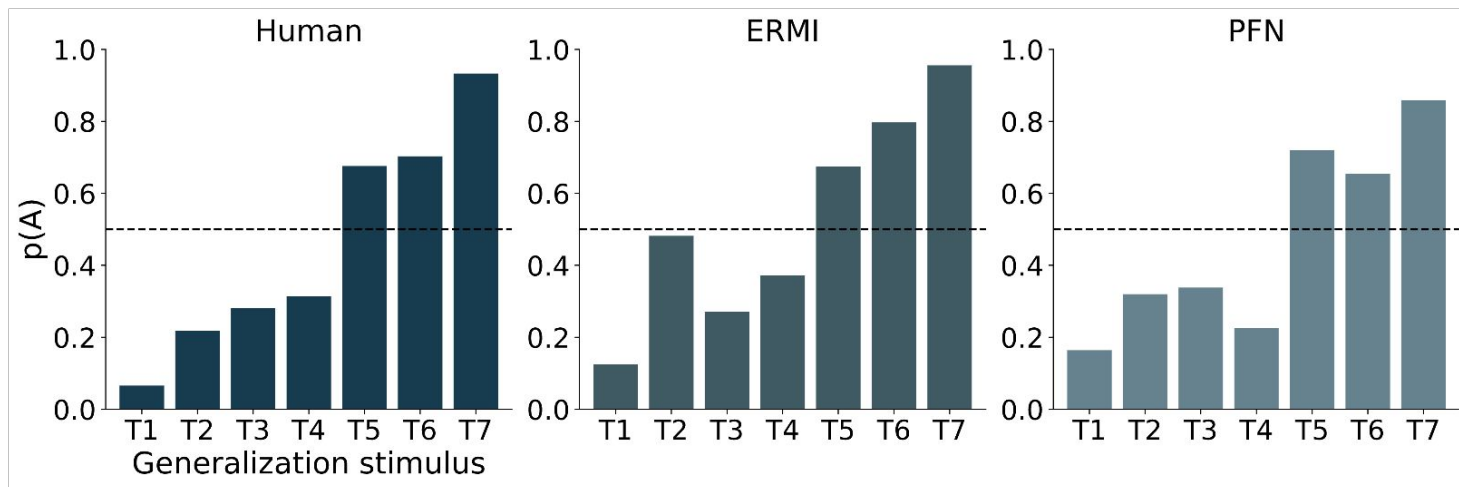
## ERMI explains human data from Devraj et al. 2022 better than cognitive models



Does ERMI generalize to unseen stimuli like humans?



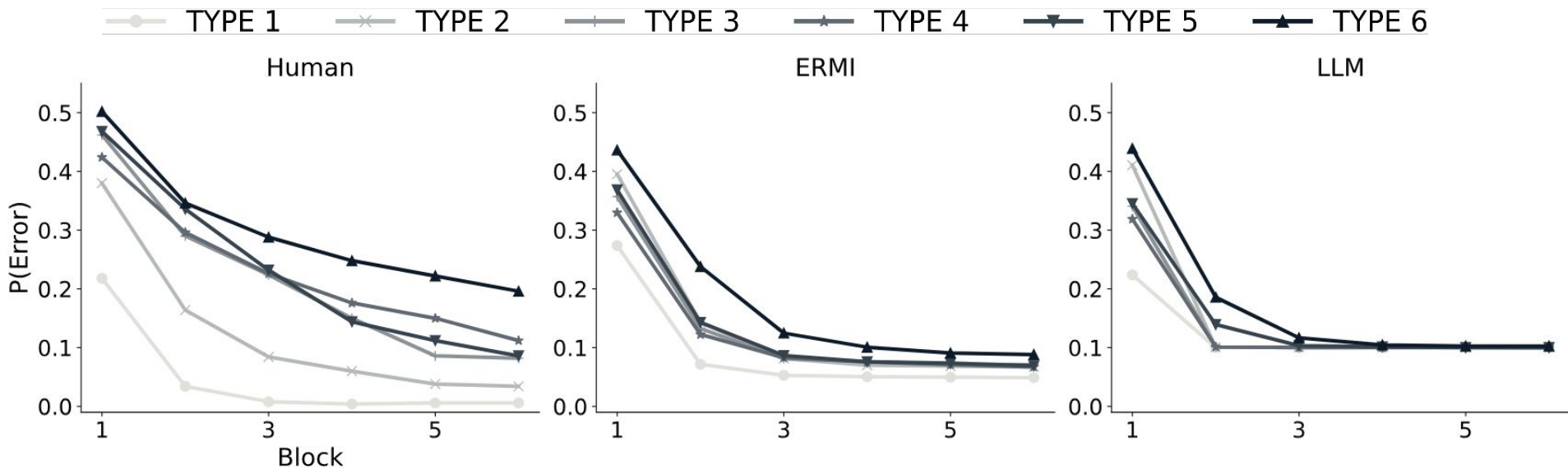
T1: [1 0 1 1]  
 T2: [1 0 1 0]  
 T3: [0 1 1 1]  
 T4: [1 1 0 1]  
 T5: [1 1 0 0]  
 T6: [0 1 1 0]  
 T7: [0 0 0 0]



T1: [1 0 1 1]  
 T2: [1 0 1 0]  
 T3: [0 1 1 1]  
 T4: [1 1 0 1]  
 T5: [1 1 0 0]  
 T6: [0 1 1 0]  
 T7: [0 0 0 0]

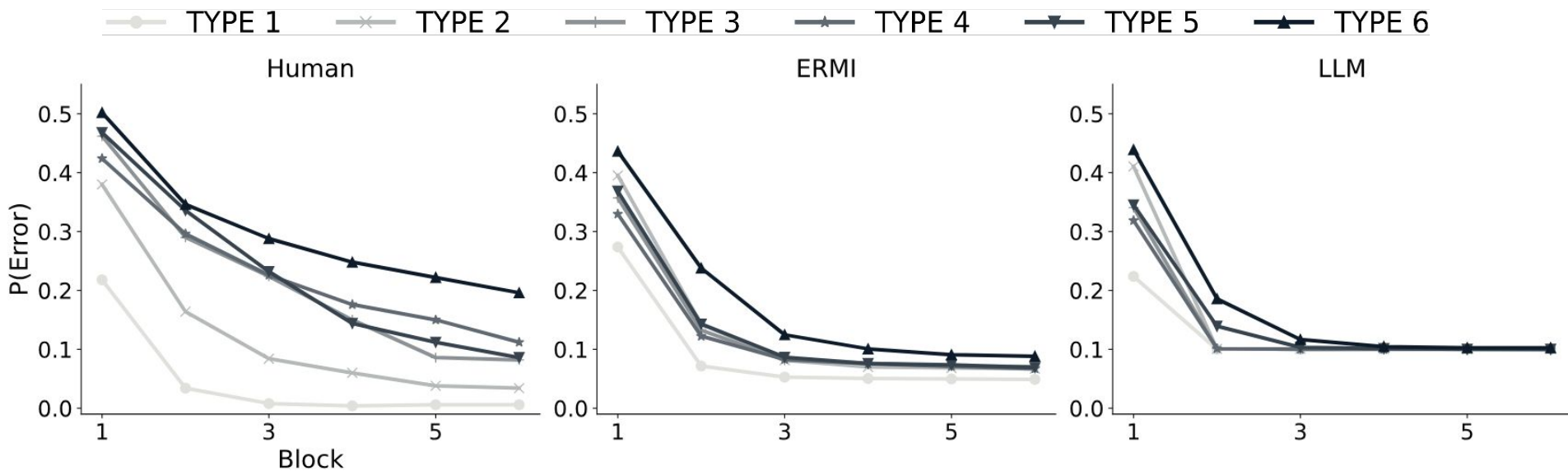
What if we directly let an LLM perform a category learning task?

# ERMI is better fit to humans both qualitatively and quantitatively in Shepard et al. 1961



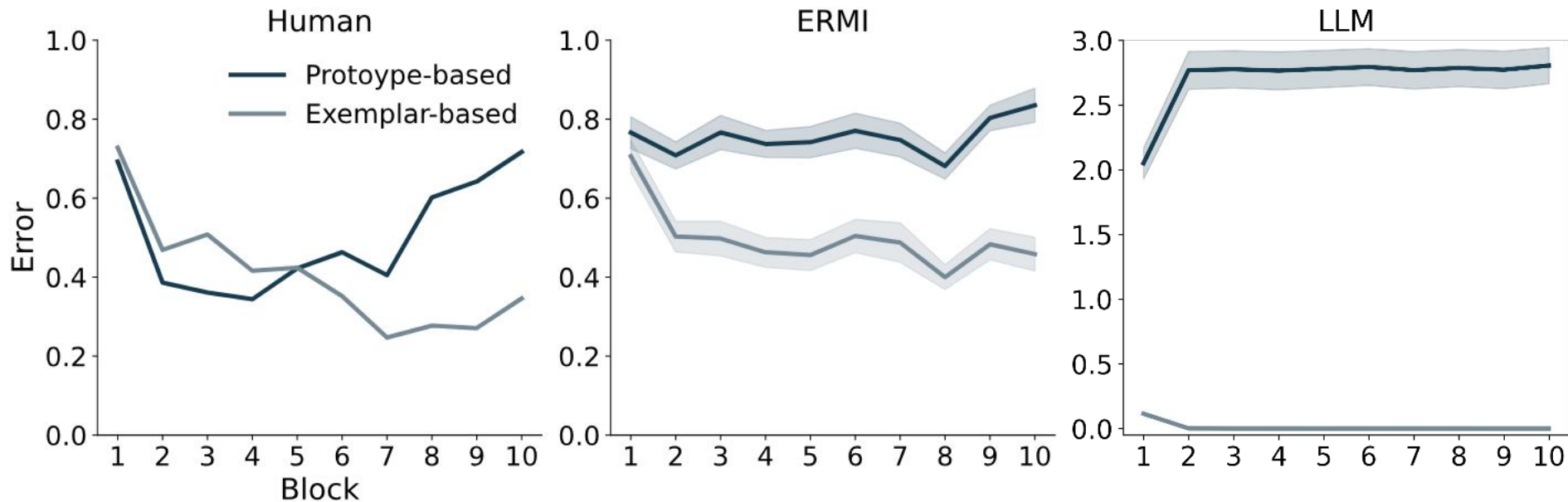


## ERMI is better fit to humans both qualitatively and quantitatively in Shepard et al. 1961



ERMI also offers a better fit to humans than LLM in term of BIC

## ERMI is better fit to humans both qualitatively and quantitatively to Smith and Minda, 1998



ERMI also offers a better fit to humans than LLM in term of BIC

Does ERM perform well on ML tabular classification tasks?

## ERMI achieves state-of-the-art performance on machine learning benchmarks

Tabular classification benchmark based on OpenML-CC18:

- 23 binary classification tasks (less than 100 features)
- 30 training points and 70 testing points (Müller et al. 2022)
- reduced the input dimensionality to 4 (based on highest ANOVA F-value to the target)

*Table 1. Performance metrics on OpenML-CC18 benchmark.*

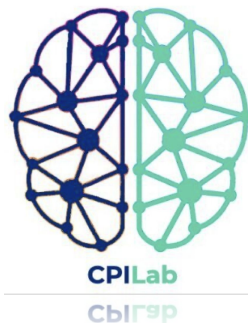
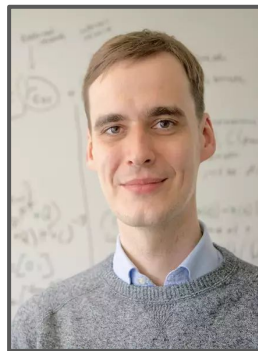
MEAN	SVM	XGBOOST	TABPFN	ERMI
ACC.	69.29%	70.17%	70.51%	<b>70.95%</b>
RANK	2.76	2.61	2.85	<b>2.26</b>

## Summary

- Large language models can generate ecologically valid data
- A class of models called ecologically rational meta-learned inference (ERMI)
- ERMI displays human-like category learning in three experiments
  - human-like learning difficulties
  - human-like learning strategies
  - human-like generalization
- ERMI achieves state-of-the-art performance on ML Benchmarks

**Surprising how far we can go just by meta-learning on the right data!**

# Thank you :)



Collaborators and colleagues from the CPI Lab

# Scraping real-world classification tasks

## OpenML Curated Classification Benchmark (CC-18)

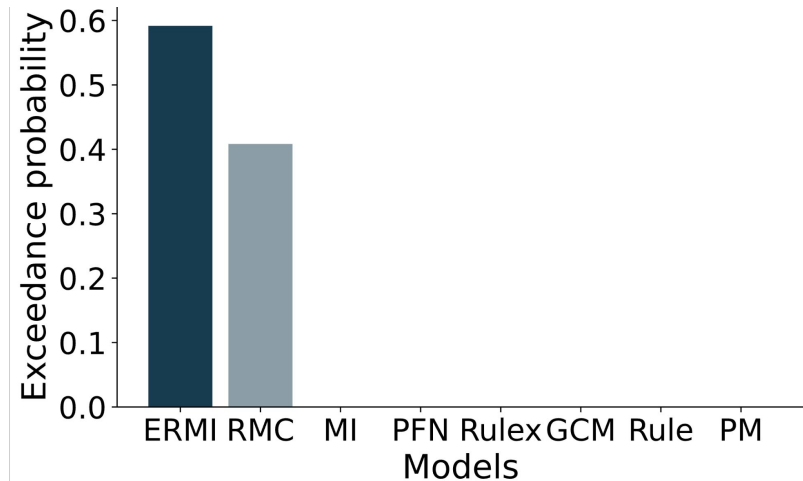
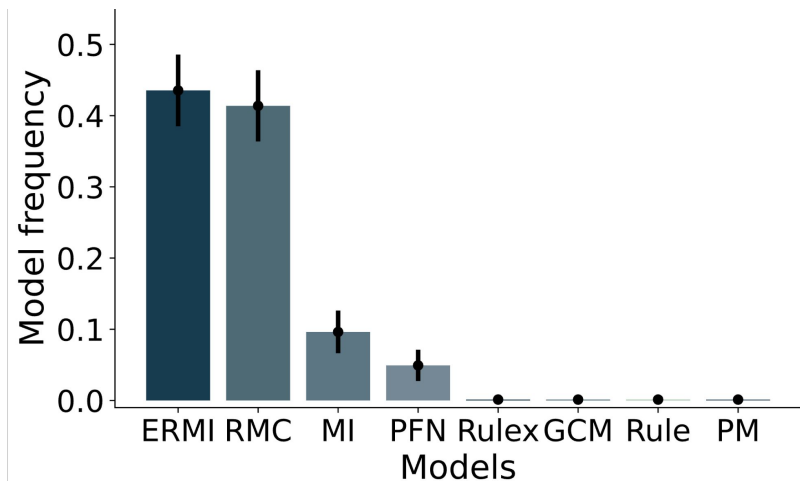
- real referenced datasets
- $500 \leq \text{\#observations} \leq 10000$
- $\text{\#classes} \geq 2$ , each with at least two observations for each
- not predictable using single feature or simple decision tree

## Filtering criteria

- binary classification
- not Nan
- $\text{\#features} < 100$

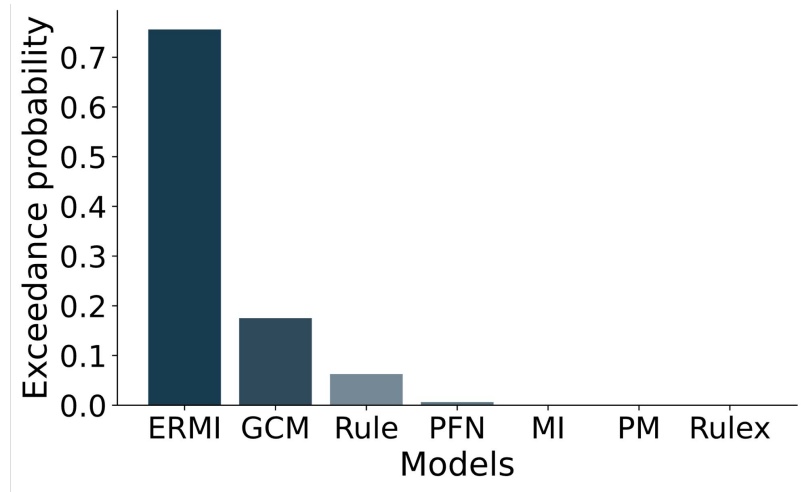
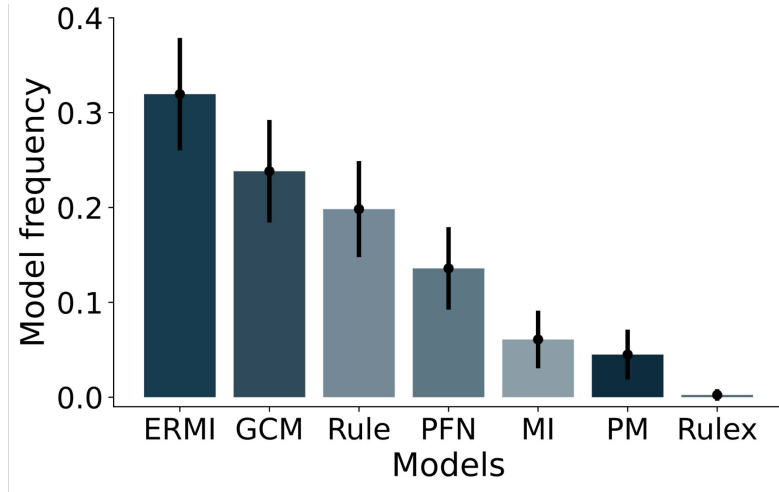
Resulted in a total of 24 datasets

## ERMI explains human data from Badham et al. 2017 better than cognitive models





## ERMI explains human data from Devraj et al. 2022 better than cognitive models



### Prompt for Badham et al. 2017 study

In this experiment, you will be shown examples of geometric objects. Each object has three different features: size, color, and shape. Your job is to learn a rule based on the object features that allows you to tell whether each example belongs in the {A} or {B} category. As you are shown each example, you will be asked to make a category judgment and then you will receive feedback. At first you will have to guess, but you will gain experience as you go along. Try your best to gain mastery of the {A} and {B} categories.

- In trial 1, you picked category {A} for Big Black Square and category {A} was correct.
- In trial 2, you picked category {A} for Small Black Triangle and category {B} was correct

Human: What category would a Small Black Triangle belong to? (Give the answer in the form “Category ⟨your answer⟩”).

Assistant: Category

### Prompt for Devraj et al. 2022 study

In this experiment, you will be shown examples of nonsense word stimuli. Look carefully at each word and decide if it belongs to group {U} or group {M}. Respond with {U} if you think it is a group {U} word and {M} if you think it is a group {M} word. You will receive feedback about the correct group after each of your response. At first, the task will seem quite difficult, but with time and practice, you should be able to answer correctly.

- In trial 1, you picked group {M} for wafuzi and group {U} was correct.
- In trial 2, you picked group {M} for gyfuzi and group {U} was correct.

Human: What group would the word gyfuzi belong to? (Give the answer in the form “Group ⟨your answer⟩”).

Assistant: Group

DATA SET	LOG. REG.	SVM	XGBoost	TABPFN	ERMI
KR-VS-KP CLASSIFICATION	0.8257	0.8514	0.7986	<b>0.8664</b>	0.8450
CREDIT-G CLASSIFICATION	<b>0.6421</b>	0.6357	0.6350	0.6036	0.6150
DIABETES CLASSIFICATION	0.6771	<b>0.7079</b>	0.6786	0.6886	0.6950
SPAMBASE CLASSIFICATION	0.5407	0.7664	0.7536	<b>0.7993</b>	0.7757
TIC-TAC-TOE CLASSIFICATION	0.5536	0.5950	<b>0.6071</b>	0.5914	<b>0.6071</b>
ELECTRICITY CLASSIFICATION	0.5543	0.6007	<b>0.7036</b>	0.6871	0.6436
PC4 SOFTWARE DEFECT PREDICTION	0.7136	0.7521	<b>0.7886</b>	0.7707	0.7714
PC3 SOFTWARE DEFECT PREDICTION	0.6514	0.7264	<b>0.7357</b>	0.7279	0.7107
KC2 SOFTWARE DEFECT PREDICTION	0.5893	0.7314	<b>0.7257</b>	<b>0.7257</b>	<b>0.7257</b>
KC1 SOFTWARE DEFECT PREDICTION	0.6271	0.6707	<b>0.6743</b>	0.6679	0.6521
PC1 SOFTWARE DEFECT PREDICTION	0.5336	0.5964	0.6514	0.6064	0.6493
WDBC CLASSIFICATION	0.9121	0.9207	0.9014	<b>0.9221</b>	0.9093
PHONEME CLASSIFICATION	0.5793	<b>0.7314</b>	0.6979	0.6921	0.7200
QSAR-BIODEG CLASSIFICATION	0.5779	0.7014	0.6850	0.6921	<b>0.7064</b>
ILPD CLASSIFICATION	0.5493	<b>0.6386</b>	0.6229	0.6121	0.6286
OZONE-LEVEL-8HR CLASSIFICATION	0.6614	0.6907	0.6707	0.6471	<b>0.6950</b>
BANKNOTE-AUTHENTICATION CLASSIFICATION	0.7721	0.9229	0.8457	<b>0.9657</b>	0.9379
BLOOD-TRANSFUSION-SERVICE-CENTER	0.4714	0.5493	0.5879	0.5671	<b>0.6186</b>
PHISHING WEBSITES CLASSIFICATION	0.7929	0.8071	<b>0.8157</b>	<b>0.8157</b>	0.8129
BANK-MARKETING CLASSIFICATION	0.5829	0.5614	<b>0.7386</b>	0.7350	0.7171
WILT CLASSIFICATION	0.5171	0.5736	0.6393	0.6371	<b>0.6507</b>
NUMERA128.6 CLASSIFICATION	0.4857	0.4779	<b>0.5029</b>	0.4779	0.4986
CHURN CLASSIFICATION	0.6321	0.7271	0.6800	0.7186	<b>0.7329</b>
<b>MEAN ACC.</b>	62.80 $\pm$ 0.66	69.29 $\pm$ 0.62	70.17% $\pm$ 0.52	70.51% $\pm$ 0.63	<b>70.95% <math>\pm</math> 0.54</b>
<b>MEAN RANK</b>	4.52 $\pm$ 0.21	2.76 $\pm$ 0.26	2.61 $\pm$ 0.30	2.85 $\pm$ 0.27	<b>2.26 <math>\pm</math> 0.22</b>