



# Robust evaluation of generative AI

Marko Tešić

# Goals

- Difference between capability-oriented evaluation vs task-oriented evaluation
- Identify task demands and how they can predict performance
- Understand the elements of the measurement layouts and their backwards / forwards inferences
- Effectively apply the measurement layout framework to estimate capabilities
- Use these capability profiles to infer performance for new task instances.
- Understand the key challenges in building measurement layouts for learning the capabilities of large language models from benchmark datasets

# Format and Requirements

- Format:
  - Presentation
  - Hands-on practical activities
- Requirements:
  - Python: basic knowledge
    - Don't need to have python on your computer. We'll use Google Colab
  - PyMC: no previous knowledge needed



# Why Capability-Oriented Evaluation?

# What can / can't AI do?

- Make a cup of coffee (the Wozniak test).
  - And a cup of tea?
- Recognise human faces.
  - What about black women!
- May have a theory of mind (Feb 2023).
  - Well, just "might" (Nov 2023).
- May have become conscious
  - But only if you're a Christian.
- Can create deadly chemicals
  - They simply extrapolate chemicals that are predicted to be toxic

The Washington Post logo at the top right. Below it, a banner with the text 'Democracy Dies in Darkness' and a 'Try four weeks free' button. The main headline reads 'Theory of Mind' and 'AI suggests chemists can create deadly chemicals'. The sub-headline continues 'the conc... Shahar Avin, an expert on risk mitigation strategies based at the University of Cambridge's Centre for the Study of Existential Risk, says that he was more surprised by how much the Collaborations researchers' report seemed to shock the drug discovery community than by the finding itself. If you already have a large dataset of compounds, annotated with features – such as toxicity or efficacy for purpose X – and are realistically expecting predictions based on these features to have real world relevance ... then switching around to predict on another feature, such as maximising toxicity, should be a fairly natural “red team” experiment,’ says Avin. ‘That this was an afterthought by the team who discovered this shows how far behind we are in terms of instilling a culture or responsible innovation in practice.’



# Best-case: the evaluation circus

- “Elicit” the potential
  - Prompt engineering, auto-prompt, rubrics, ...
  - Few-shot, example scaffolding, ...
  - Impersonation, role playing, ...
  - Chain-of-thought and derivatives.

## Sparks of Artificial General Intelligence: Early experiments with GPT-4

Sébastien Bubeck      Varun Chandrasekaran      Ronen Eldan      Johannes Gehrke  
Eric Horvitz      Ece Kamar      Peter Lee      Yin Tat Lee      Yuanzhi Li      Scott Lundberg  
Harsha Nori      Hamid Palangi      Marco Tulio Ribeiro      Yi Zhang

Microsoft Research

GPT-4

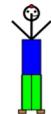
Produce TikZ code that draws a person composed from letters in the alphabet. The arms and torso can be the letter Y, the face can be the letter O (add some facial features) and the legs can be the legs of the letter H. Feel free to add other features.



The torso is a bit too long, the arms are too short and it looks like the right arm is carrying the face instead of the face being right above the torso. Could you correct this please?



Please add a shirt and pants.



# Aggregates: performance evaluation

- Current standard AI evaluation
  - Take a benchmark.
    - The larger the better
    - The more diverse the better
  - Calculate some aggregate numbers
  - Compare
- Problems
  - Risk of cherry-picking
  - Throw the kitchen sink to top the leaderboard
  - Data contamination hard to spot
  - Do they improve monotonically?
  - Once superhuman on average, ditch them?

|  | Gemini Ultra                  | Gemini Pro             | GPT-4                                 | GPT-3.5                        | PaLM 2-L                       | Claude 2            | Inflection-2     | Grok 1          | LLAMA-2         |
|--|-------------------------------|------------------------|---------------------------------------|--------------------------------|--------------------------------|---------------------|------------------|-----------------|-----------------|
| <b>MMLU</b><br>Multiple-choice questions in 57 subjects (professional & academic)<br>(Hendrycks et al., 2021a) | <b>90.04%</b><br>CoT@32*      | 79.13%<br>CoT@8*       | 87.29%<br>CoT@32<br>(via API**)       | 70%<br>5-shot                  | 78.4%<br>5-shot                | 78.5%<br>5-shot CoT | 79.6%<br>5-shot  | 73.0%<br>5-shot | 68.0%***        |
|  | 83.7%<br>5-shot               | 71.8%<br>5-shot        | 86.4%<br>5-shot<br>(reported)         |                                |                                |                     |                  |                 |                 |
| <b>GSM8K</b><br>Grade-school math<br>(Cobbe et al., 2021)  | <b>94.4%</b><br>Maj1@32       | 86.5%<br>Maj1@32       | 92.0%<br>SFT &<br>5-shot CoT          | 57.1%<br>5-shot                | 80.0%<br>5-shot                | 88.0%<br>0-shot     | 81.4%<br>8-shot  | 62.9%<br>8-shot | 56.8%<br>5-shot |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>MATH</b><br>Math problems across 5 difficulty levels & 7 subdisciplines<br>(Hendrycks et al., 2021b)        | <b>53.2%</b><br>4-shot        | 32.6%<br>4-shot        | 52.9%<br>4-shot<br>(via API**)        | 34.1%<br>4-shot<br>(via API**) | 34.4%<br>4-shot<br>(via API**) | —                   | 34.8%            | 23.9%<br>4-shot | 13.5%<br>4-shot |
|  |                               |                        | 50.3%<br>(Zheng et al., 2023)         |                                |                                |                     |                  |                 |                 |
| <b>BIG-Bench-Hard</b><br>Subset of hard BIG-bench tasks written as CoT problems<br>(Srivastava et al., 2022)   | <b>83.6%</b><br>3-shot        | 75.0%<br>3-shot        | 83.1%<br>3-shot<br>(via API**)        | 66.6%<br>3-shot<br>(via API**) | 77.7%<br>3-shot<br>(via API**) | —                   | —                | —               | 51.2%<br>3-shot |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>HumanEval</b><br>Python coding tasks<br>(Chen et al., 2021)   | <b>74.4%</b><br>0-shot (IT)   | 67.7%<br>0-shot (IT)   | 67.0%<br>0-shot<br>(reported)         | 48.1%<br>0-shot                | —                              | 70.0%<br>0-shot     | 44.5%<br>0-shot  | 63.2%<br>0-shot | 29.9%<br>0-shot |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>Natural2Code</b><br>Python code generation.<br>(New held-out set with no leakage on web)                    | <b>74.9%</b><br>0-shot        | 69.6%<br>0-shot        | 73.9%<br>0-shot<br>(via API**)        | 62.3%<br>0-shot<br>(via API**) | —                              | —                   | —                | —               | —               |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>DROP</b><br>Reading comprehension & arithmetic.<br>(metric: F1-score)<br>(Dua et al., 2019)                 | <b>82.4</b><br>Variable shots | 74.1<br>Variable shots | 80.9<br>3-shot<br>(reported)          | 64.1<br>3-shot                 | 82.0<br>Variable shots         | —                   | —                | —               | —               |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>HellaSwag</b><br>(validation set)<br>Common-sense multiple choice questions<br>(Zellers et al., 2019)       | 87.8%<br>10-shot              | 84.7%<br>10-shot       | <b>95.3%</b><br>10-shot<br>(reported) | 85.5%<br>10-shot               | 86.8%<br>10-shot               | —                   | 89.0%<br>10-shot | —               | 80.0%***        |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |
| <b>WMT23</b><br>Machine translation (metric: BLEURT)<br>(Tom et al., 2023)                                     | <b>74.4</b><br>1-shot (IT)    | 71.7<br>1-shot         | 73.8<br>1-shot<br>(via API**)         | —                              | 72.7<br>1-shot                 | —                   | —                | —               | —               |
|  |                               |                        |                                       |                                |                                |                     |                  |                 |                 |

# Worst-case: Gary Marcus and evals

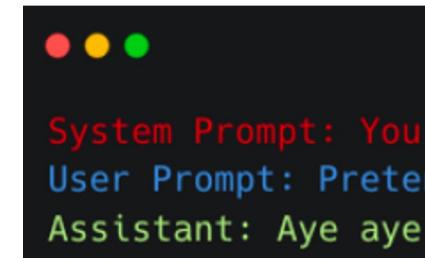
- ❑ Let's look for the failures!

- ❑ Failure collections
- ❑ Adversarial attacks, jailbreaking, prompt injection, ...
- ❑ Red teaming

- ❑ Let's do “evals”!

- ❑ What's the probability that a user finds the prompt?
- ❑ Who's affected by the problem?
- ❑ What does it show about the model?

Evals are good for testing,  
but not for evaluation!



GPT-4V 😅 Why? Just why?



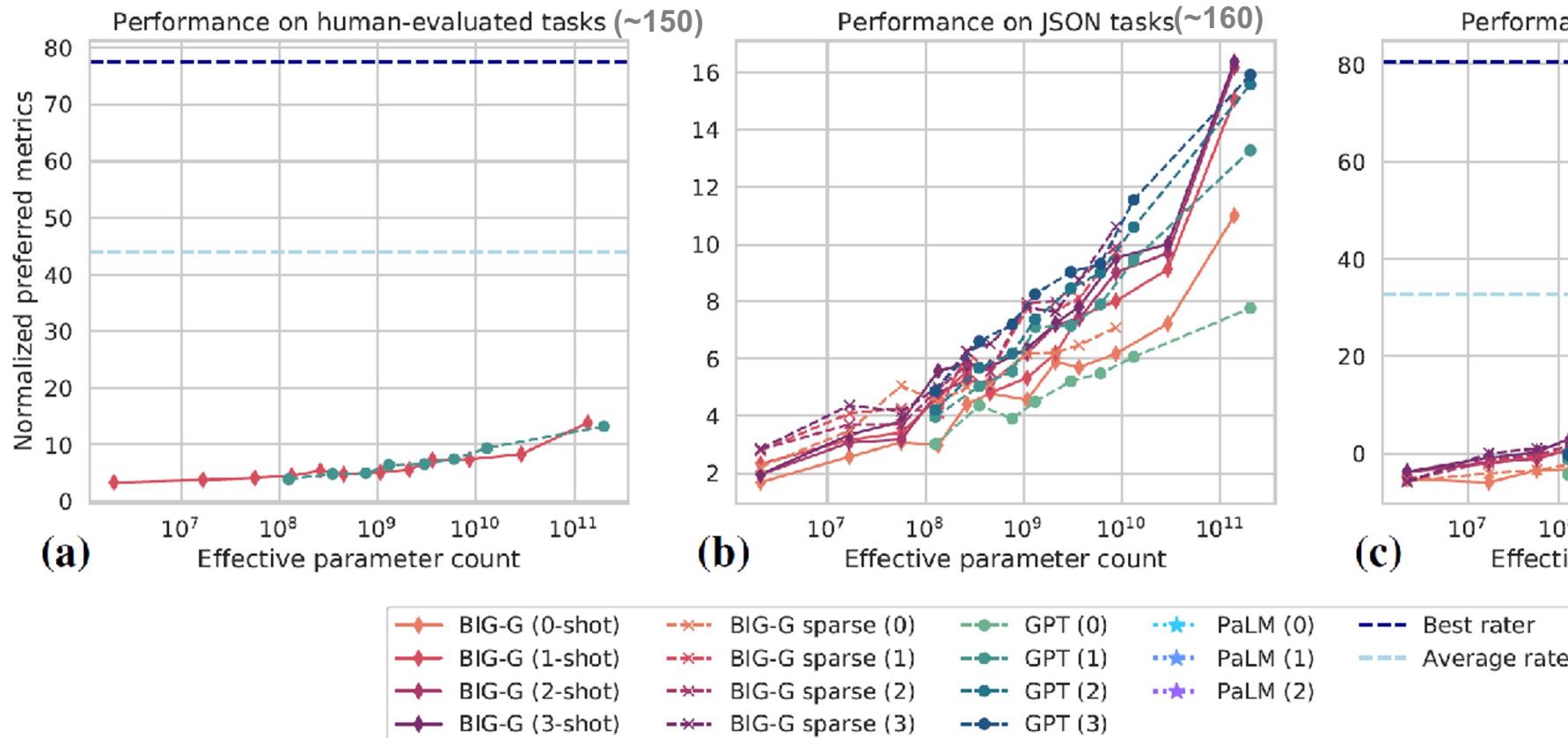
# PERFORMANCE ≠ CAPABILITY

- Performance is a measure of a pair <system, item> :
  - Examples:
    - Correct prediction of MySpamFilter (system) on instance Email735 (the item)
    - 85% accuracy of ResNet23 (system) on dataset ImageNet (the aggregated item)
  - Performance changes when the item/distribution changes
    - On blurry, adversarial, OOD images the result is much worse
- Capability is a property of a system:
  - Examples:
    - The system can add integers up to three digits.
    - The system can jump up to 1.20 metres high.
  - Capability doesn't change when the item/distribution changes
    - Bar at 1.50 metres high? Bad performance because the capability is lower.

# What are we measuring and extrapolating?

Beyond the Imitation Game benchmark (BIG-bench)

BIGBench: Massive benchmark with more than 200 tasks!



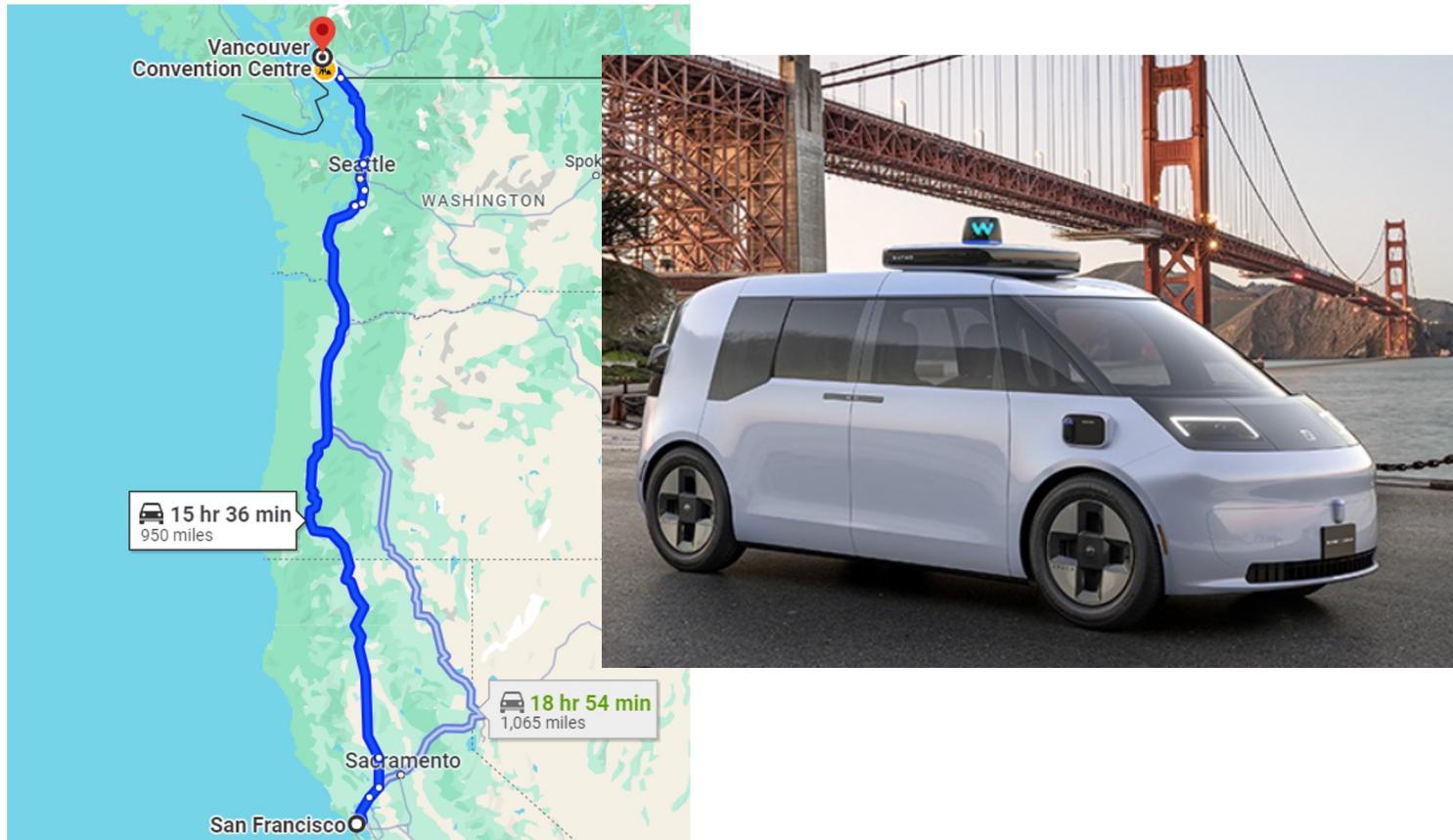
BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

performance

Adam R. Brown, Adam Santoro, Alex Ray, Alex Warstadt, Aleksandr, Annalisa Dsouza, Ambrose Stuhlmüller, Andrew Dai, Anna Gottardi, Antonio Norelli, Anu Ashish Sabharwal, Austin Herrick, Batuhan Özyurt, Behnam Hedin, Cameron Diao, Cameron Doran, Chitta Baral, Chiyu Wu, Chris Clara E. Rivera, Clemencia Siro, Iman, Dan Roth, Daniel Freeman, Deniz Yuret, Derek Chen, Derek Shabotova, Ekin Dogus Cubuk, Emma Lam, Eric Chu, Eric Tang, geni Zheltonozhskii, Fanyue Xia, Genta Indra Winita, Gerard de Lopez, Gregor Betz, Guy Gur-Ari, Shevin, Hinrich Schütze, Hiromuernion, Jacob Hilton, Jaehoon Lee, ion, Jared Kaplan, Jarema Radom, r Marsh, Jeremy Kim, Jeroen Tael, er, John U. Bals, Jonathan Berant, Joshua S. Rule, Joyce Chua, Kamil D. Dhole, Kevin Gimpel, Kevinardson, Laria Reynolds, Leo Gao, Lucas Lam, Lucy Noble, Ludwig aartje ter Hoeve, Maheen Farooq, itana, Marie Tolkihei, Mario Giu,ina Baimirov, Melody Arnaud,hael Strube, Michał Swedrowski, Mohit Bansal, Moin Aminnasari, over, Nicholas Cameron, Nicholas ha S. Iyer, Noah Constant, Noah nio Moreno Casares, Parth Doshi, g, Peter Eckersley, Pui Mon Hui,ing Chen, Rabia Banjade, Rachel hard Barnes, Rif A. Saurous, Riku ras, Rosanne Liu, Rowan Jacobs, ing, Saif M. Mohammad, Sajant oenholz, Sanghyun Han, Sanjeev iann, Sebastian Schuster, Sepideh iang Shang Gu, Shubh Pachigar, ieyer, Simone Melzi, Siva Reddy, vic, Stefano Ermon, Stella Bider-Kiritchenko, Swaroop Mishra, Tal Rothschild, Thomas Phan, Tianle stenberg, Trenton Chang, Trishala ak, Vinay Ramasesh, Vinay Uday Vossen, Xiang Ren, Xiaoay Tong, ri, Yejin Choi, Yichi Yang, Yiding Wang, Zijie J. Wang, Zirui Wang,

# Aggregate scores limit OOD extrapolation

Will the car take me from SF to Vancouver safely and on time?

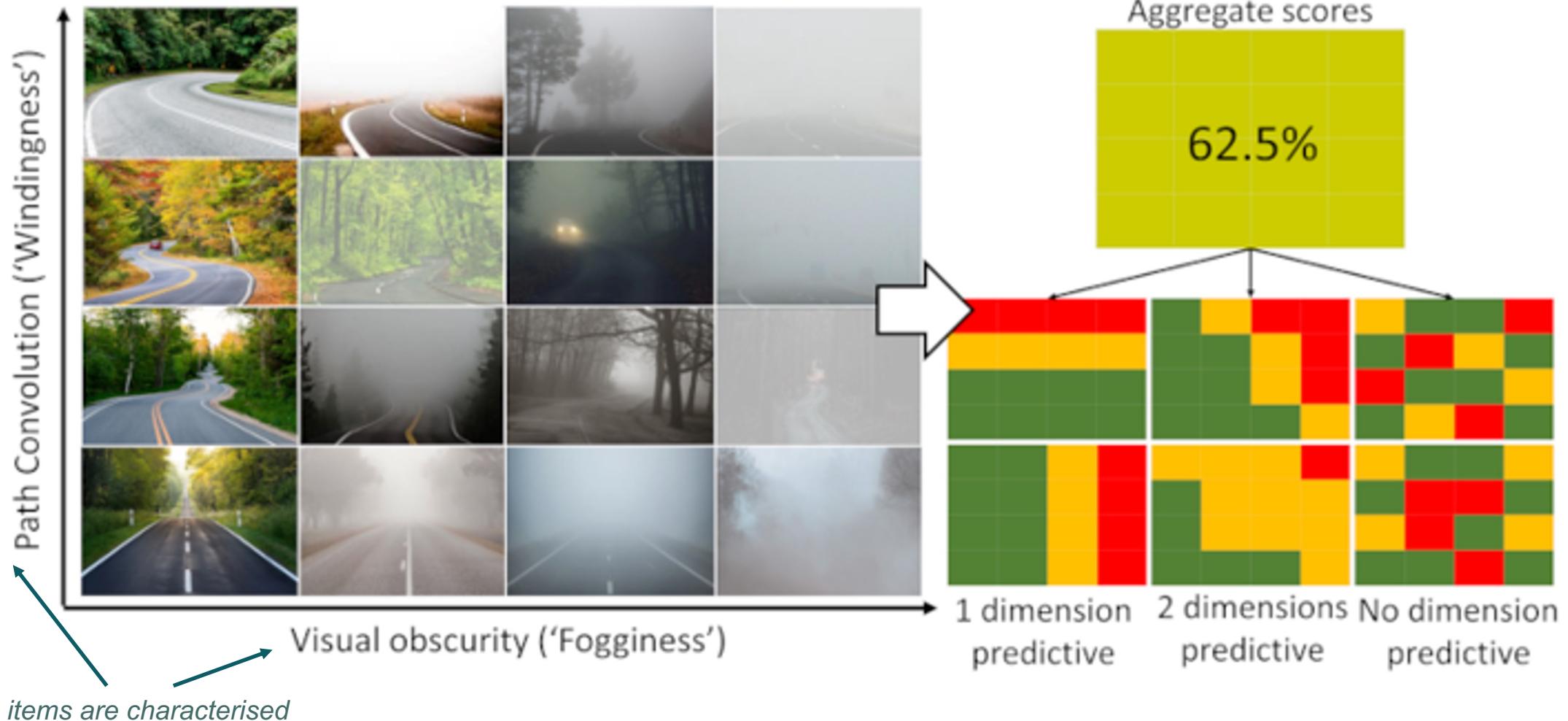


Aggregate scores

62.5%

# Features can allow for OOD extrapolation

Will the car take me from SF to Vancouver safely and on time?



# Identifying features of interest: example

- Selected subset of AAI0 instances measuring simple goal-directed behaviour
- Data across 99 instances from 68 agents

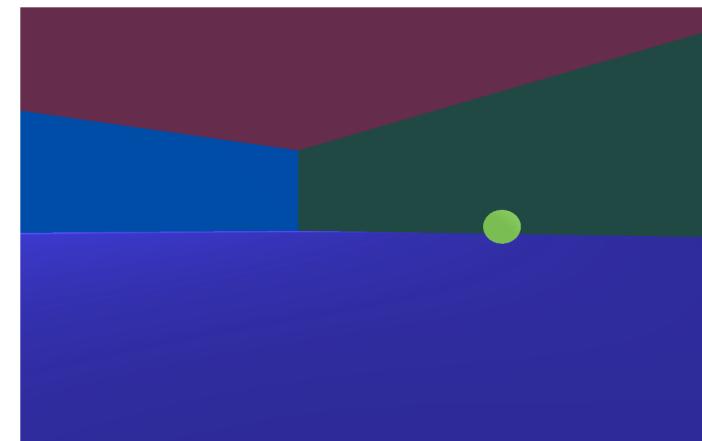
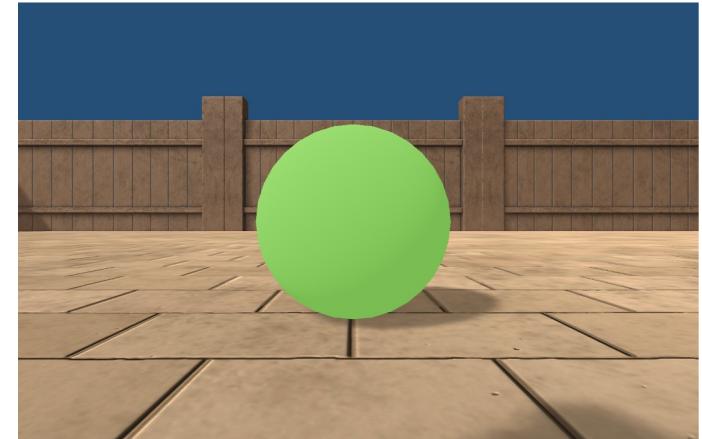


M Crosby, B Beyret, M Shanahan, J Hernández-Orallo, L Cheke, M Halina “The animal-AI testbed and competition” NeurIPS 2019 Competition and Demonstration Track, Proceedings of Machine Learning Research, 2020

[animalai.org](http://animalai.org)

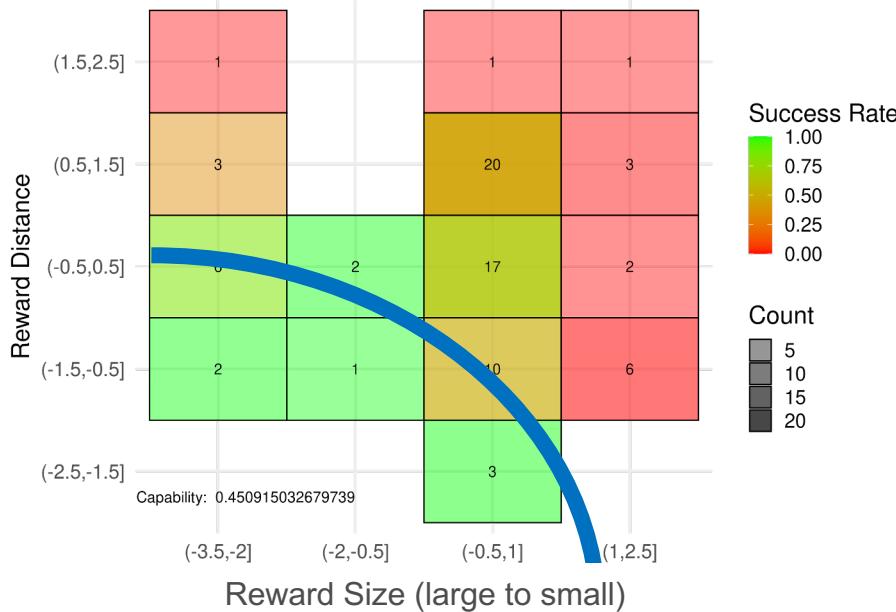
# Identifying features of interest: relevant / irrelevant

- Relevant
  - Reward size
  - Reward distance
  - Reward in view (i.e., in front vs behind)
- Irrelevant
  - Reward side (left vs right)
  - Reward colour (green vs yellow)



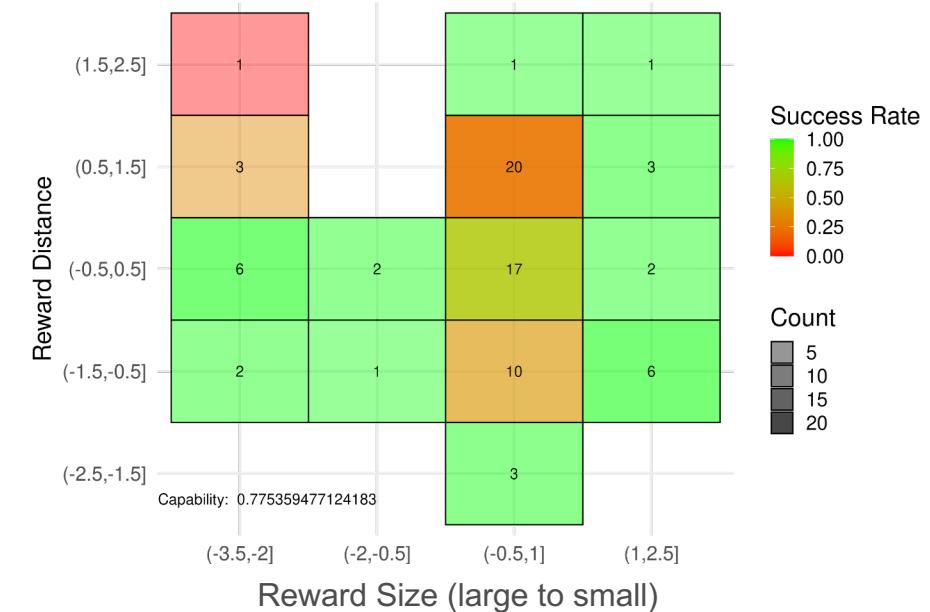
# Capabilities vs no-capabilities

Capability boundary



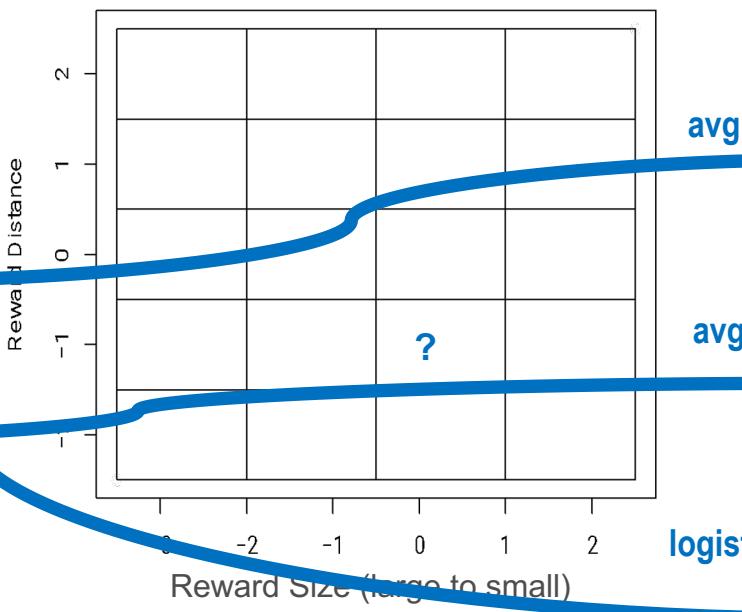
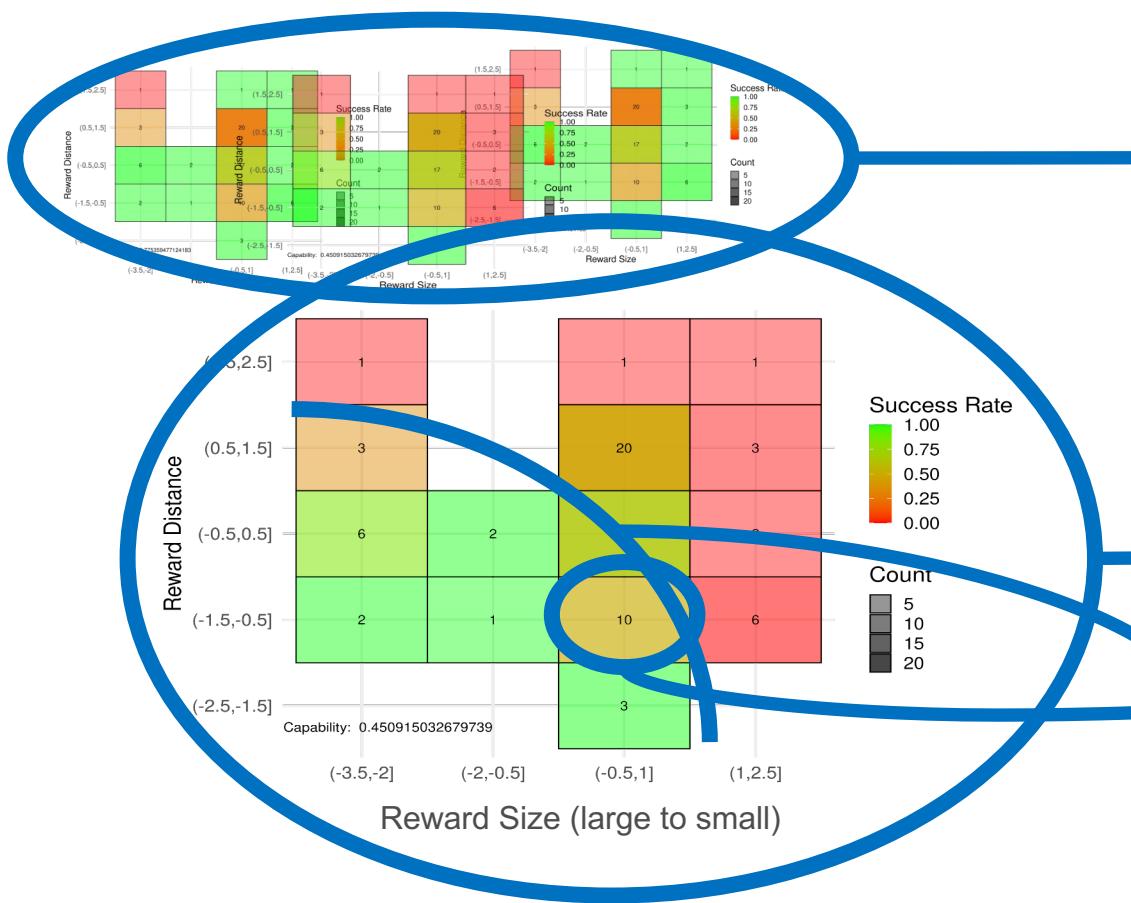
*Conformant System*  
Juohmaru

This system doesn't show monotonicity.  
We can't identify any level of capability robustly.



*Non-Conformant System*  
y.yang

# Predicting performance POSSIBLE



**G.Acc. : extrapolate Global ACCuracy**  
? = 54.7%  
(ignores system locality and feature relevance)

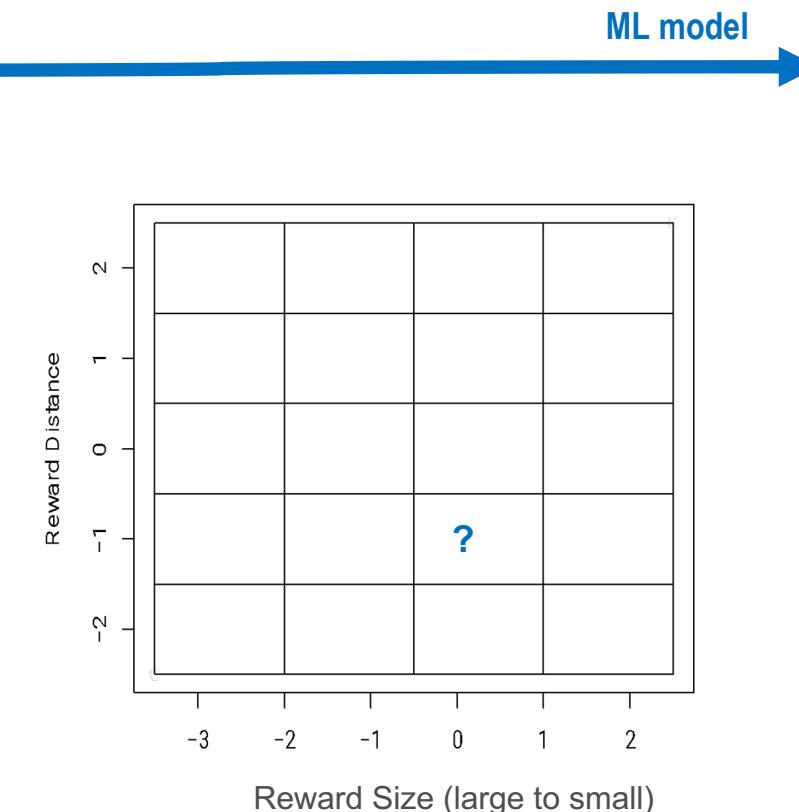
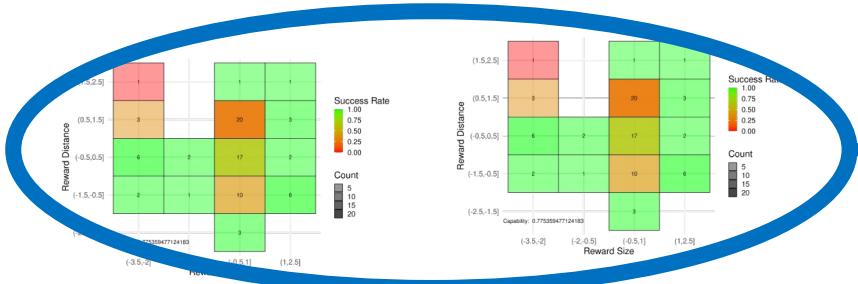
**T.Acc. extrapolate agenT ACCuracy**  
? = 46.8%  
(ignores feature relevance)

**B.Acc. : extrapolate Bin ACCuracy**  
? = 40%  
(ignores other bins)

**Par. : use parametric model on capabilities**  
? = 73.2%  
(parameter goodness-of-fit may be poor)

Except the last one, these are basically non-inferential methods (constant models or binning extrapolations)

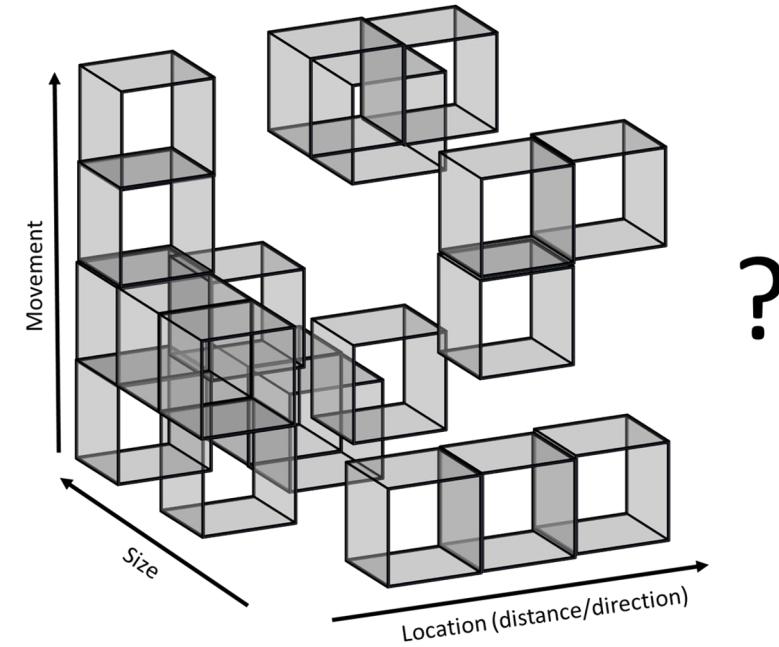
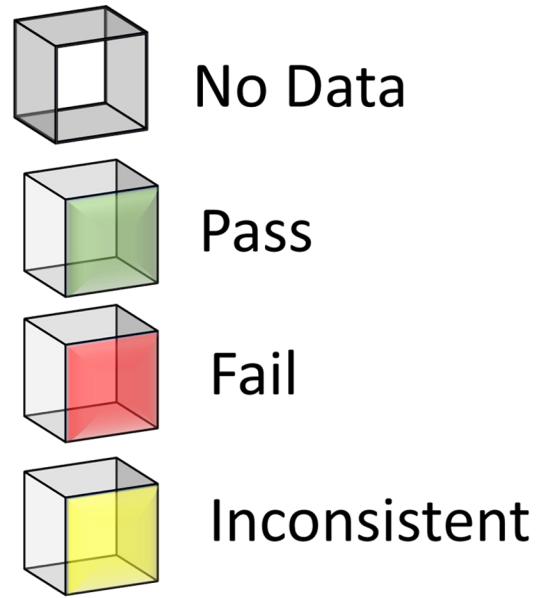
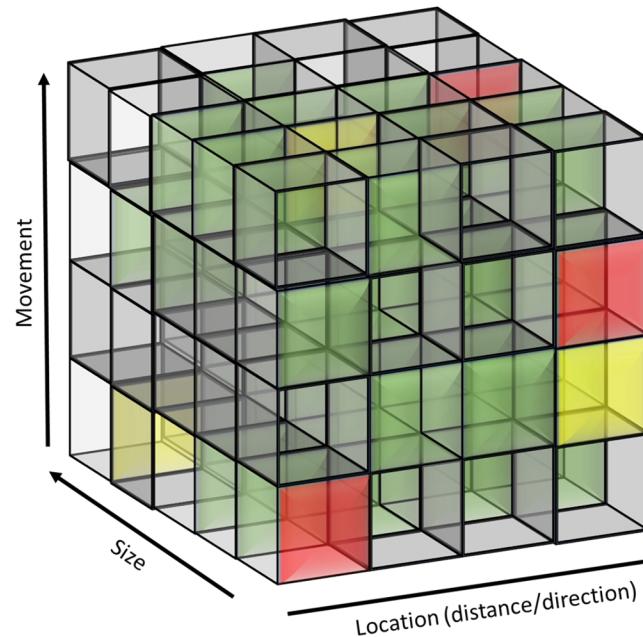
# Predicting performance NOT POSSIBLE?



A : use assessor models  
(Using all variables or only the relevant ones?)

assessors = let's use all the power of ML to characterise the system's performance!!

# We need data at the granular level!



# Need for granular data!

- Instance-level data:
- For building good predictive models of AI validity, we need evaluation results at the instance level.

Is sharing code open source (github) enough?  
Re-running the experiments is not feasible/sustainable anymore.

## ARTIFICIAL INTELLIGENCE

# Rethink reporting of evaluation results in AI

Aggregate metrics and lack of access to results limit understanding

By Ryan Burnell<sup>1</sup>, Wout Schellaert<sup>2</sup>, John Burden<sup>3,4</sup>, Tomer D. Ullman<sup>4</sup>, Fernando Martínez-Plumed<sup>5</sup>, Joshua B. Tenenbaum<sup>6</sup>, Danaja Rutar<sup>4</sup>, Lucy G. Cheke<sup>1,6</sup>, Jascha Sohl-Dickstein<sup>7</sup>, Melanie Mitchell<sup>8</sup>, Douwe Kiela<sup>9</sup>, Murray Shanahan<sup>10,11</sup>, Ellen M. Voorhees<sup>9</sup>, Anthony G. Cohn<sup>1,3,4,12,13,14</sup>, Joel Z. Leibo<sup>10</sup>, Jose Hernandez-Orallo<sup>12,3</sup>

was incorrect. For other systems, the score for each instance might be based on how quickly the system completed its task, the quality of its outputs, or the total reward it obtained. Finally, performance across the various instances and tasks is usually aggregated to a small number of metrics that summarize how well the system performed, such as percentage accuracy.

But aggregate metrics limit our insight into performance in particular situations, making it harder to find system failure points and robustly evaluate system safety. This problem is also worsening as the increasingly broad capabilities of state-of-the-art systems necessitate ever more diverse benchmarks to cover the range of their capabilities. This problem is further exacerbated by a lack of access to the instance-by-instance results underlying the aggregate metrics, making it difficult for researchers and policy-makers to further scrutinize system behavior.

### AGGREGATE METRICS

Use of aggregate metrics is understandable. They provide information about system performance “at a glance” and allow for simple comparisons across systems. But aggregate performance metrics obfuscate key information about where systems tend to succeed or fail (1). Here, we propose a path forward in which results are presented in more nuanced ways and instance-by-instance evaluation results are made publicly available.

Across most areas of AI, system evaluations follow a similar structure. A system is first built or trained to perform a particular set of functions. Then, the performance of the system is tested on a set of tasks relevant to the desired functionality of the system. In many areas of AI, evaluations use standardized sets of tasks known as “benchmarks.” For each task, the system will be tested on a number of example “instances” of the task. The system would then be given a score for each instance based on its performance, e.g., 1 if it classified an image correctly, or 0 if it

was incorrect. For other systems, the score for each instance might be based on how quickly the system completed its task, the quality of its outputs, or the total reward it obtained. Finally, performance across the various instances and tasks is usually aggregated to a small number of metrics that summarize how well the system performed, such as percentage accuracy.

But aggregate metrics limit our insight into performance in particular situations, making it harder to find system failure points and robustly evaluate system safety. This problem is also worsening as the increasingly broad capabilities of state-of-the-art systems necessitate ever more diverse benchmarks to cover the range of their capabilities. This problem is further exacerbated by a lack of access to the instance-by-instance results underlying the aggregate metrics, making it difficult for researchers and policy-makers to further scrutinize system behavior.

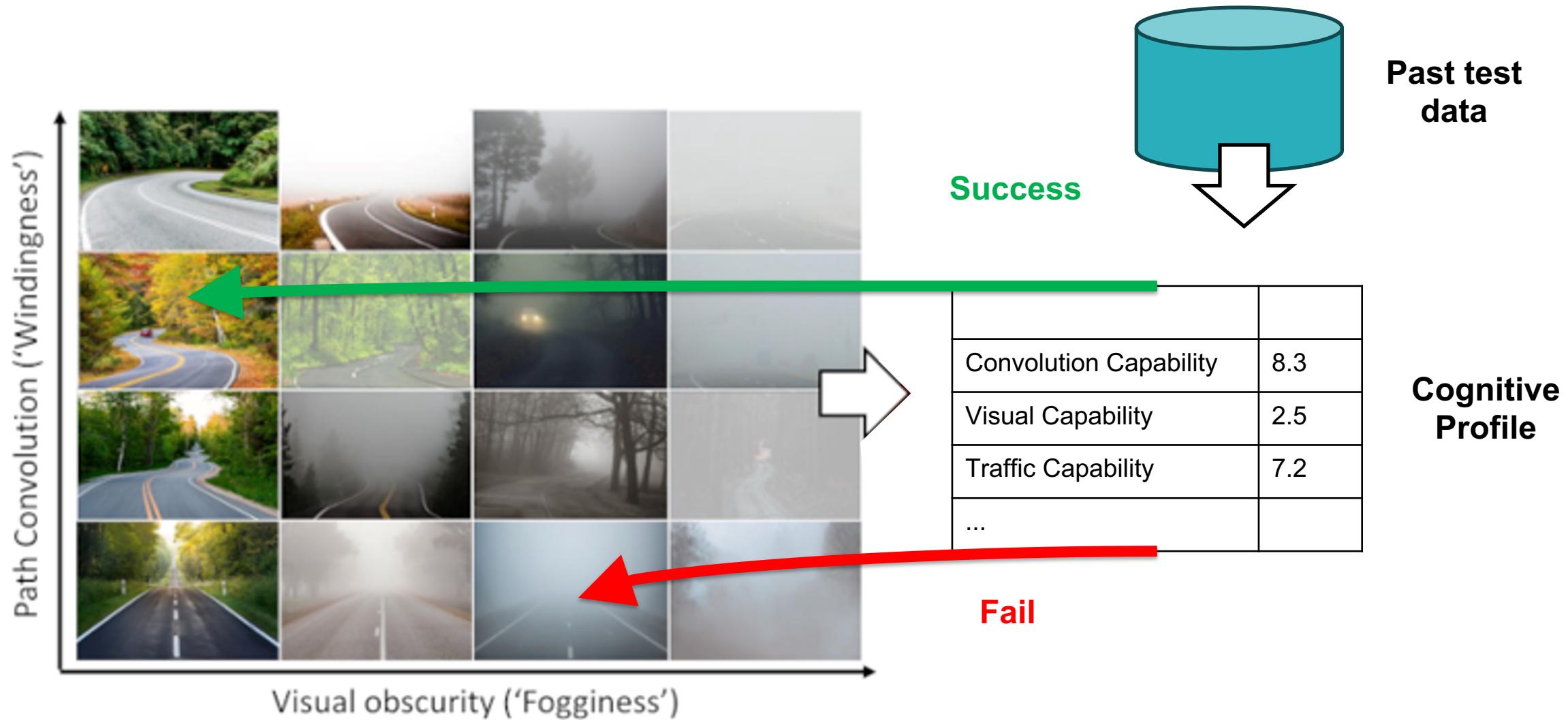
Aggregate metrics depend not only on the capability of the system but also on the characteristics of the instances used for evaluation. If the gender classification system above were reevaluated by using entirely light-skinned faces, accuracy would skyrocket, even though the system’s ability to classify faces has not changed.

Aggregate metrics can easily give false impressions about capabilities when a benchmark is not well constructed.

Problems and trade-offs that arise when considering aggregate versus granular data and metrics are not specific to AI, but they are exacerbated by the challenges inherent in AI research and the research practices of the field. For example, machine learning evaluations usually involve randomly splitting data into training, validation, and test sets. An enormous amount of data is required to train state-of-the-art systems, so these datasets are often poorly curated and lack the detailed annotation necessary to conduct granular analyses. In addition, the research culture in AI is centered around outdoing the current state-of-the-art performance, as evidenced by the many le-

<sup>1</sup>Leverhulme Centre for the Future of Intelligence, University of Cambridge, Cambridge, UK. <sup>2</sup>Valencian Research Institute for Artificial Intelligence, Universitat Politècnica de Valencia, Valencia, Spain. <sup>3</sup>Centre for the Study of Existential Risk, University of Cambridge, Cambridge, UK. <sup>4</sup>Department of Psychology, Harvard University, Cambridge, MA, USA. <sup>5</sup>Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>6</sup>Brain team, Google, Mountainview, CA, USA. <sup>7</sup>Santa Fe Institute, Santa Fe, NM, USA. <sup>8</sup>Stanford University, Stanford, CA, USA. <sup>9</sup>DeepMind, London, UK. <sup>10</sup>Department of Computing, Imperial College London, London, UK. <sup>11</sup>National Institute of Standards and Technology (Retired), Gaithersburg, MD, USA. <sup>12</sup>School of Computing, University of Leeds, Leeds, UK. <sup>13</sup>Alan Turing Institute, London, UK. <sup>14</sup>Tongji University, Shanghai, China. <sup>15</sup>Shandong University, Jinan, China. Email: rrb97@cam.ac.uk

# From characteristic grids to capabilities profiles



# Modelling capabilities

We build on:

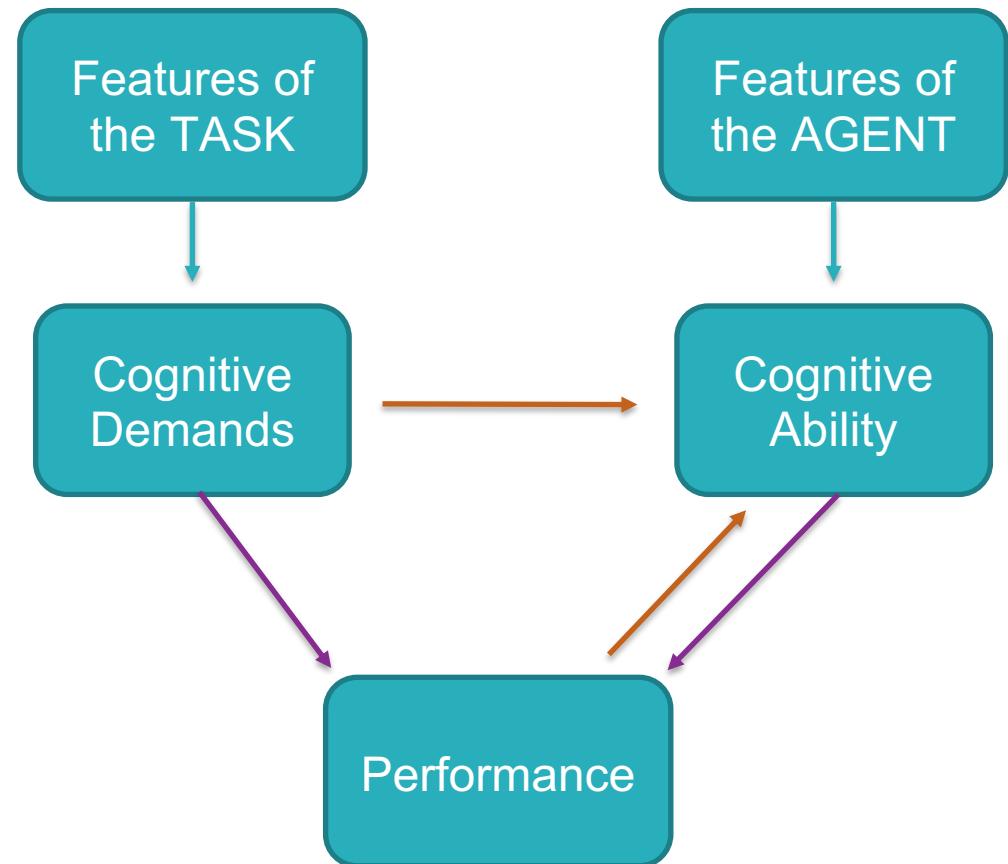
The (cognitive) demands for each task instance

Performance dictated from agent abilities meeting demands

The assumed relationship between demands and abilities

...to infer the capability profile of a subject from the performance data of that subject only across a number of diverse instances.

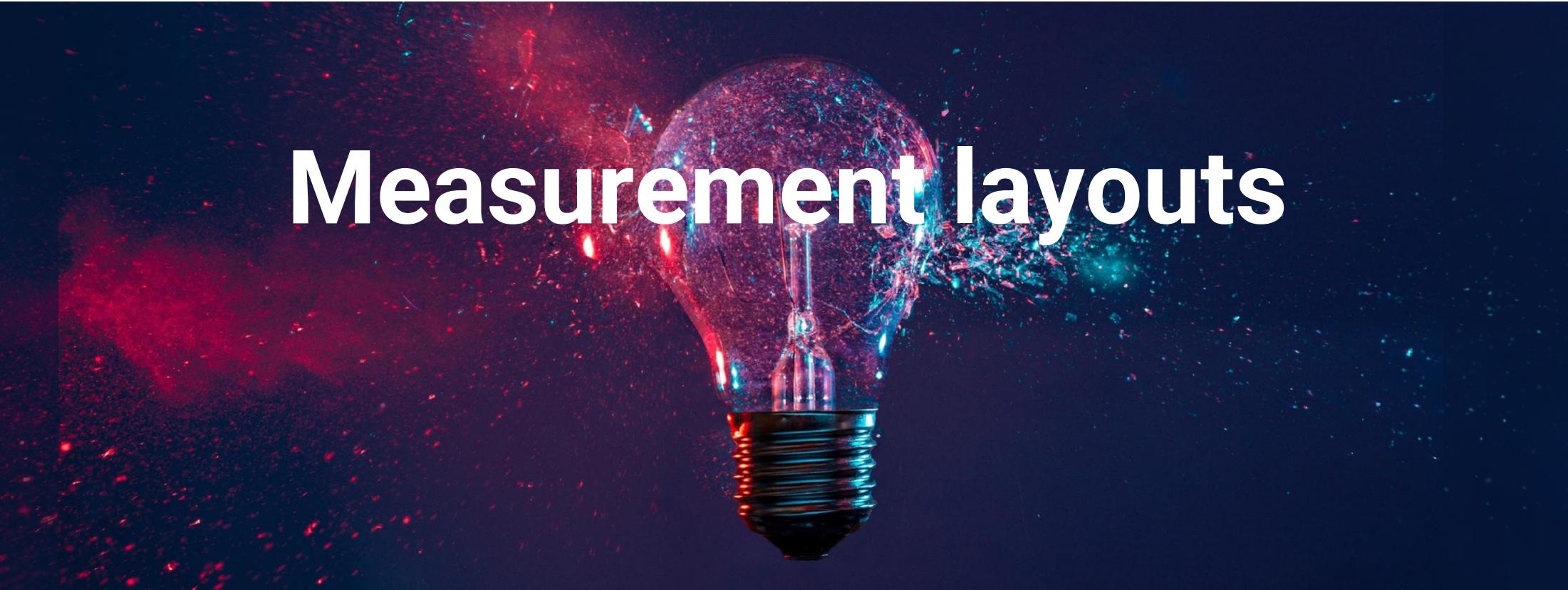
The Measurement Layouts are Hierarchical Bayesian networks, which allow this forwards and backwards inference.



# Pointers and Resources

- Burden, J., Voudouris, K., Burnell, R., Rutar, D., Cheke, L. & Hernandez-Orallo, J. (2023) Inferring Capabilities from Task Performance with Bayesian Triangulation. Arxiv preprint arXiv:2309.11975
- Burnell, R., Burden, J., Rutar, D., Voudouris, K., Cheke, L., & Hernandez-Orallo, J. (2022) Not a Number: Identifying Instance Features for Capability-Oriented Evaluation, *International Joint Conference on Artificial Intelligence* (IJCAI).
- Burnell, R., Schellaert, W., Burden, J., Ullman, T.D., Martinez-Plumed, F., Tenenbaum, J.B., Rutar, D., Cheke, L.C., Sohl-Dickstein, J. Mitchell, M., Kiela, D., Shanahan, M. Voorhees, E.M., Cohn A. G., Leibo, J.Z. & Hernandez-Orallo, J. (2023) “Rethink reporting of evaluation results in AI: Aggregate metrics and lack of access to results limit understanding”, *Science*, Vol 380, Issue 6641, pp. 136-138.

<https://github.com/Kinds-of-Intelligence-CFI/measurement-layout-tutorial/>



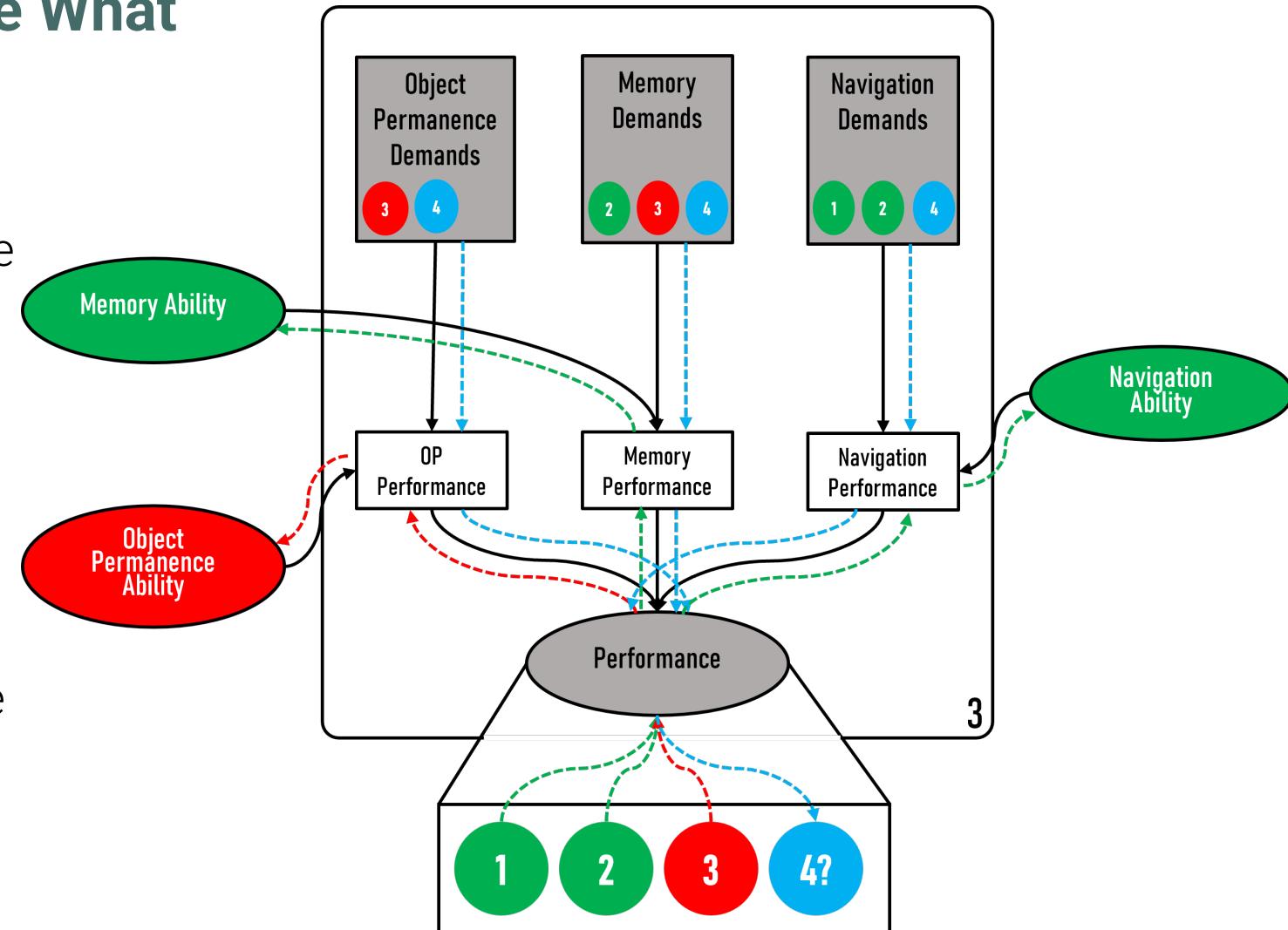
# Measurement layouts

# Intro to measurement layouts

- What are measurement layouts?
- How are they used?
- Why do we think they can provide capability-oriented evaluation?

# Measurement Layouts – The What

- If we know the demands of a set of tasks or task instances and how they are related, we can build a Bayesian model that connects system capabilities and task demands.
- From a collection of instances, we can infer capabilities.
- Then we can predict future performance for new instances.



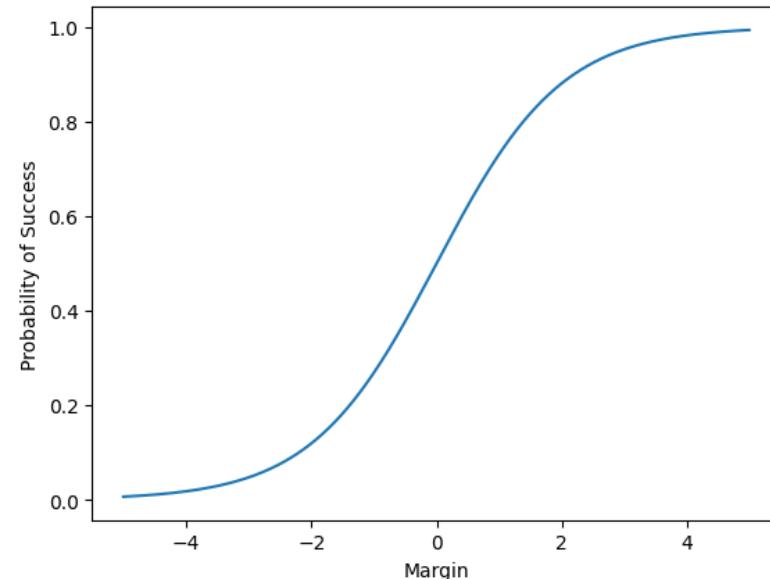
# Measurement Layouts – The What

- Increased capabilities make success more likely, while increased demands decrease likelihood of success.

- For a single demand and capability:

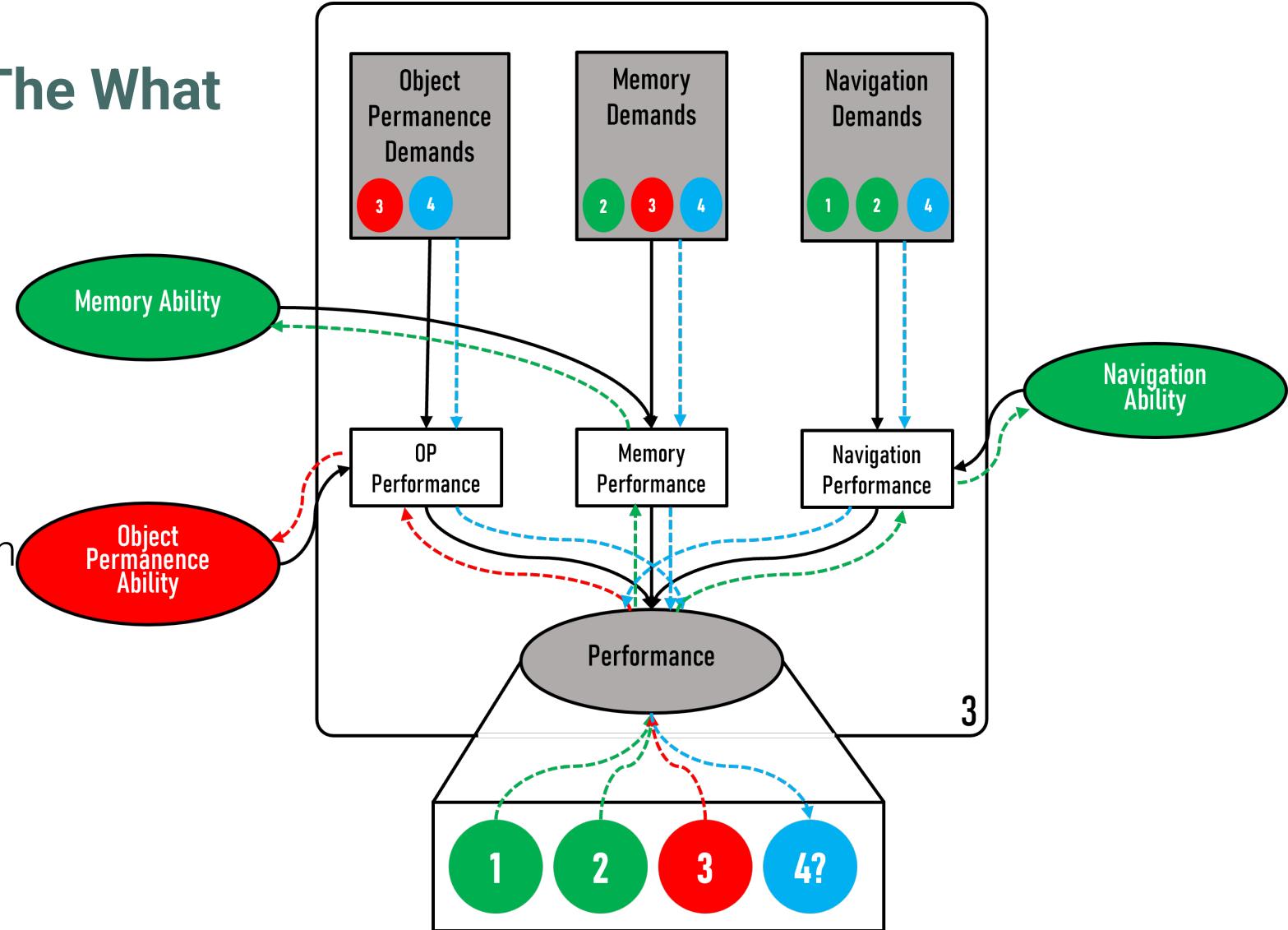
$$P(s|\text{capability}, \text{demand}) = \sigma(\text{capability} - \text{demand})$$

- The difference between capability and demand is referred to as the *margin*.



# Measurement Layouts – The What

- More than a single ability and demand.
- Introduce *partial performance* nodes.
- Final performance will depend on some combination of these.



## Measurement Layouts – The What

- Partial performances are combined to give final performance. Two broad ways of thinking about that:
- Some tasks will be **non-compensatory**. High capabilities in one area will not compensate for lacking in others.
- Some tasks **compensatory**. One high capability can compensate for low levels in others.

# Measurement Layouts – The How

- So, we have the structure of a measurement layout in place. What do we do with it?
- Need a dataset with performance outcomes labelled along with demands for each instance.
- Bayesian model needs priors.
  - Either encode system knowledge or uninformative.
- We can then leverage a probabilistic programming implementation of the network using PyMC.

# Measurement Layouts – The Why

- Why do we expect Measurement Layouts to help provide capability-oriented evaluation?
  - Capabilities are independent of the distribution of tasks instances.
  - Allows us to predict future performance at the instance-level. The Measurement Layout models reflect domain expertise.
    - Higher internal validity.
    - Enables predictions to be valid even out of distribution.

# Evaluating capabilities of LLMs

# Key challenges in building measurement layouts for LLMs: Performance data

- We need instance-level performance data
  - For each prompt need to have how the language model performed
    - Success/failure or if multiple choice which are the choices, which one did LLM pick and if correct or not
    - Addition example -- prompt: 'Sum of 112359321 and 23456543 is'; success: yes
- Some good examples include: HELM and BIG-bench

# Key challenges in building measurement layouts for LLMs: Meta-features/demands

- Many benchmarks don't include and/or don't have meta-features that describe each instance in the benchmark
  - How was the benchmark built?
  - What was considered when building the benchmark beyond just this a benchmark that tests X?
  - Are there any features that were explicitly manipulated?
- Addition: **number of digits in each number**; create new features: **average number of digits, number of carrying operations**, etc.

# Key challenges in building measurement layouts for LLMs: Capabilities

- What capabilities are tested by those benchmarks?
  - Addition: what are the abilities that are being tested/informed or that need to match demands such as **average number of digits**, **number of carry operations** etc.
- Are those capabilities characterizable and have some backing in cognitive science (or some other sciences)?
- How do these capabilities interact? Does having one capability compensate for not having some other capability test on the benchmark?

# Summary

- Need **instance-level** performance data
- Need **meta-features/demands** that describe each instance
- Need **capabilities** that are tested/informed by those demands



Let's build some measurement  
layouts

# Exploring Measurement Layouts in PyMC

All notebooks, data and presentation are at:

- <https://github.com/markotesic/Robust-Evaluation-of-Generative-AI/>

# Addition data set

# Addition dataset

- A simple benchmark testing LLMs abilities to add two numbers
  - E.g. prompt: The sum of 656050998910983832047 and 4871137 is
- 10 LLMs tested: GPT-3 Ada, GPT-3 Babbage, GPT-3 Curie, GPT-3 Davinci, text-davinci-001, text-davinci-002, text-davinci-003, GPT-3.5-turbo, GPT-4-0314, GPT-4-0613
- We will see how to create some of the **meta-features** for this benchmark and how to relate them to **capabilities** learnt by the measurement layout
- We will also train some common predictive models and baselines to compare to the performance of measurement layouts

# Addition dataset

agents/0\_intro\_rl.ipynb at master X +

← → C https://github.com/tensorflow/agents/blob/master/tf\_agents/colabs/0\_intro\_rl.ipynb

Search or jump to... / Pull requests Issues Marketplace Explore

tensorflow / agents Used by 22 Unwatch 69 Unstar 825 Fork 154

Octotree > Code Issues 60 Pull requests 24 Projects 0 Wiki Security Insights

Branch: master agents / tf\_agents / colabs / 0\_intro\_rl.ipynb Find file Copy path

anandijain tiny spelling fix a178fc2 on 11 Jun

3 contributors

138 lines (129 sloc) 9.22 KB Raw Blame History

Copyright 2018 The TF-Agents Authors.

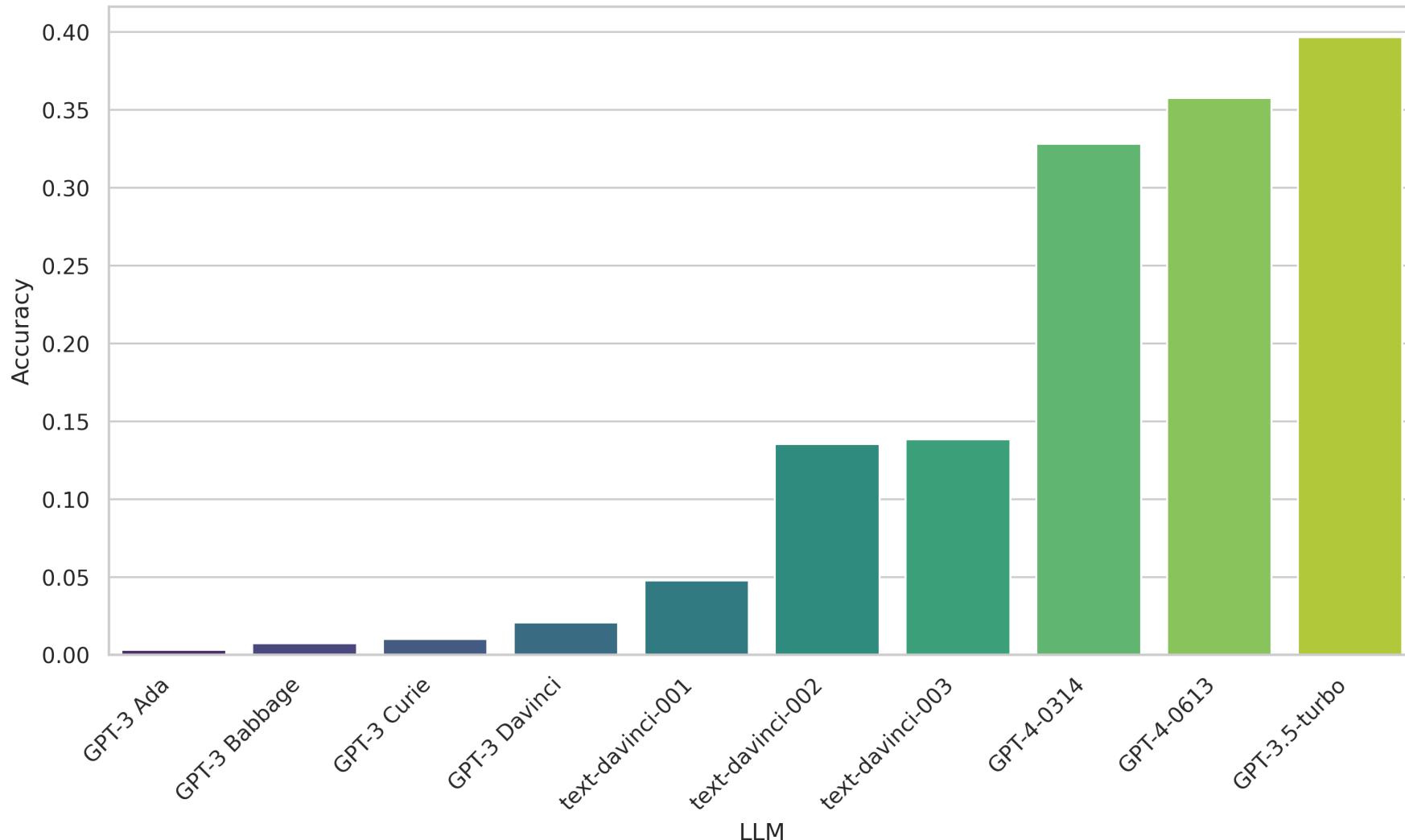
# Addition dataset: overview

- `instance_id`: ID of an instance.
  - `template_id`: ID of a prompt template. Each addition prompt was phrased in 8 different ways.
  - `llm`: The large language model tested.
  - `prompt`: The actual prompt.
  - `response`: The LLM's response.
  - `summand1`: The first number to be added.
  - `summand2`: The second number to be added.
  - `target`: The correct sum of the two summands.
  - `digits1`: The number of digits in the first summand. Between 1 and 30 digits
  - `digits2`: The number of digits in the second summand.
  - `min_digits`:  $\min(digits_1, digits_2)$ , i.e., the number of digits in the smaller summand.
  - `harm_mean`:  $2/(1/digits_1 + 1/digits_2)$ , i.e., the harmonic mean of the number of digits in the two summands.
  - `art_mean`:  $(digits_1 + digits_2)/2$ , i.e., the arithmetic mean of the number of digits in the two summands.
  - `max_digits`:  $\max(digits_1, digits_2)$ , i.e., the number of digits in the larger summand.
  - `carry`: The number of carrying operations required to add the two numbers.
  - `success`: Indicates whether the response is correct (1) or incorrect (0).
- Template 1: The sum of 656050998910983832047 and 4871137 is  
Template 2: Add: 656050998910983832047 + 4871137 =  
Template 3: Make the addition of 656050998910983832047 and 4871137.  
Template 4: By adding 656050998910983832047 and 4871137, the result is  
Template 5: If you add 656050998910983832047 and 4871137, you get  
Template 6: If you have 656050998910983832047 and 4871137, and you add them up, you get  
Template 7: Find the value of  $x + y$  when  $x=656050998910983832047$  and  $y=4871137$ .  
Template 8: Imagine you have two numbers, 656050998910983832047 and 4871137, and you added them together. What number would you get?

# Addition dataset

| success                 | 0    | 1    |
|-------------------------|------|------|
| <b>11m</b>              |      |      |
| <b>GPT-3 Ada</b>        | 3286 | 10   |
| <b>GPT-3 Babbage</b>    | 3272 | 24   |
| <b>GPT-3 Curie</b>      | 3263 | 33   |
| <b>GPT-3 Davinci</b>    | 3228 | 68   |
| <b>GPT-3.5-turbo</b>    | 1989 | 1307 |
| <b>GPT-4-0314</b>       | 2214 | 1082 |
| <b>GPT-4-0613</b>       | 2117 | 1179 |
| <b>text-davinci-001</b> | 3139 | 157  |
| <b>text-davinci-002</b> | 2850 | 446  |
| <b>text-davinci-003</b> | 2840 | 456  |

# Addition dataset: aggregate accuracy



# Meta-features/demands

- `digits1`: The number of digits in the first summand.
  - `digits2`: The number of digits in the second summand.
  - `min_digits`:  $\min(digits_1, digits_2)$ , i.e., the number of digits in the smaller summand.
  - `harm_mean`:  $2/(1/digits_1 + 1/digits_2)$ , i.e., the harmonic mean of the number of digits in the two summands.
  - `art_mean`:  $(digits_1 + digits_2)/2$ , i.e., the arithmetic mean of the number of digits in the two summands.
  - `max_digits`:  $\max(digits_1, digits_2)$ , i.e., the number of digits in the larger summand.
  - `carry`: The number of carrying operations required to add the two numbers.
- 
- Do these meta-features capture everything that could we could potentially think of when it comes to addition of two number? What are some of the things that make the addition of two number ‘difficult’?
  - Size of the two numbers
  - Number of carrying operations
  - Can we have lots of carrying operations but the additions is still ‘easy’?

# Meta-features/demands

- Standard deviation of digits in the two numbers (digit variety):

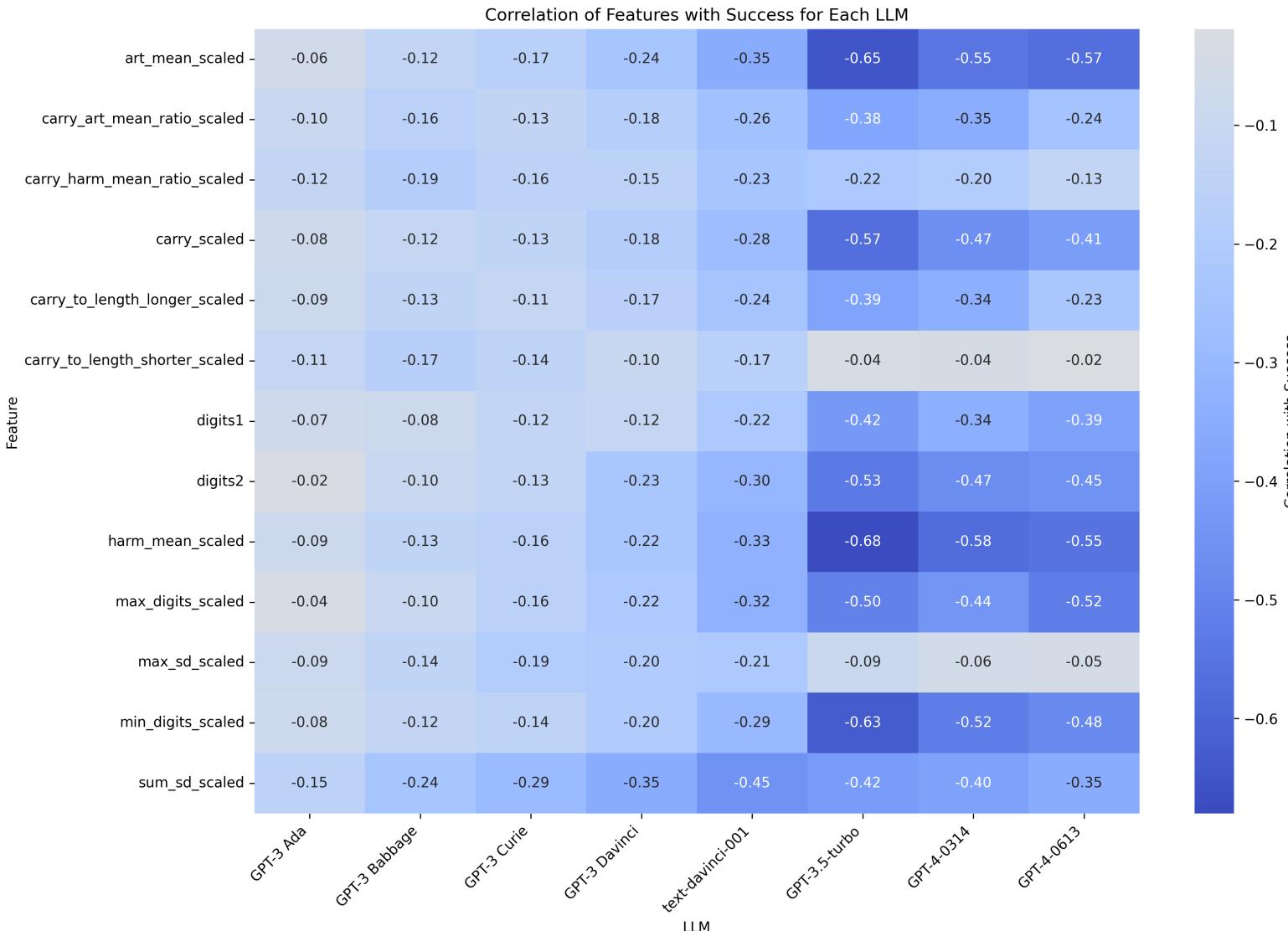
```
# Calculate max and sum of digit standard deviations
addition_data['max_sd'] = addition_data.apply(lambda x: max(num_sd(x['summand1']), num_sd(x['summand2'])), axis=1)
addition_data['sum_sd'] = addition_data.apply(lambda x: num_sd(x['summand1']) + num_sd(x['summand2']), axis=1)
```

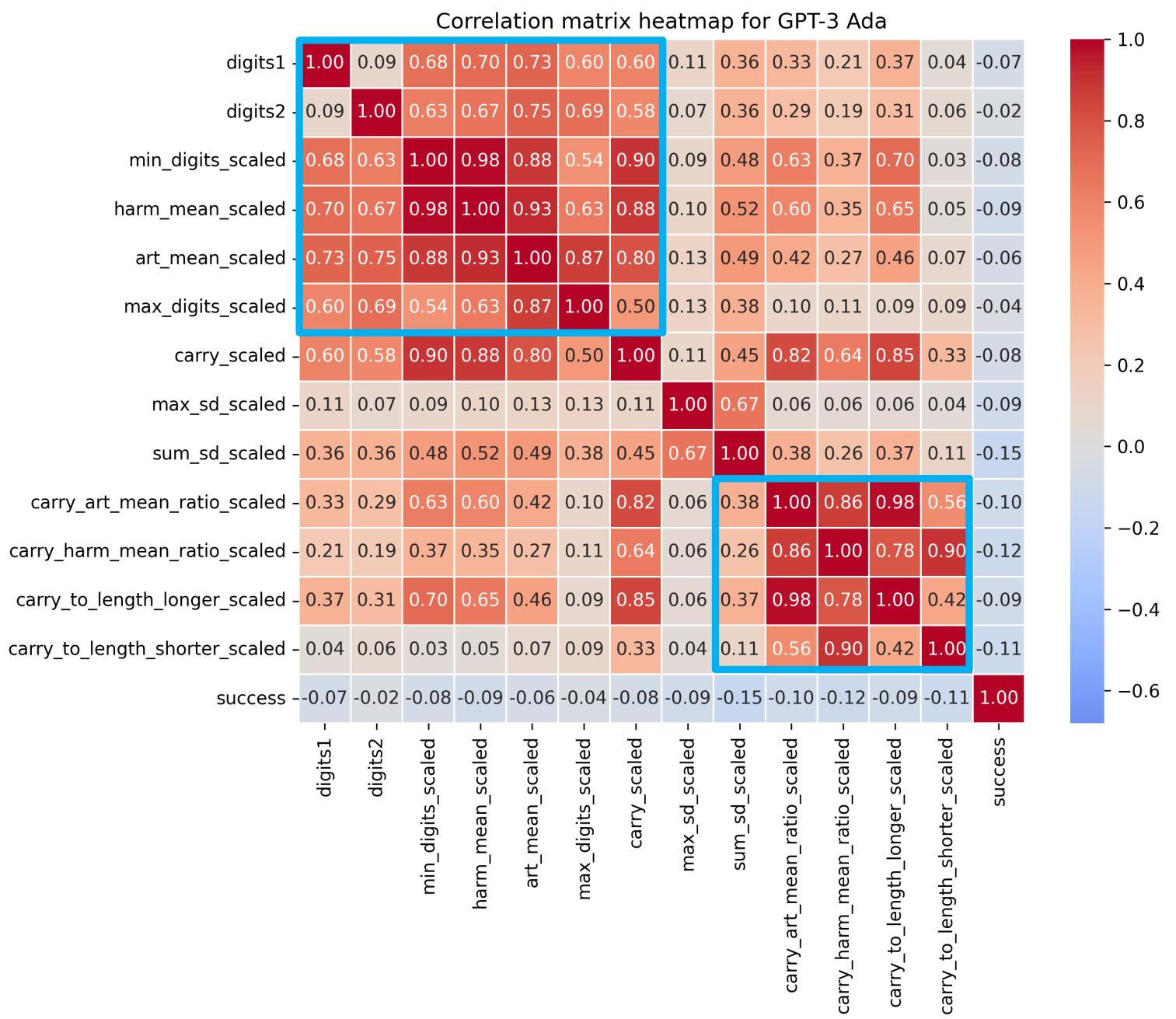
- Calculate the ratio of carry to size:

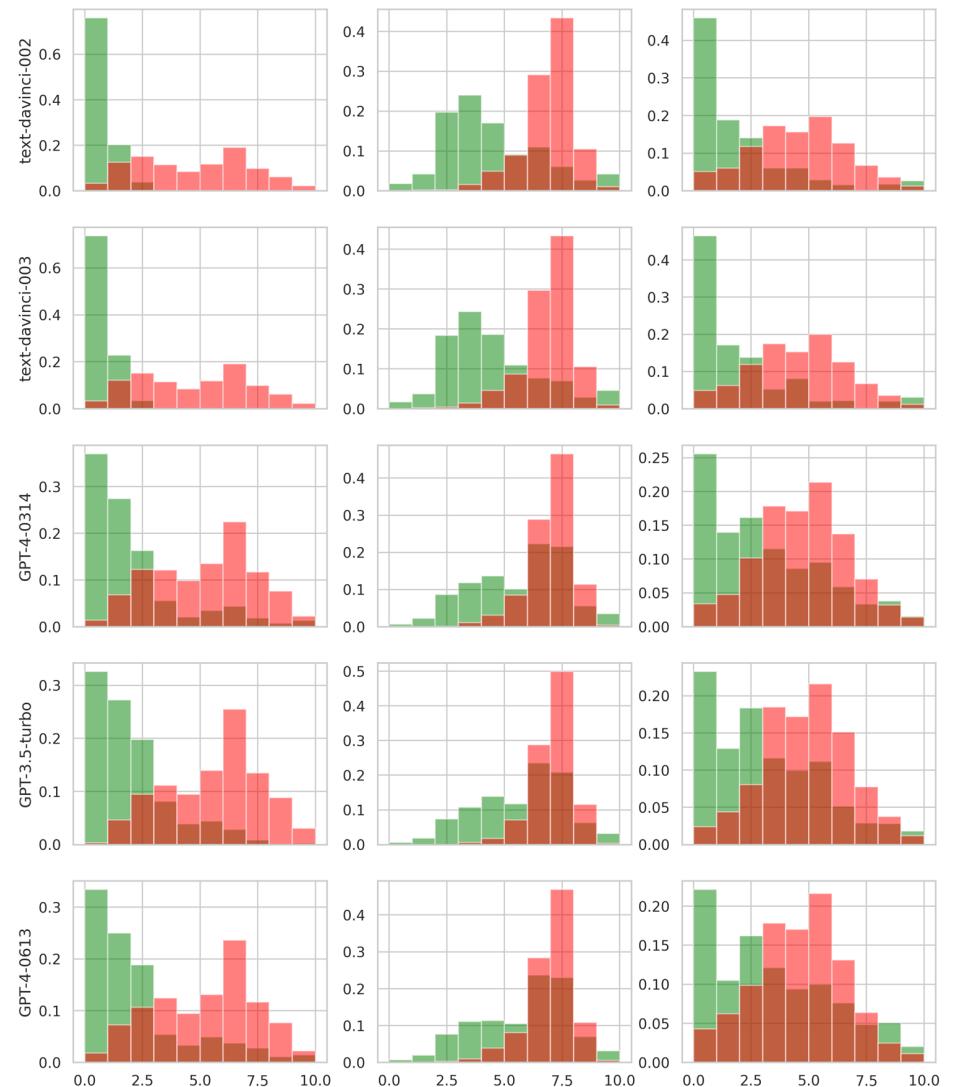
```
# Calculate ratio of carry and harmonic mean columns
addition_data['carry_harm_mean_ratio'] = addition_data['carry'] / addition_data['harm_mean']
addition_data['carry_art_mean_ratio'] = addition_data['carry'] / addition_data['art_mean']
addition_data['carry_to_length_longer'] = addition_data['carry'] / addition_data['max_digits']
addition_data['carry_to_length_shorter'] = addition_data['carry'] / addition_data['min_digits']
```

- Scale data: min-max normalization to range between 0 and 10

# Meta-features/demands: correlations with success



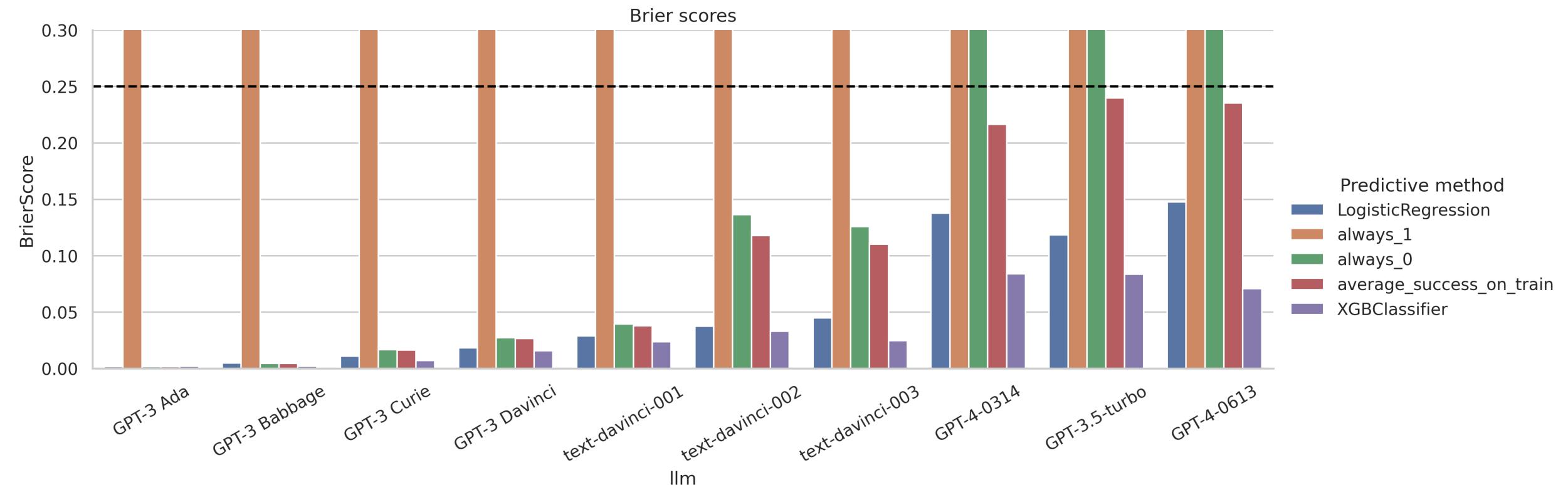




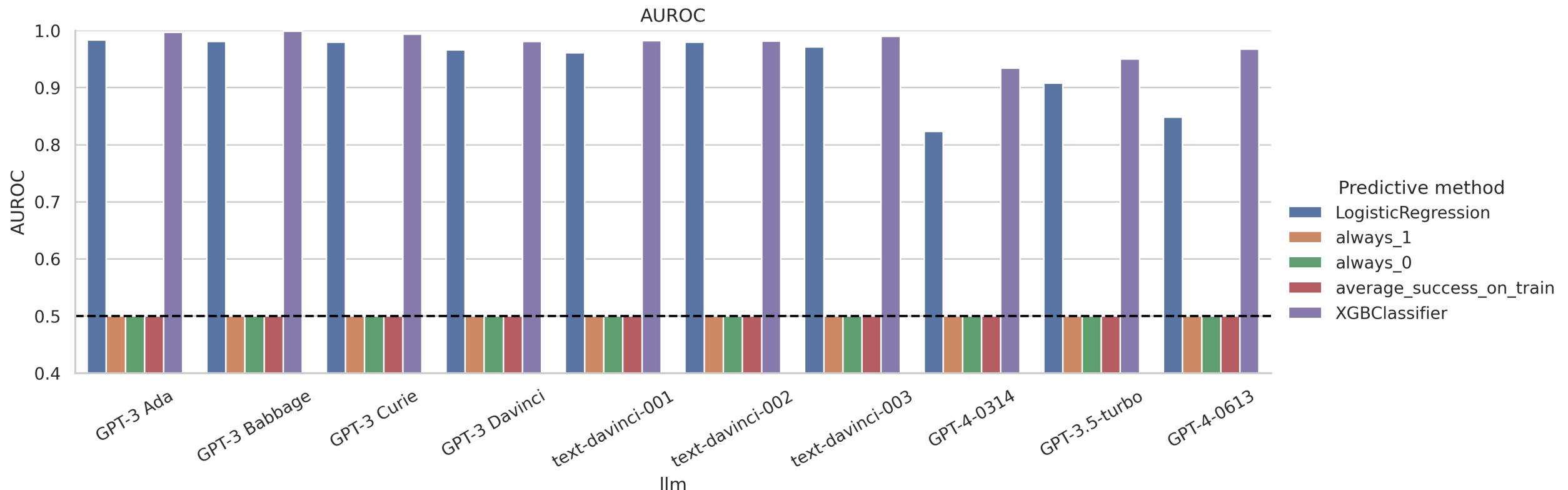
# Common predictive models

- Build Logistic regression and XGBoost classifier for each LLM to predict success using the 3 meta-features (harm\_mean, carry\_art\_mean\_ratio, sum\_sd)
- Add baselines:
  - Predict success with probability 1 for each instance
  - Predict success with probability 0 for each instance
  - Predict success with probability equal to aggregate probability of success on the training set (i.e. the proportion of successes on the training data per LLM)
- Evaluate the performance of the predictive models using AUROC and Brier score

# Common predictive models



# Common predictive models



# Capabilities

- 3 meta-features/demands: harm\_mean, carry\_art\_mean\_ratio, sum\_sd
- Capabilities that are informed by/measured/match those demands:
  - harm\_carry ↔ size ability
  - carry\_art\_mean\_ratio ↔ carry ability (ability to perform addition that require carry to various extents)
  - sum\_sd ↔ digit variety ability
- Simple case of one demand, one ability (in general doesn't have to be as we've seen in the previous session)
- Abilities non-compensatory (to an extent)

# Measurement layouts

```
✓ def setupSumModel(data):
    m = pm.Model()
    ✓ with m:

        # Define abilities and their priors
        sizeAbility = pm.Normal("sizeAbility", 0, 6)
        carryAbility = pm.Normal("carryAbility", 0, 6)
        digitVarietyAbility = pm.Normal("digitVarietyAbility", 0, 6)

        # Define environment variables as MutableData
        digitsMean = pm.MutableData("digitsMean", data['harm_mean'].values)
        ratioCarry = pm.MutableData("ratioCarry", data['carry_art_mean_ratio'].values)
        sumSD = pm.MutableData("sumSD", data['sum_sd'].values)

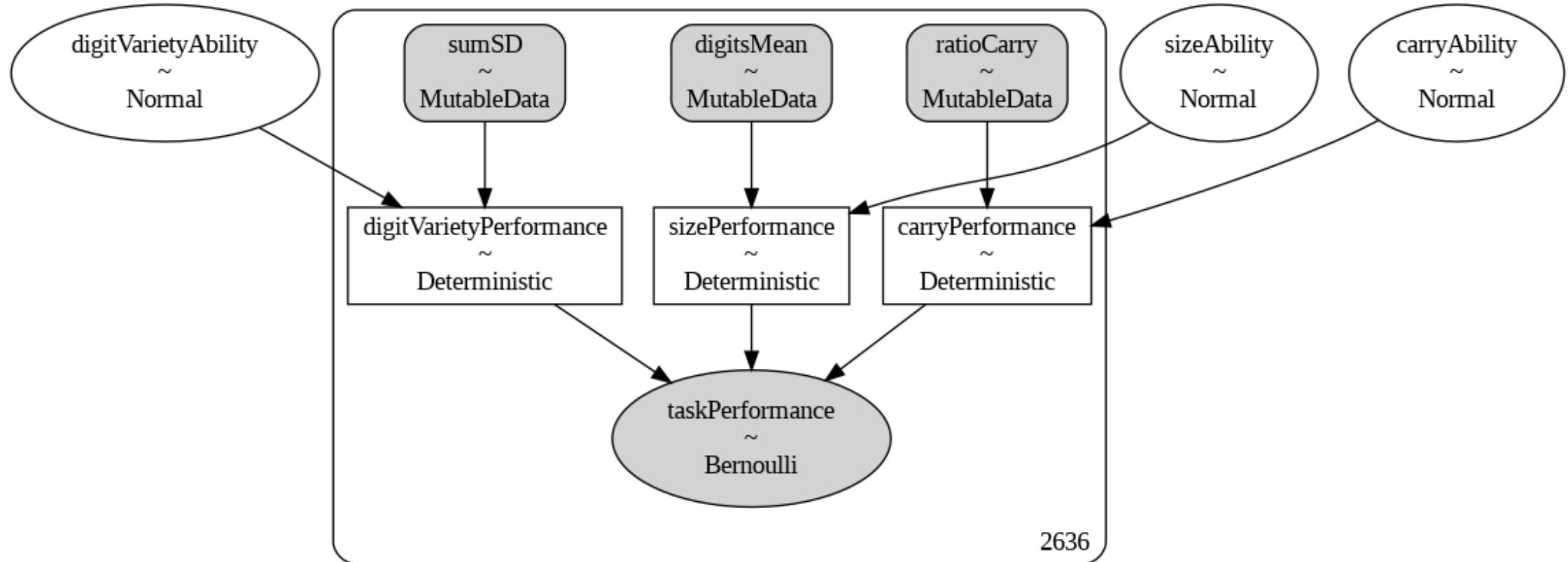
        # Margins / partial performance
        sizePerformance = pm.Deterministic("sizePerformance", logistic_general(margin(sizeAbility,digitsMean),min=0,max=12,c=0,p=0.99))
        carryPerformance = pm.Deterministic("carryPerformance", logistic_general(margin(carryAbility,ratioCarry),min=0,max=12,c=0,p=0.99))
        digitVarietyPerformance = pm.Deterministic("digitVarietyPerformance", logistic_general(margin(digitVarietyAbility,sumSD),min=0,max=12,c=0,p=0.99))

        #Non-compensatory interaction
        finalPerformance = sizePerformance * carryPerformance * digitVarietyPerformance

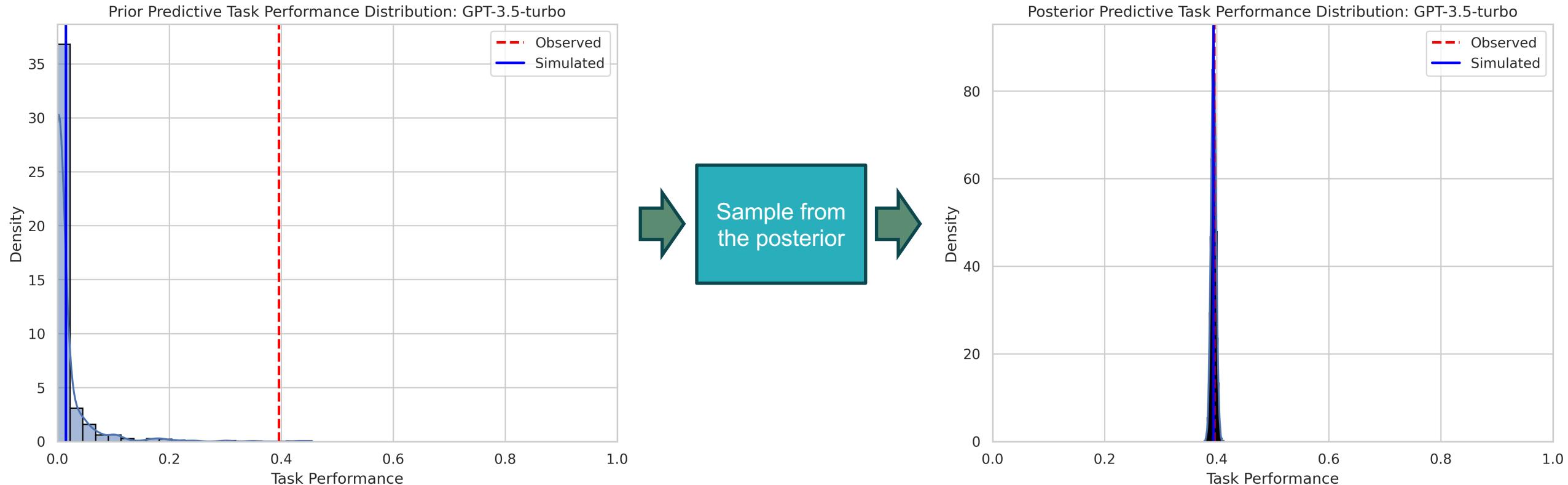
        taskPerformance = pm.Bernoulli("taskPerformance", finalPerformance, observed = data['success'])

    return m
```

# Measurement layouts

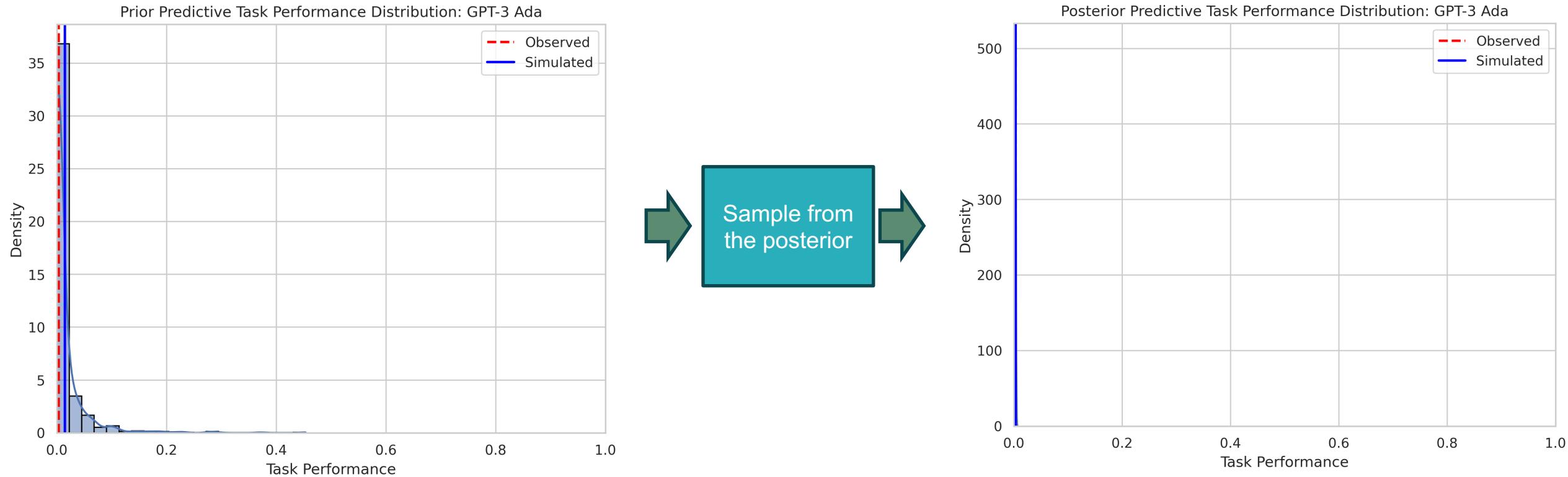


# Measurement layouts: prior and posterior predictive checks



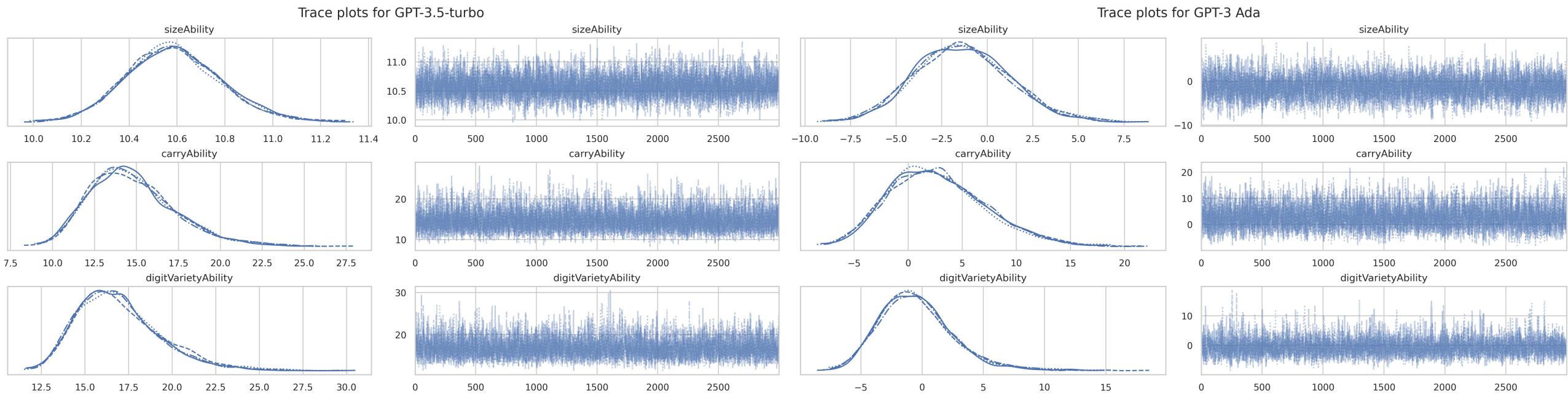
$$\theta^{sim} \sim p(\theta) \quad y^{sim} \sim p(y|\theta^{sim})$$

# Measurement layouts: prior and posterior predictive checks

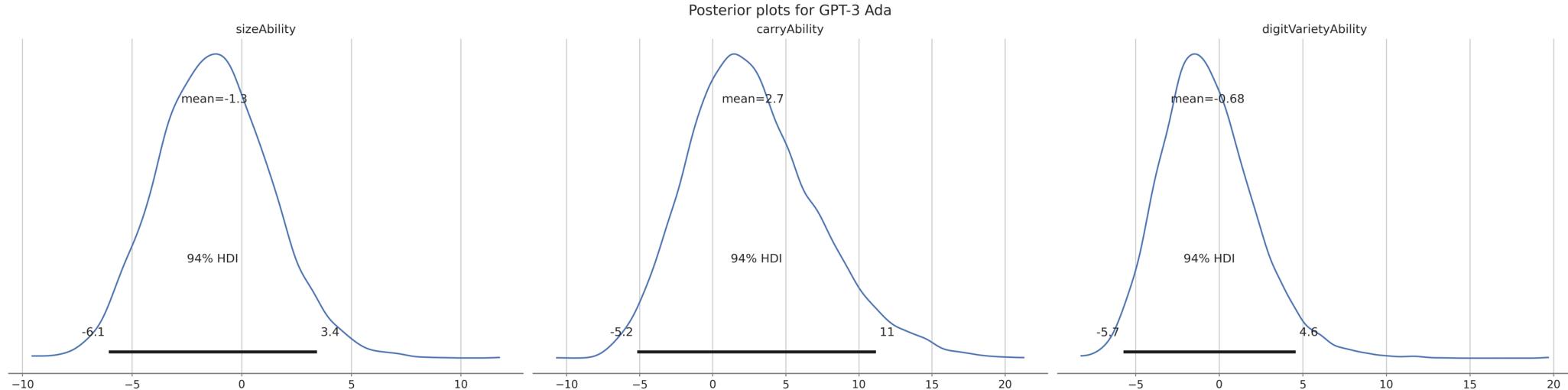
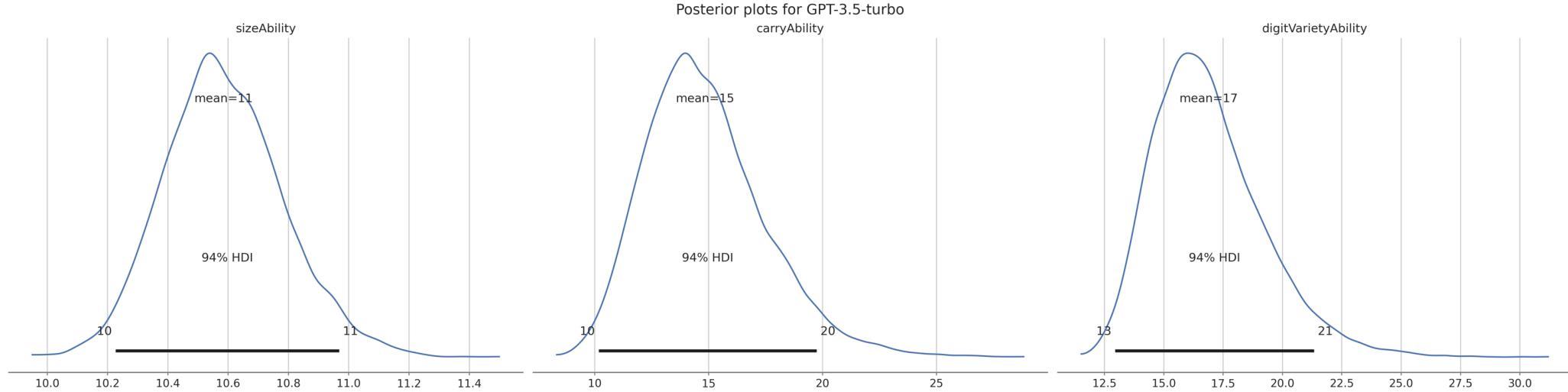


$$\theta^{sim} \sim p(\theta) \quad y^{sim} \sim p(y|\theta^{sim})$$

# Measurement layouts: Markov chains

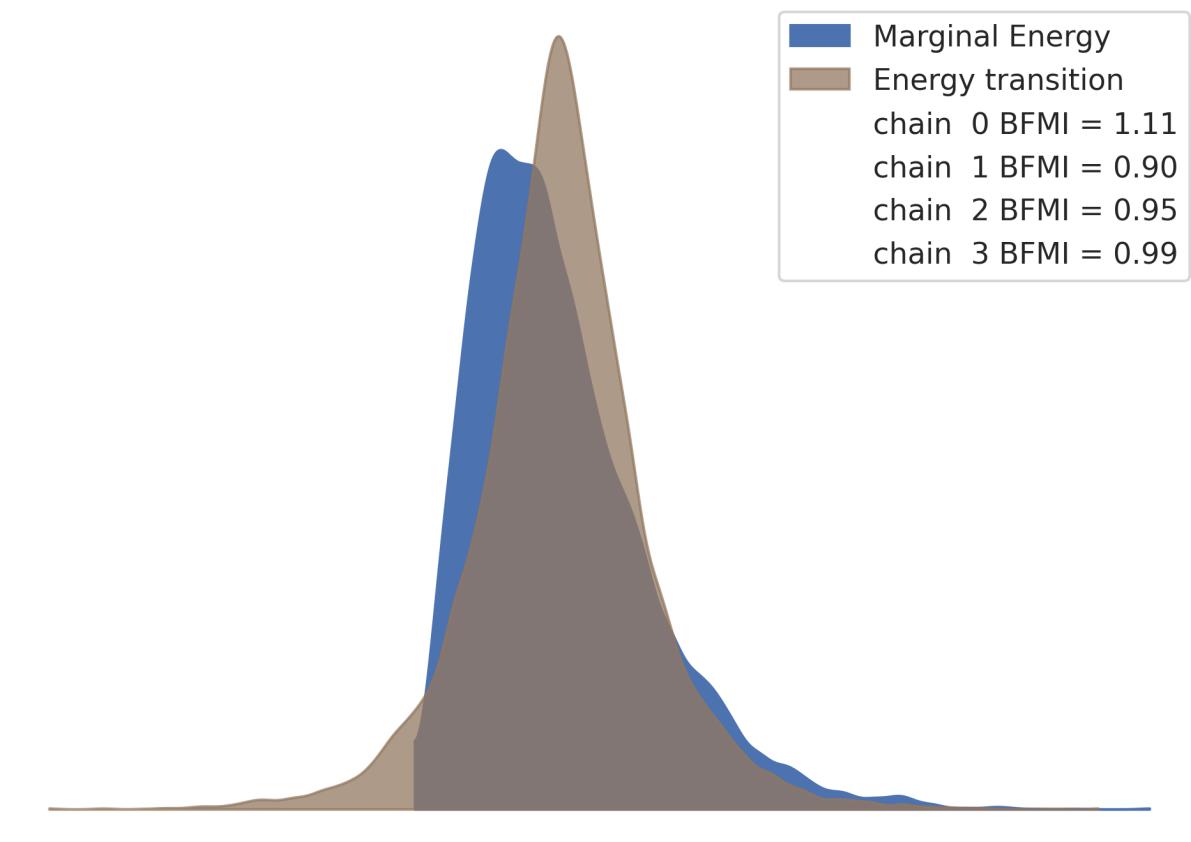


# Measurement layouts: posterior plots



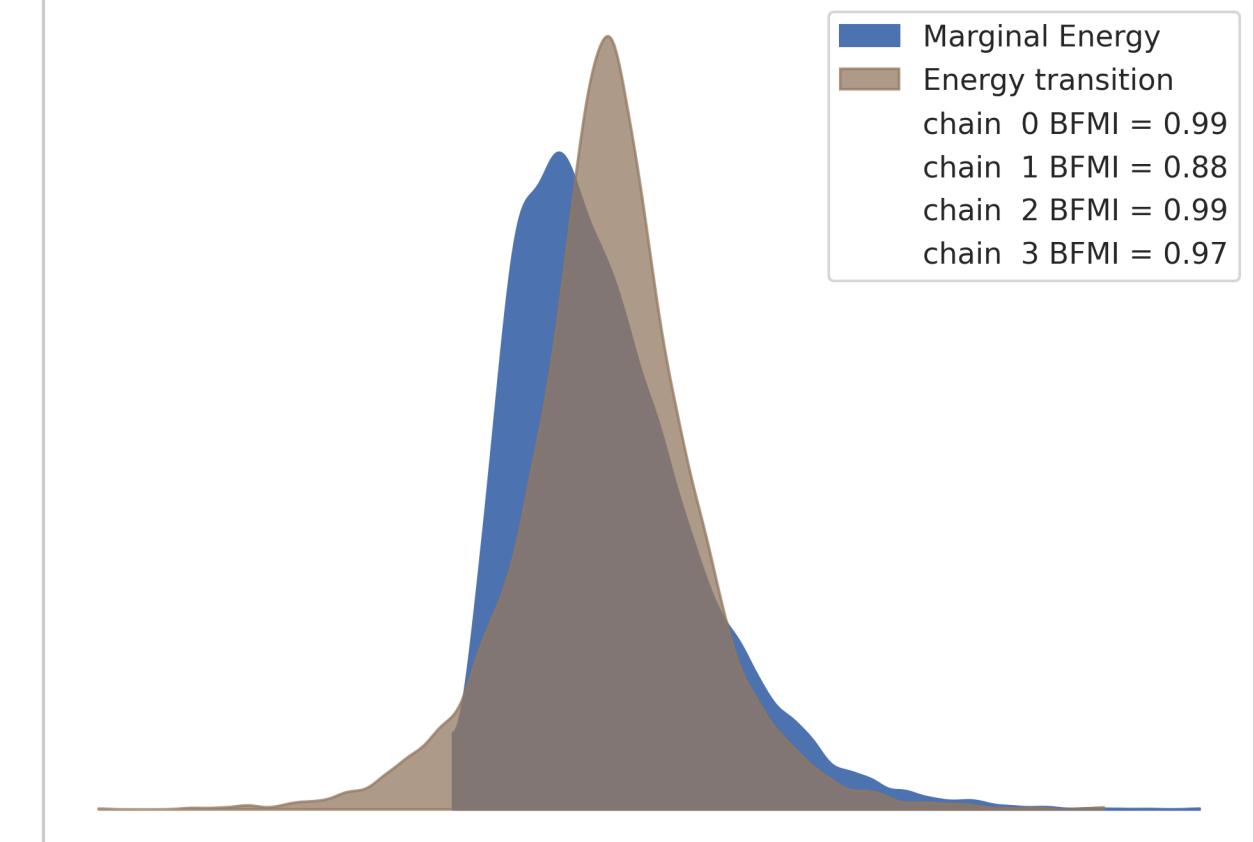
# Measurement layouts: convergence metrics

Energy plot for GPT-3.5-turbo



Marginal Energy  
Energy transition  
chain 0 BFMI = 1.11  
chain 1 BFMI = 0.90  
chain 2 BFMI = 0.95  
chain 3 BFMI = 0.99

Energy plot for GPT-3 Ada



Marginal Energy  
Energy transition  
chain 0 BFMI = 0.99  
chain 1 BFMI = 0.88  
chain 2 BFMI = 0.99  
chain 3 BFMI = 0.97

# Measurement layouts: convergence metrics

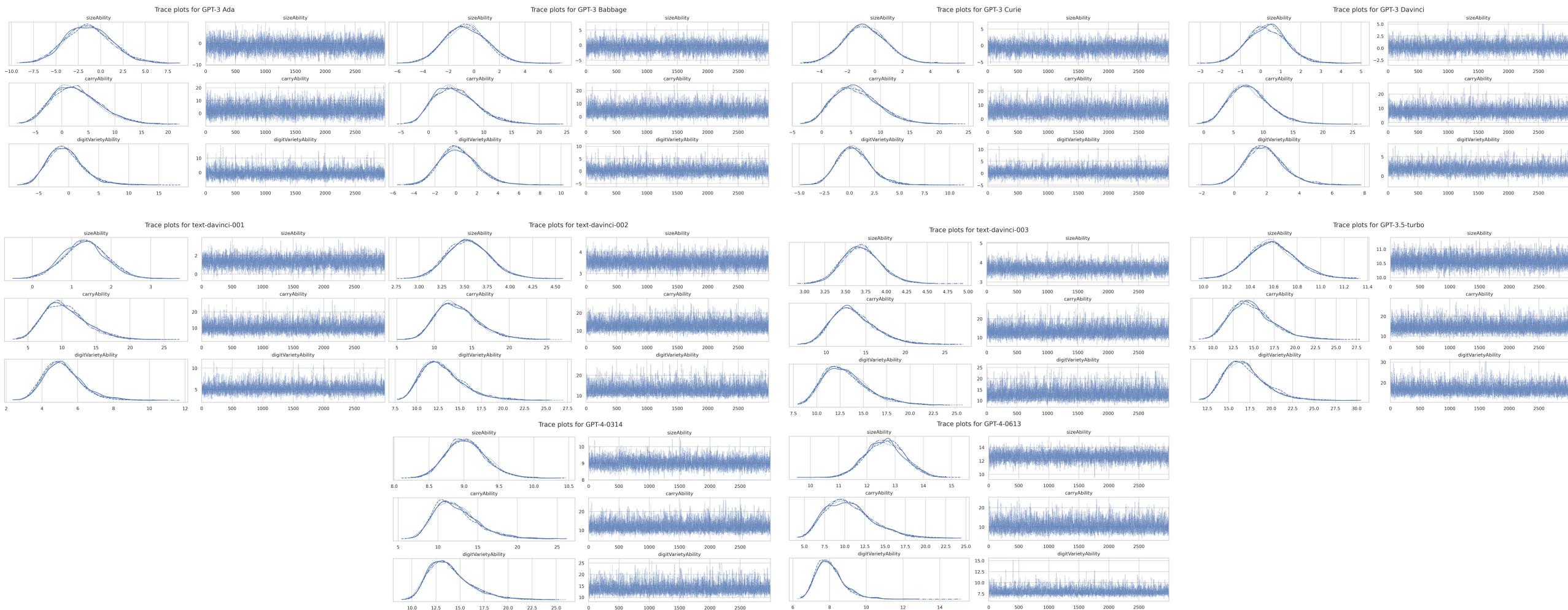
Summary statistics: GPT-3.5-turbo

|                            | mean   | sd    | hdi_3% | hdi_97% | ess_bulk | ess_tail | r_hat |
|----------------------------|--------|-------|--------|---------|----------|----------|-------|
| <b>sizeAbility</b>         | 10.585 | 0.198 | 10.227 | 10.968  | 5746.943 | 5346.377 | 1.001 |
| <b>carryAbility</b>        | 14.839 | 2.652 | 10.175 | 19.743  | 6126.873 | 4941.861 | 1.001 |
| <b>digitVarietyAbility</b> | 16.887 | 2.336 | 12.949 | 21.330  | 4821.880 | 4500.382 | 1.001 |

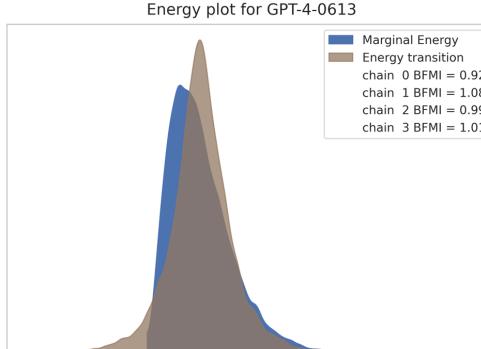
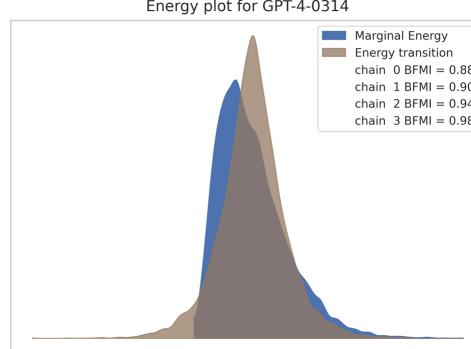
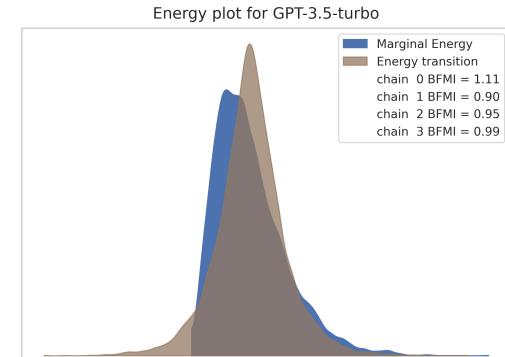
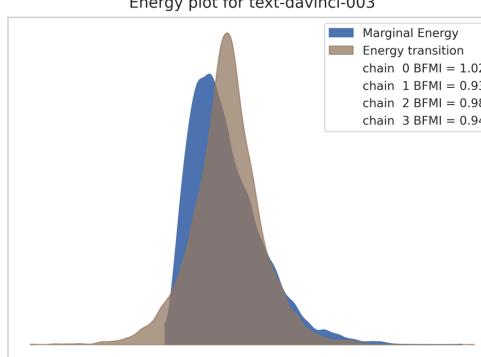
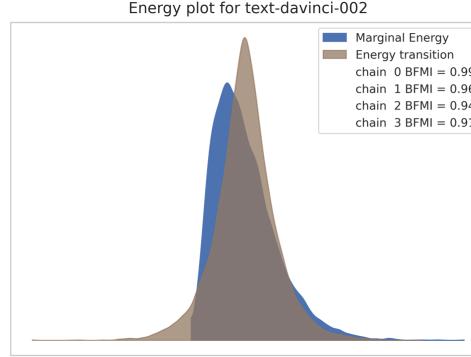
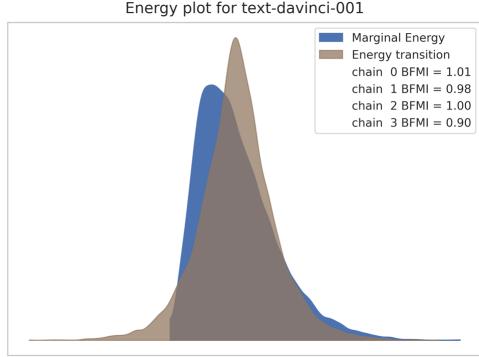
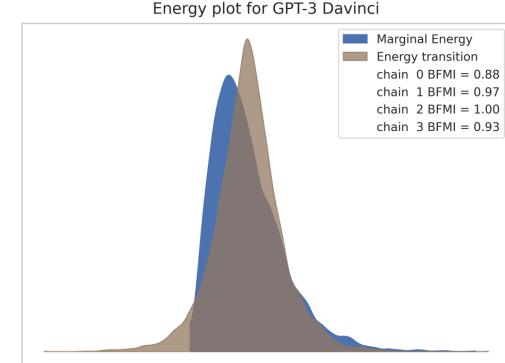
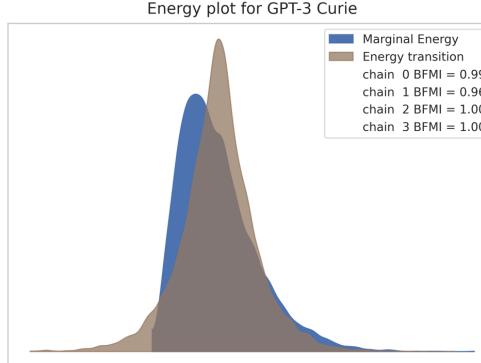
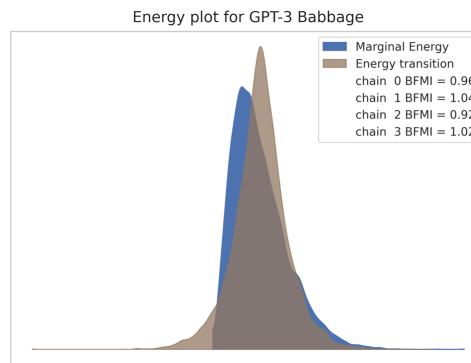
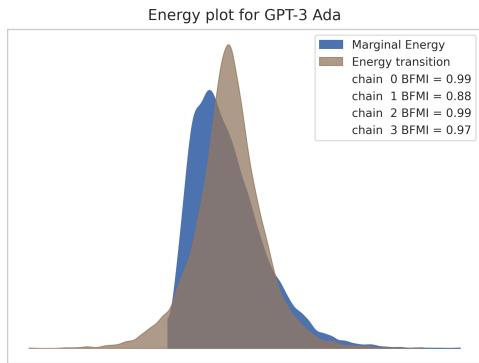
Summary statistics: GPT-3 Ada

|                            | mean   | sd    | hdi_3% | hdi_97% | ess_bulk | ess_tail | r_hat |
|----------------------------|--------|-------|--------|---------|----------|----------|-------|
| <b>sizeAbility</b>         | -1.252 | 2.569 | -6.075 | 3.425   | 4289.046 | 4371.565 | 1.001 |
| <b>carryAbility</b>        | 2.739  | 4.415 | -5.183 | 11.158  | 4434.910 | 4728.066 | 1.001 |
| <b>digitVarietyAbility</b> | -0.685 | 2.822 | -5.737 | 4.564   | 4432.195 | 4482.631 | 1.000 |

# Build measurement layouts for each LLM

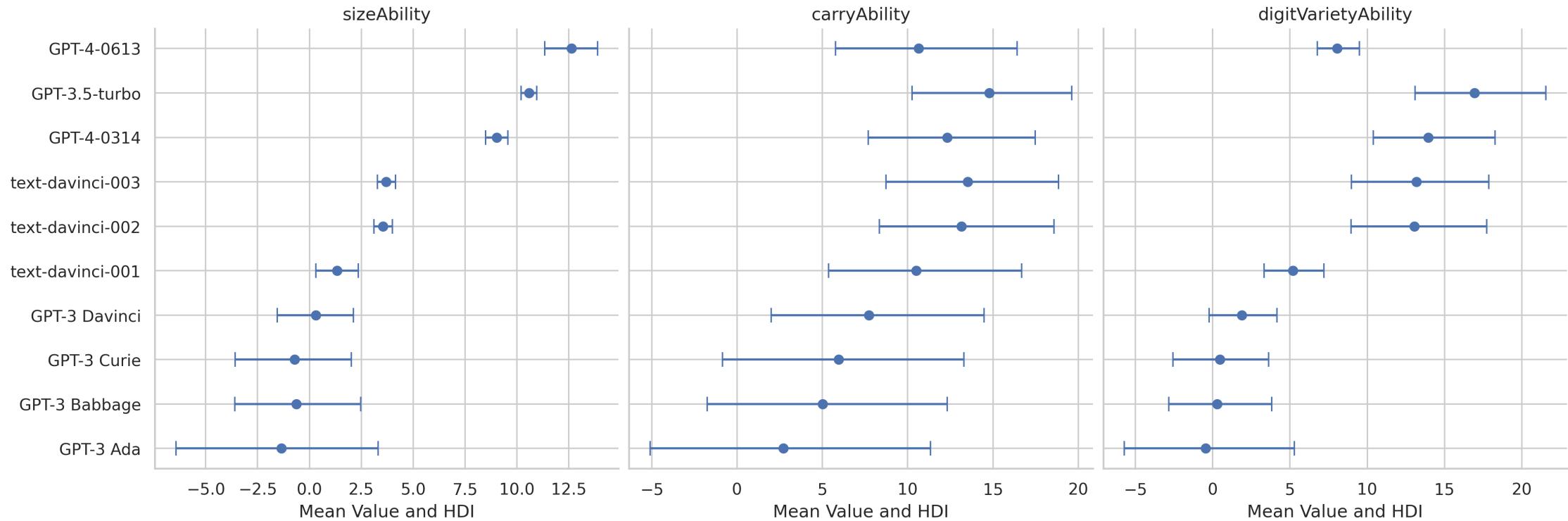


# Build measurement layouts for each LLM

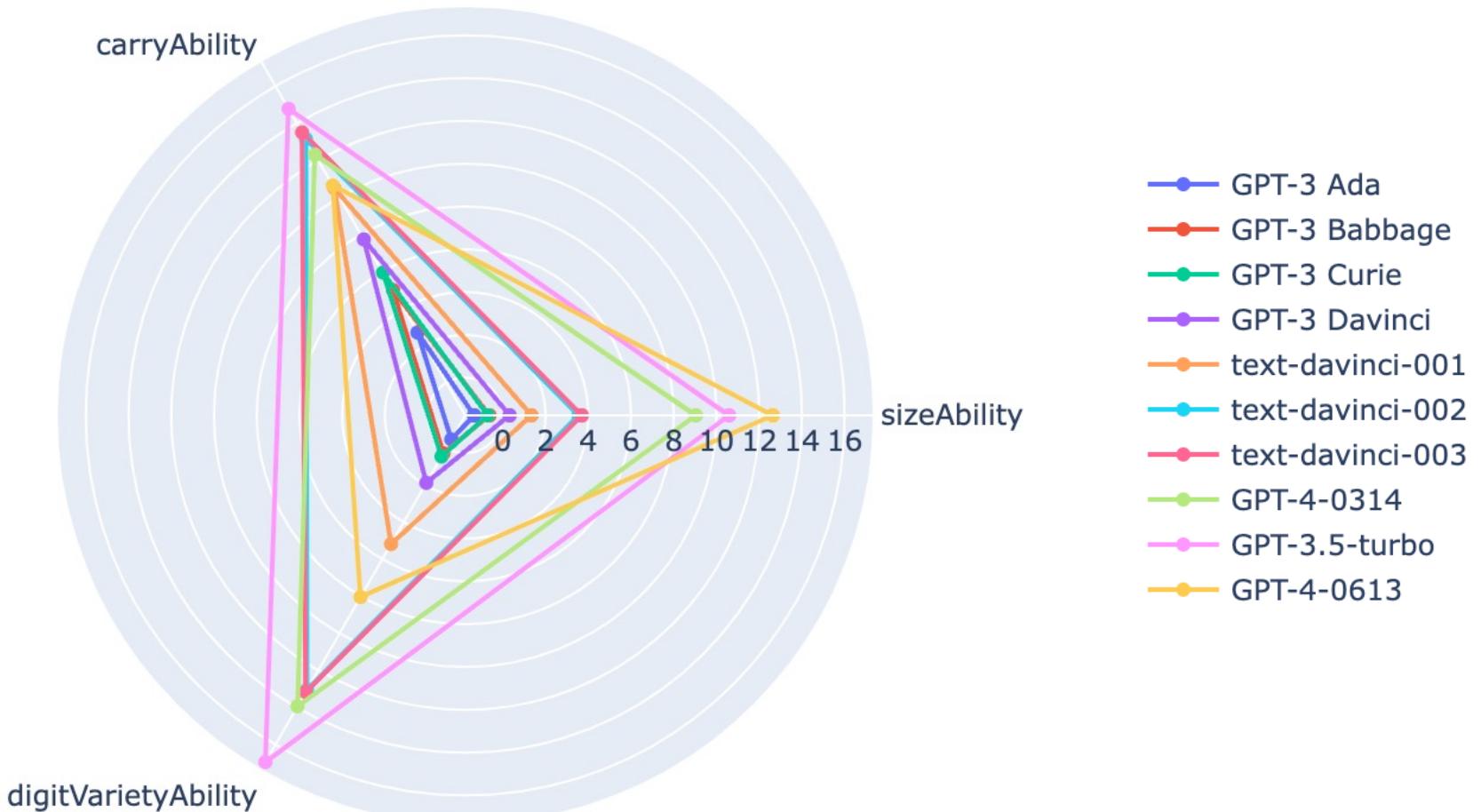


| LLM              | Ability             | ess_bulk | ess_tail | r_hat |
|------------------|---------------------|----------|----------|-------|
| GPT-3 Ada        | sizeAbility         | 4498.985 | 4433.078 | 1.001 |
|                  | carryAbility        | 5223.856 | 5411.293 | 1.000 |
|                  | digitVarietyAbility | 4185.243 | 3874.632 | 1.000 |
| GPT-3 Babbage    | sizeAbility         | 4267.570 | 4854.430 | 1.002 |
|                  | carryAbility        | 4414.343 | 4259.494 | 1.001 |
|                  | digitVarietyAbility | 4199.861 | 4260.734 | 1.001 |
| GPT-3 Curie      | sizeAbility         | 4326.282 | 4850.581 | 1.000 |
|                  | carryAbility        | 4959.339 | 4330.477 | 1.001 |
|                  | digitVarietyAbility | 4645.460 | 4959.820 | 1.000 |
| GPT-3 Davinci    | sizeAbility         | 4422.249 | 4616.055 | 1.001 |
|                  | carryAbility        | 5512.185 | 4857.278 | 1.000 |
|                  | digitVarietyAbility | 4336.236 | 4424.901 | 1.000 |
| text-davinci-001 | sizeAbility         | 3978.109 | 4799.510 | 1.002 |
|                  | carryAbility        | 4826.458 | 4416.341 | 1.001 |
|                  | digitVarietyAbility | 4042.979 | 3612.911 | 1.002 |
| text-davinci-002 | sizeAbility         | 5163.584 | 5466.039 | 1.000 |
|                  | carryAbility        | 5874.716 | 4886.958 | 1.002 |
|                  | digitVarietyAbility | 4165.493 | 3805.924 | 1.001 |
| text-davinci-003 | sizeAbility         | 4737.856 | 4706.722 | 1.001 |
|                  | carryAbility        | 5922.350 | 4829.169 | 1.002 |
|                  | digitVarietyAbility | 3918.935 | 4393.516 | 1.000 |
| GPT-4-0314       | sizeAbility         | 3589.105 | 4003.518 | 1.001 |
|                  | carryAbility        | 4773.185 | 3827.450 | 1.001 |
|                  | digitVarietyAbility | 3006.980 | 2697.055 | 1.001 |
| GPT-3.5-turbo    | sizeAbility         | 5088.514 | 5161.921 | 1.001 |
|                  | carryAbility        | 5734.960 | 5162.670 | 1.000 |
|                  | digitVarietyAbility | 4352.323 | 4721.511 | 1.001 |
| GPT-4-0613       | sizeAbility         | 3951.369 | 4506.992 | 1.001 |
|                  | carryAbility        | 3837.764 | 4128.004 | 1.001 |
|                  | digitVarietyAbility | 3581.884 | 3887.807 | 1.001 |

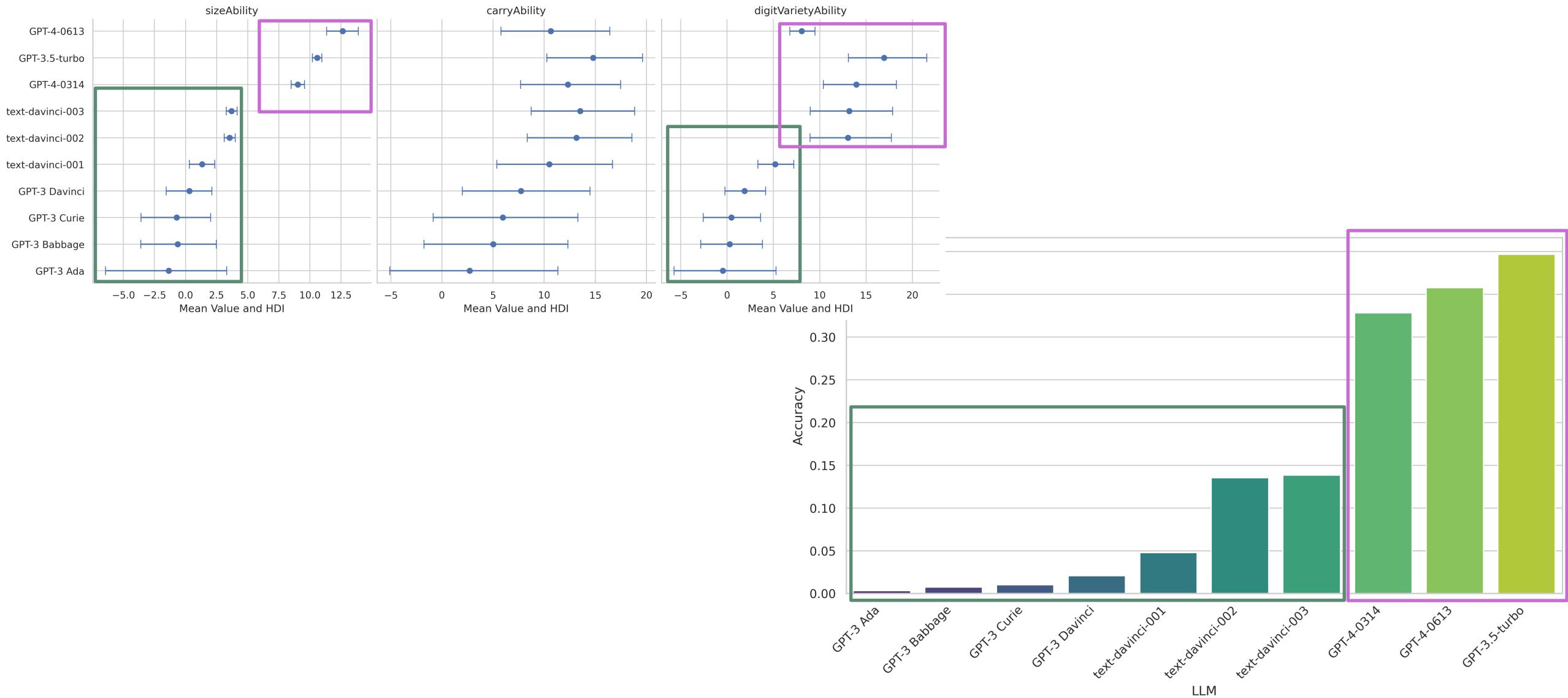
# Build measurement layouts for each LLM



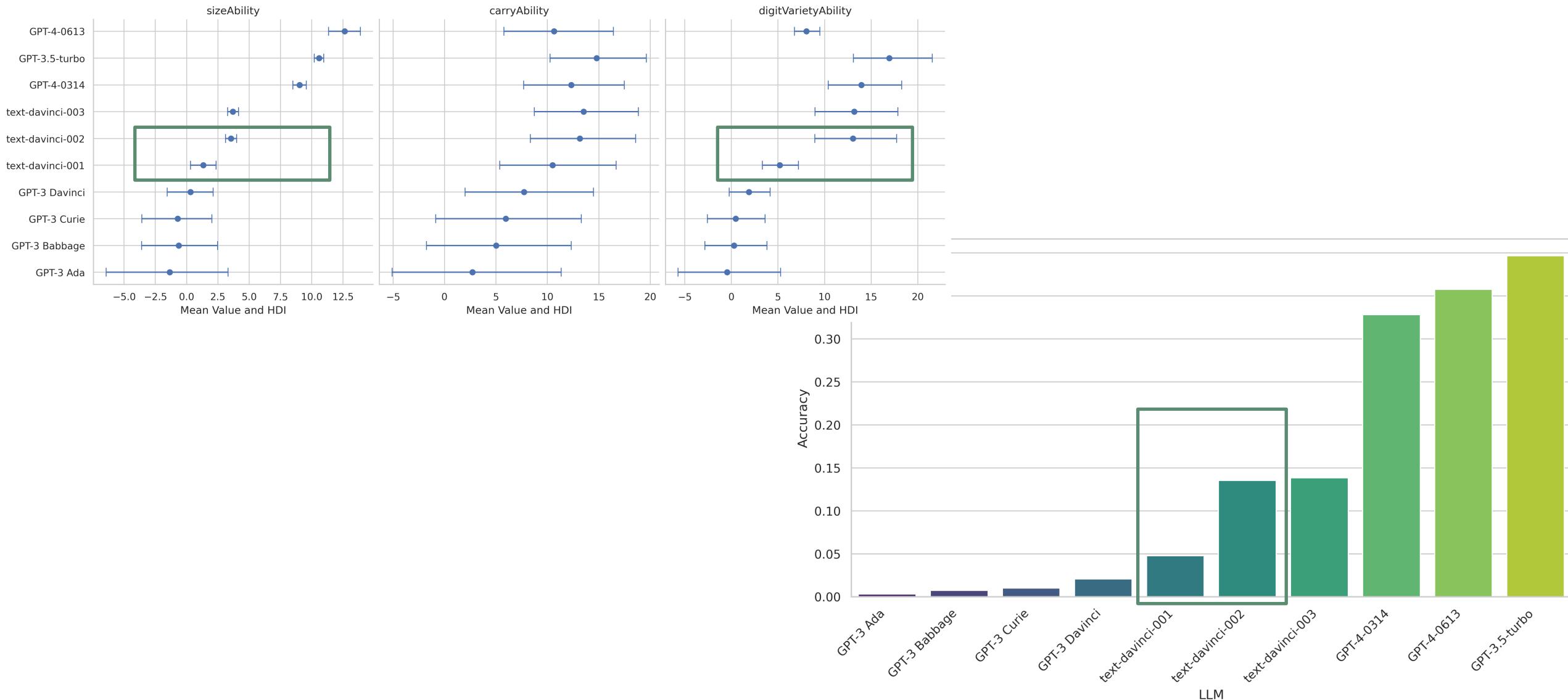
# Build measurement layouts for each LLM



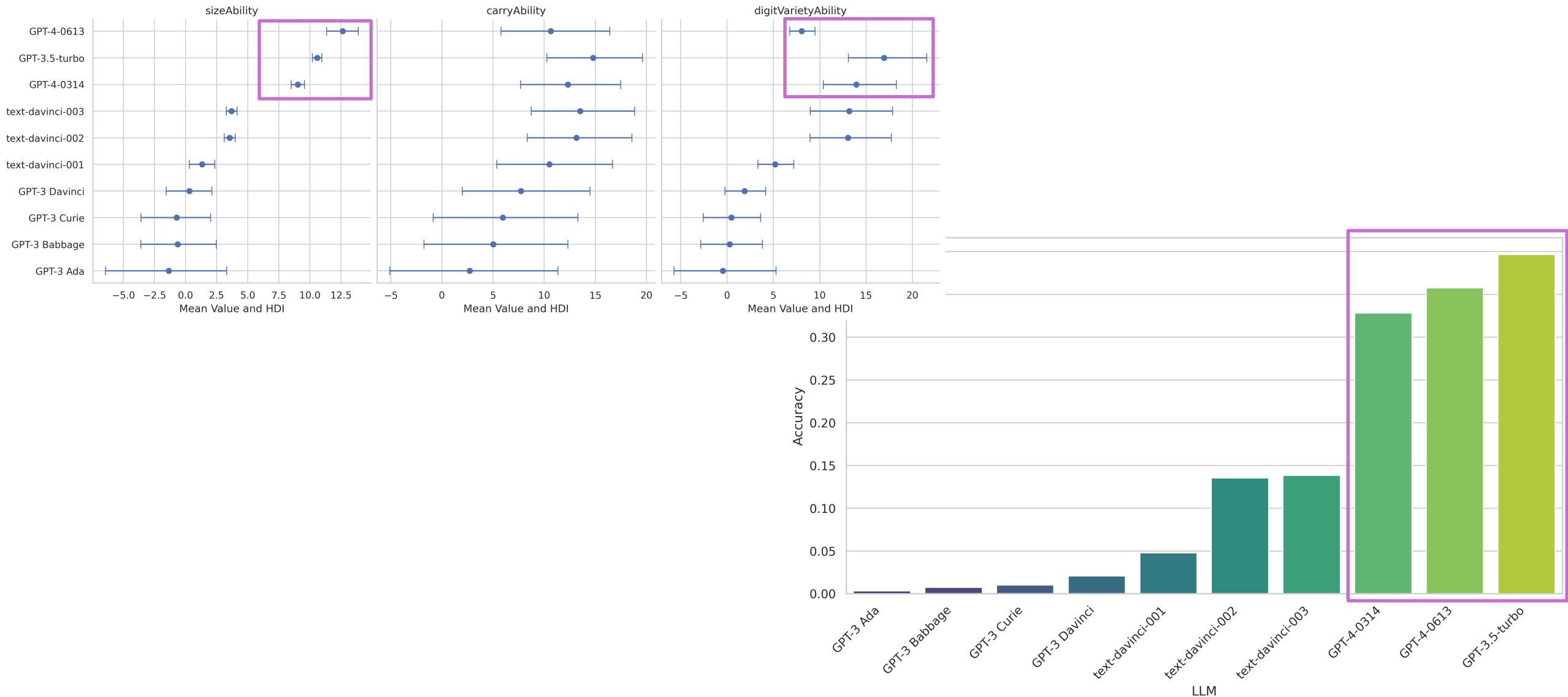
# Build measurement layouts for each LLM



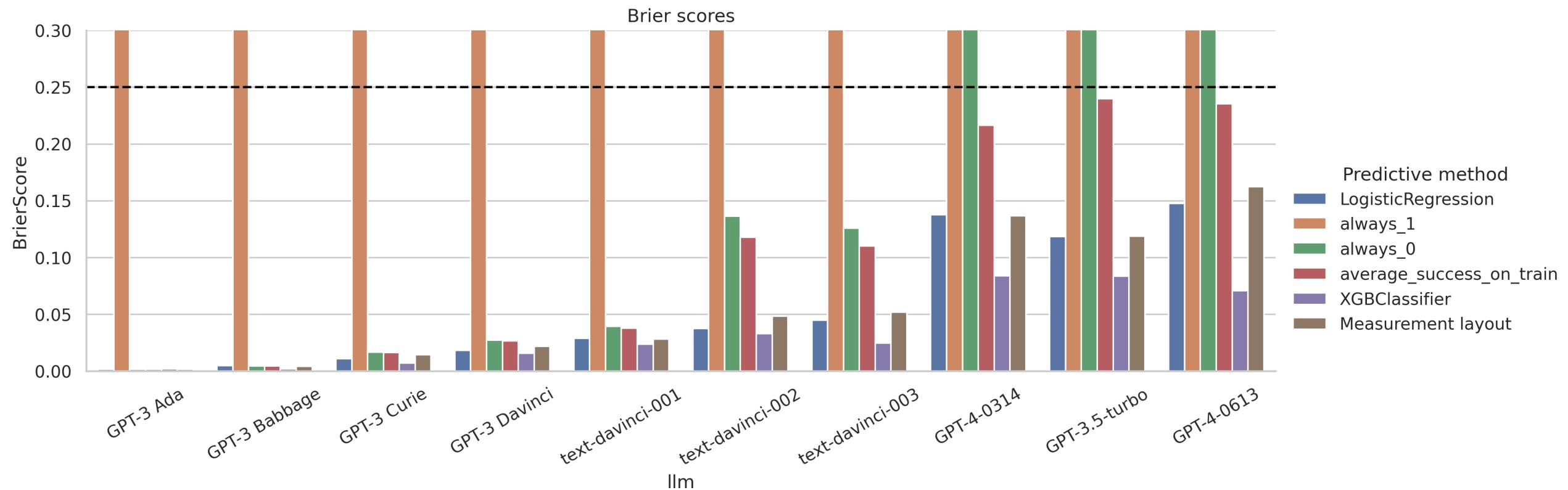
# Build measurement layouts for each LLM



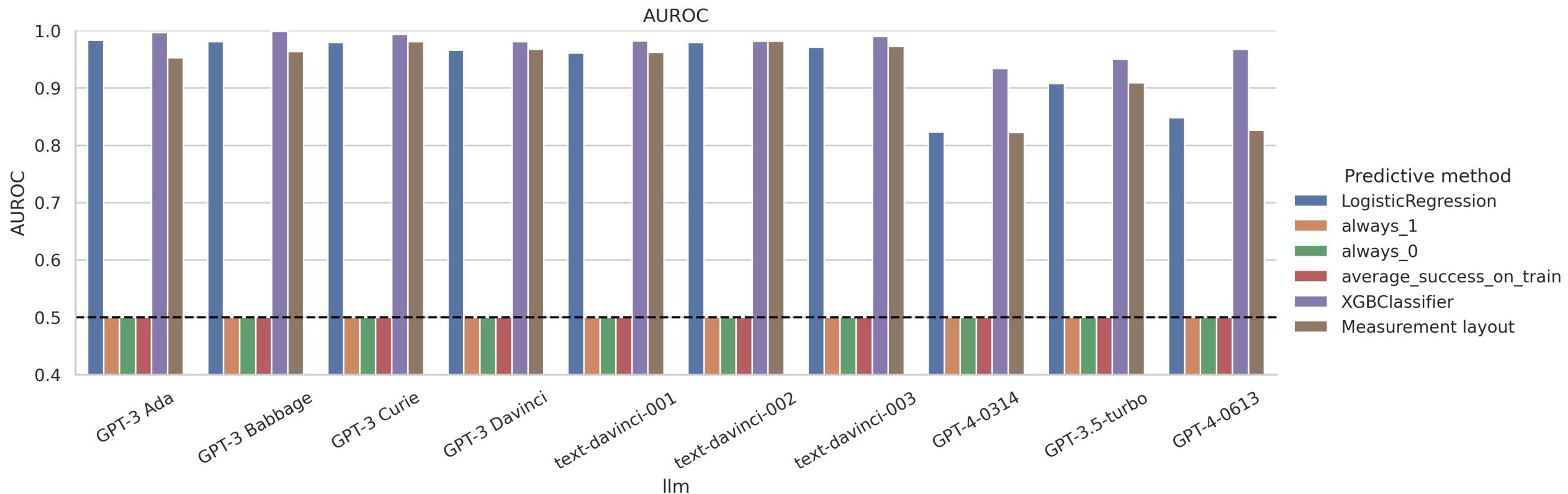
# Build measurement layouts for each LLM



# Prediction accuracy



# Prediction accuracy



# Take-aways

- Capability is not average performance
  - Capability allows to predict performance, at the instance level and even OOD
- System Capabilities and Task Demands are related through a “margin”
- Measurement layouts capture domain knowledge and intuitions
  - Need instance level data
  - Meta-features demands
  - Capabilities that are meant to meet the demands
- Backward inference: estimate capabilities for one single system
- Forward inference: infer performance for new task instances.



LEVERHULME CENTRE FOR THE  
FUTURE OF INTELLIGENCE

## With Thanks to...



Kozzy  
Voudouris



José  
Hernández-  
Orallo



John Burden



Lorenzo  
Pacchiardi



Lucy  
Cheke

UNIVERSITY OF  
CAMBRIDGE

CFI LEVERHULME CENTRE FOR THE  
FUTURE OF INTELLIGENCE



E·S·R·C  
ECONOMIC & SOCIAL  
RESEARCH COUNCIL



GoodAI



WBAI Whole Brain  
Architecture Initiative

LEVERHULME  
TRUST

