

Benchmarking Large Language Models for Bio-Image Analysis Code Generation

Robert Haase, Christian Tischer, Jean-Karim Hériché, Nico Scherf

<https://www.biorxiv.org/content/10.1101/2024.04.19.590278v3>

Introduction

In the life-sciences,

- many people need to write code for analysing data.
 - They may not have coding skills.
 - They use ChatGPT etc to generate code.
 - The risk is that they use code that does wrong things.
-
- **Question:** How good are large language models (LLMs) e.g. for code generation specifically to analyse microscopy image data?

Introduction: Code generation using LLMs

- Codex model (Chen et al 2021) was major milestone in the code generation context – it's the predecessor of Github Copilot
- Chen et al 2021 also delivered a common benchmark: HumanEval

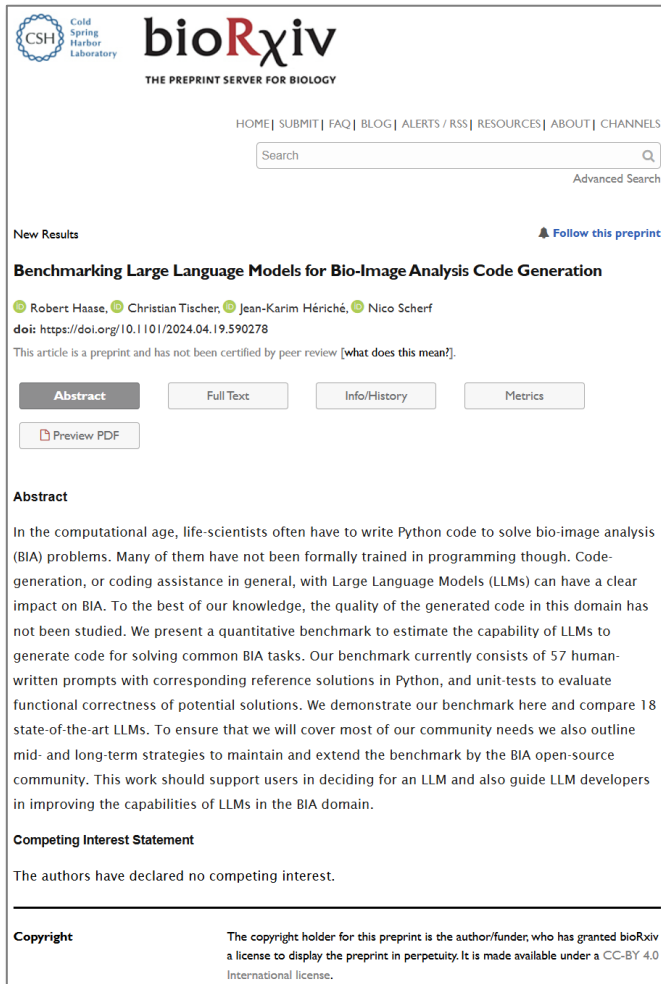
```
def incr_list(l: list):  
    """Return list with elements incremented by 1.  
    >>> incr_list([1, 2, 3])  
    [2, 3, 4]  
    >>> incr_list([5, 3, 5, 2, 3, 3, 9, 0, 123])  
    [6, 4, 6, 3, 4, 4, 10, 1, 124]  
    """  
    return [i + 1 for i in l]
```

Prompt

Reference solution

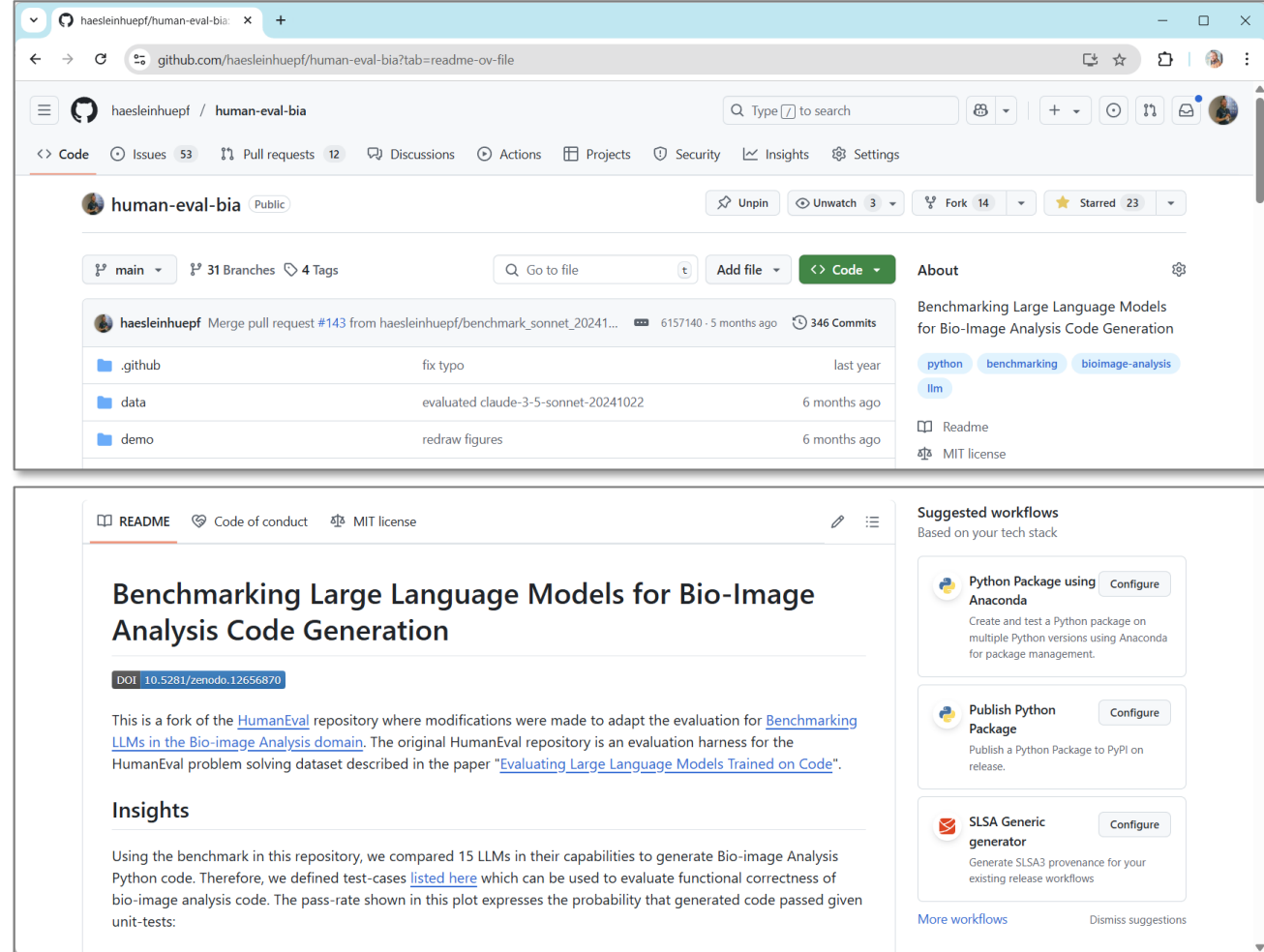
Figure 2. Three example problems from the HumanEval dataset, where the probabilities that a single sample from Codex-12B passes unit tests are 0.9, 0.17, and 0.005. The prompt provided to the model is shown with a white background, and a successful model-generated completion is shown in a yellow background. Though not a guarantee for problem novelty, all problems were hand-written and not programmatically copied from existing sources. Random problems and samples can be found in Appendix B.

Benchmarking LLMs for Bio-image Analysis



The screenshot shows the bioRxiv preprint page for the paper "Benchmarking Large Language Models for Bio-Image Analysis Code Generation". The page includes the bioRxiv logo, navigation links, a search bar, and the paper's title. The authors listed are Robert Haase, Christian Tischler, Jean-Karim Hériché, and Nico Scherf. The DOI is <https://doi.org/10.1101/2024.04.19.590278>. The abstract states that the paper presents a quantitative benchmark to estimate the capability of LLMs to generate code for solving common BIA tasks. The competing interest statement declares no competing interest. The copyright notice indicates that the preprint is available under a CC-BY 4.0 International license.

Preprint



The screenshot shows the GitHub repository page for "human-eval-bia" by haesleinhuepf. The repository is public and has 31 branches and 4 tags. It includes a table of files: ".github" (fix typo, last year), "data" (evaluated claude-3-5-sonnet-20241022, 6 months ago), and "demo" (redraw figures, 6 months ago). The README section is visible, showing the title "Benchmarking Large Language Models for Bio-Image Analysis Code Generation" and the DOI [10.5281/zenodo.12656870](https://doi.org/10.5281/zenodo.12656870). The README text describes the repository as a fork of the HumanEval repository, adapted for the evaluation of LLMs in the Bio-image Analysis domain. The Insights section mentions that the benchmark compared 15 LLMs in their capabilities to generate Bio-image Analysis Python code. The right sidebar shows suggested workflows for Python Package using Anaconda, Publish Python Package, and SLSA Generic generator.

Code

Methods: test cases

Example test-case inspired by HumaEval (Chen et al 2021)

```
[1]: def workflow_segmentation_measurement_summary(image):  
    """  
    This function implements a workflow consisting of these steps:  
    * threshold intensity input image using Otsu's method  
    * label connected components  
    * measure area of the labeled objects  
    * determine mean area of all objects  
    """  
  
    import skimage  
    import numpy as np  
    binary_image = image > skimage.filters.threshold_otsu(image)  
    label_image = skimage.measure.label(binary_image)  
    stats = skimage.measure.regionprops(label_image)  
    areas = [s.area for s in stats]  
    return np.mean(areas)
```

Prompt

Reference
solution

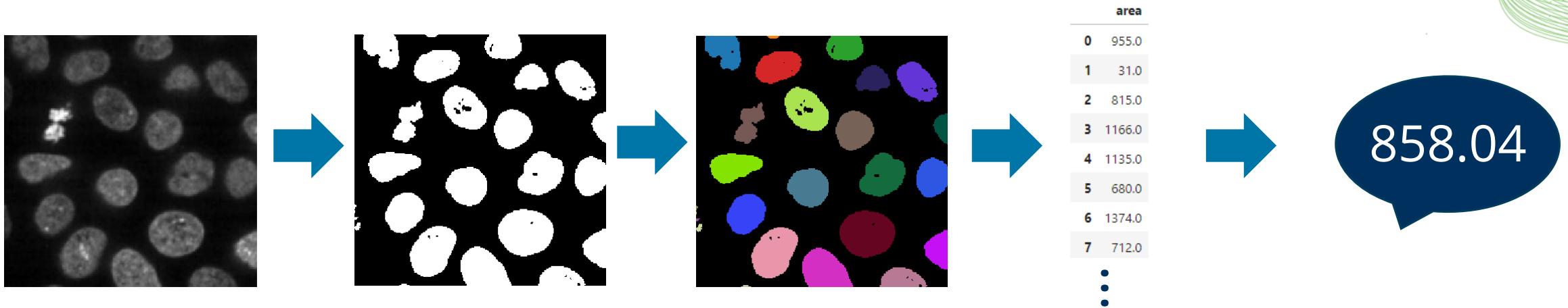
```
[2]: def check(candidate):  
    import numpy as np  
  
    assert candidate(np.asarray([  
        [0,0,0,0,0],  
        [1,1,1,0,0],  
        [1,1,1,0,0],  
        [1,1,0,0,0],  
        [0,0,0,0,0],  
    ])) == 8
```

Unit test
(excerpt)

Unit test pass rate ->
functional correctness

Methods: test cases

Use case: segment the image and measure the average area of objects.



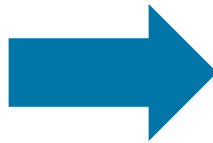
Unit-test pass-rate (n=10):

	reference	gpt-4-turbo-2024-04-09	Claude-3-opus-20240229	gpt-4-1106-preview	gpt-3.5-turbo-1106	gemini-pro	codellama
workflow_segmentation_measurement_summary	1.0	0.9	1.0	0.8	0.5	0.5	0.1

Methods: test cases

Use-case: compute the correlation matrix

	a	b	c	d	e
0	1.600000	0.100000	1.600000	1.700000	1.700000
1	2.300000	0.200000	2.300000	2.400000	2.400000
2	2.600000	0.300000	2.600000	2.400000	2.400000
3	3.700000	0.300000	3.700000	3.600000	3.600000
4	3.400000	0.400000	3.400000	3.500000	3.500000
5	3.900000	0.400000	3.900000	3.900000	3.900000
6	4.300000	0.400000	4.300000	4.400000	4.400000
7	4.300000	0.500000	4.300000	4.200000	4.200000
8	4.000000	0.500000	4.000000	4.100000	4.100000
9	5.100000	0.500000	5.100000	5.000000	5.000000
10	5.200000	0.600000	5.200000	5.100000	5.100000
11	5.300000	0.600000	5.300000	5.400000	5.400000
12	5.500000	0.600000	5.400000	5.600000	5.600000








	a	b	c	d	e
a	1.000000	0.949504	0.999775	0.995800	0.995800
b	0.949504	1.000000	0.949594	0.946039	0.946039
c	0.999775	0.949594	1.000000	0.995001	0.995001
d	0.995800	0.946039	0.995001	1.000000	1.000000
e	0.995800	0.946039	0.995001	1.000000	1.000000

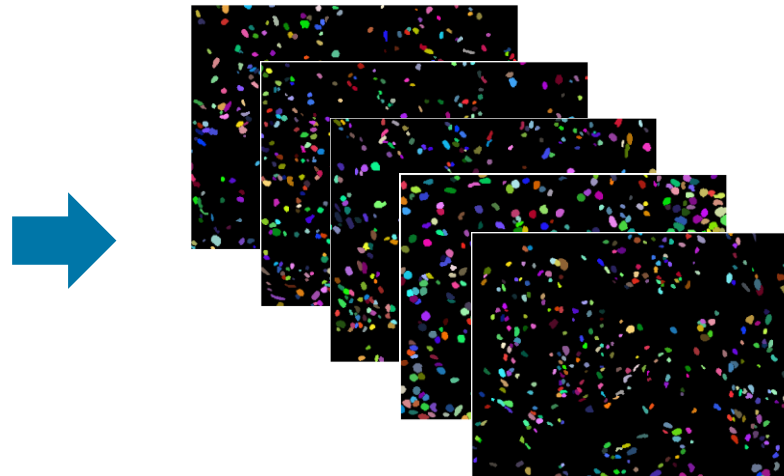
Unit-test pass-rate (n=10):

	reference	gpt-4-turbo-2024-04-09	Claude-3-opus-20240229	gpt-4-1106-preview	gpt-3.5-turbo-1106	gemini-pro	codellama
pair_wise_correlation_matrix	1.0	1.0	1.0	0.9	1.0	0.5	0.1

Methods: test cases

Use case: Count segmented objects in a folder of segmentation results.

 Ganglioneuroblastoma_0.tif
 Ganglioneuroblastoma_1.tif
 Ganglioneuroblastoma_2.tif
 Ganglioneuroblastoma_3.tif
 Ganglioneuroblastoma_4.tif



300
398
368
378
363

Unit-test pass-rate (n=10):

	reference	gpt-4-turbo-2024-04-09	Claude-3-opus-20240229	gpt-4-1106-preview	gpt-3.5-turbo-1106	gemini-pro	codellama
workflow_batch_process_folder_count_labels	1.0	0.1	0.0	0.3	0.0	0.0	0.0

Results

Unit-test pass-rate (n=10)

	reference	gpt-4-turbo-2024-04-09	Claude-3-opus-20240229	gpt-4-1106-preview	gpt-3.5-turbo-1106	gemini-pro	codellama
Statistics / tabular data wrangling							
combine_columns_of_tables	1.0	0.8	0.1	1.0	0.9	0.7	0.1
create_umap	1.0	0.8	1.0	0.9	1.0	0.8	0.0
t_test	1.0	1.0	1.0	0.9	1.0	0.5	0.3

Measurements / feature extraction

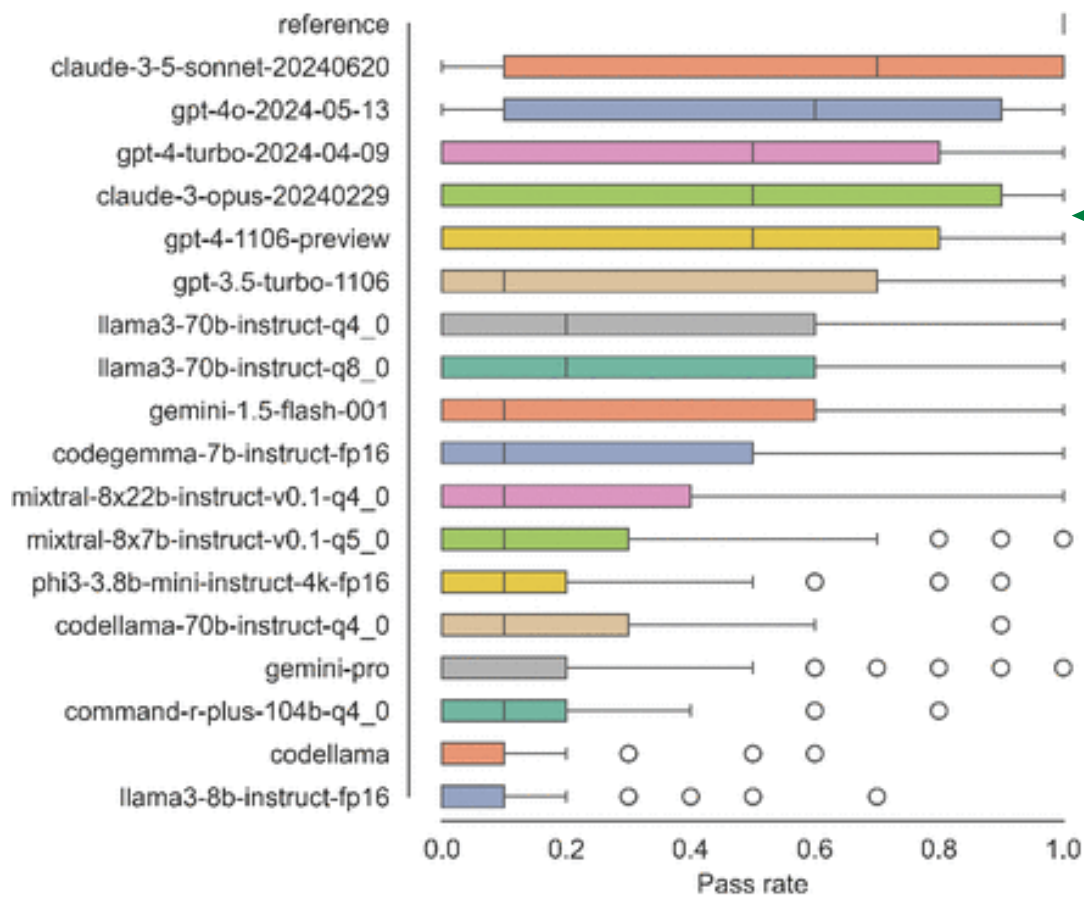
measure_intensity_over_time	1.0	0.9	0.4	0.1	0.4	0.0	0.1
measure_intensity_of_labels	1.0	0.2	0.4	0.4	0.1	0.0	0.0
measure_properties_of_regions	1.0	0.4	0.6	0.8	0.2	0.0	0.1
count_number_of_touching_neighbors	1.0	0.6	0.1	0.2	0.1	0.0	0.0

Advanced workflows / big data

tiled_image_processing	1.0	0.2	0.0	0.0	0.0	0.0	0.0
workflow_batch_process_folder_measure_intensity	1.0	0.5	0.0	0.9	0.1	0.0	0.0

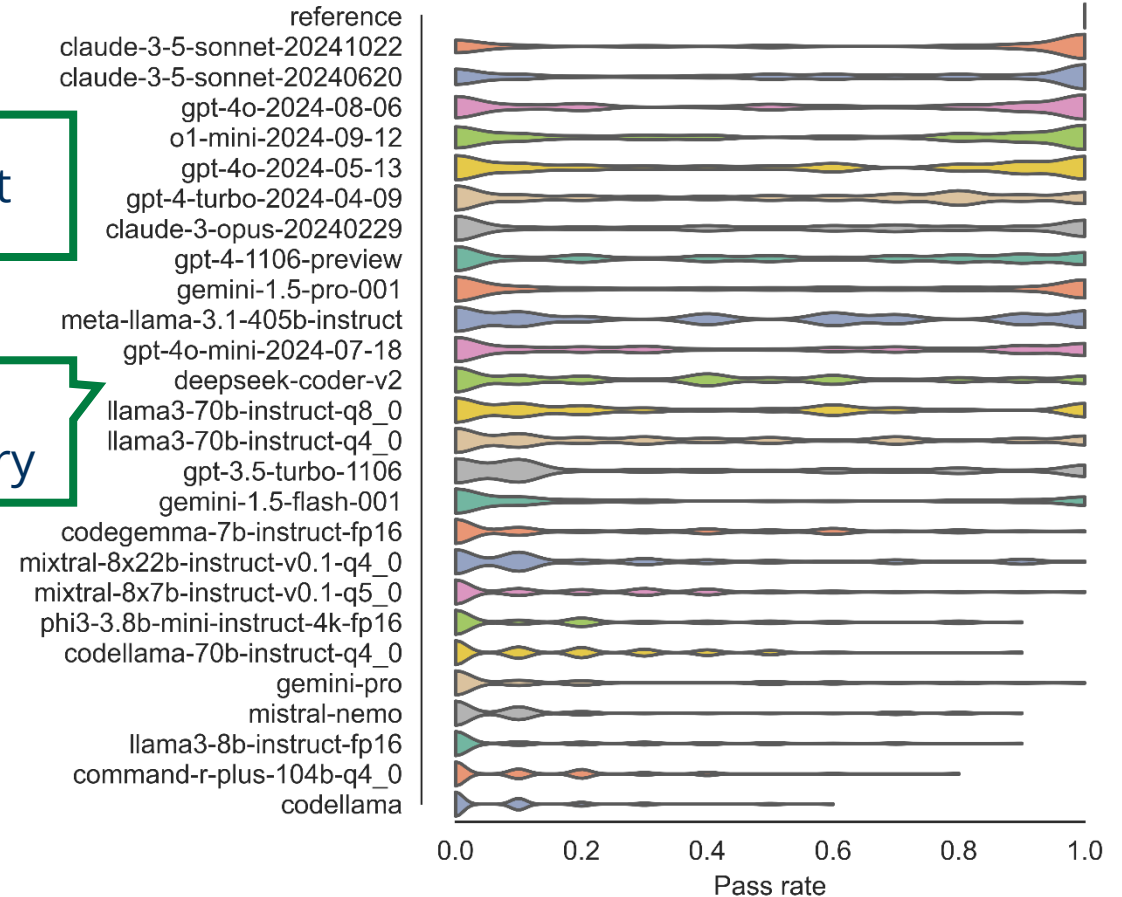
Results

Summary: 57 use-cases (yet), 26 LLMs (yet), n=10



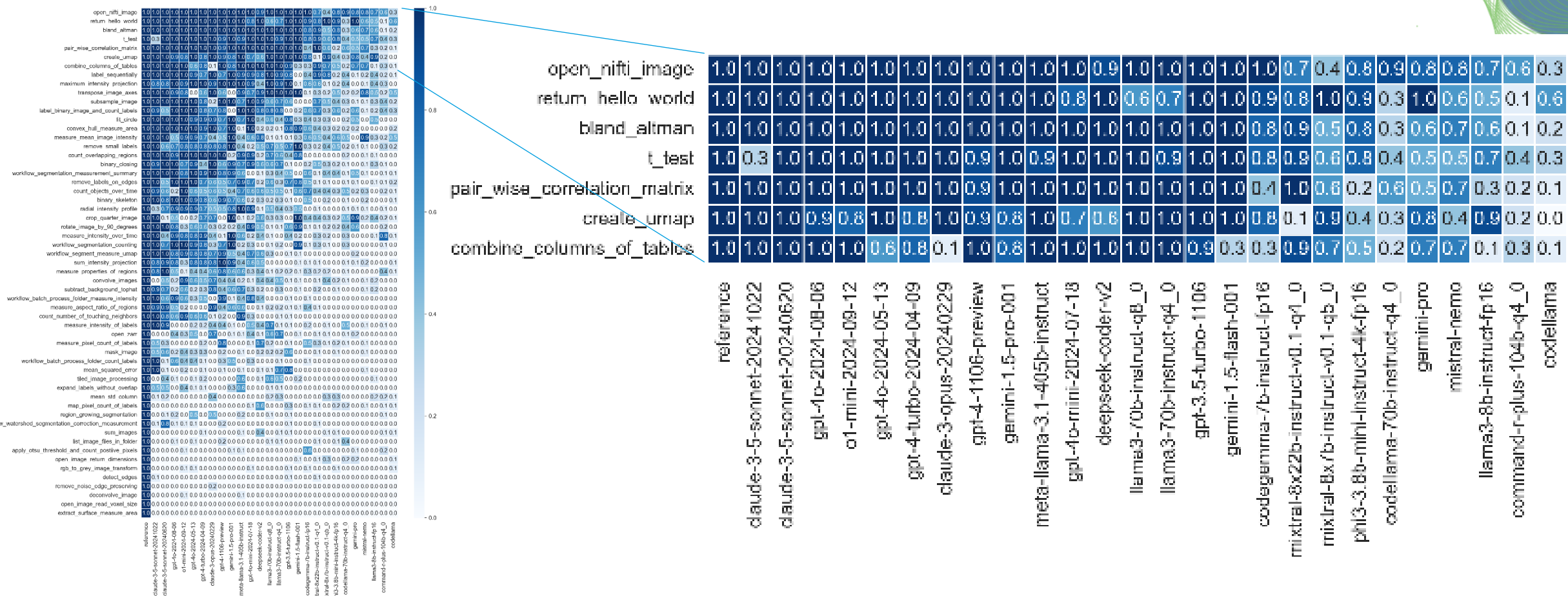
Preprint

Github repository



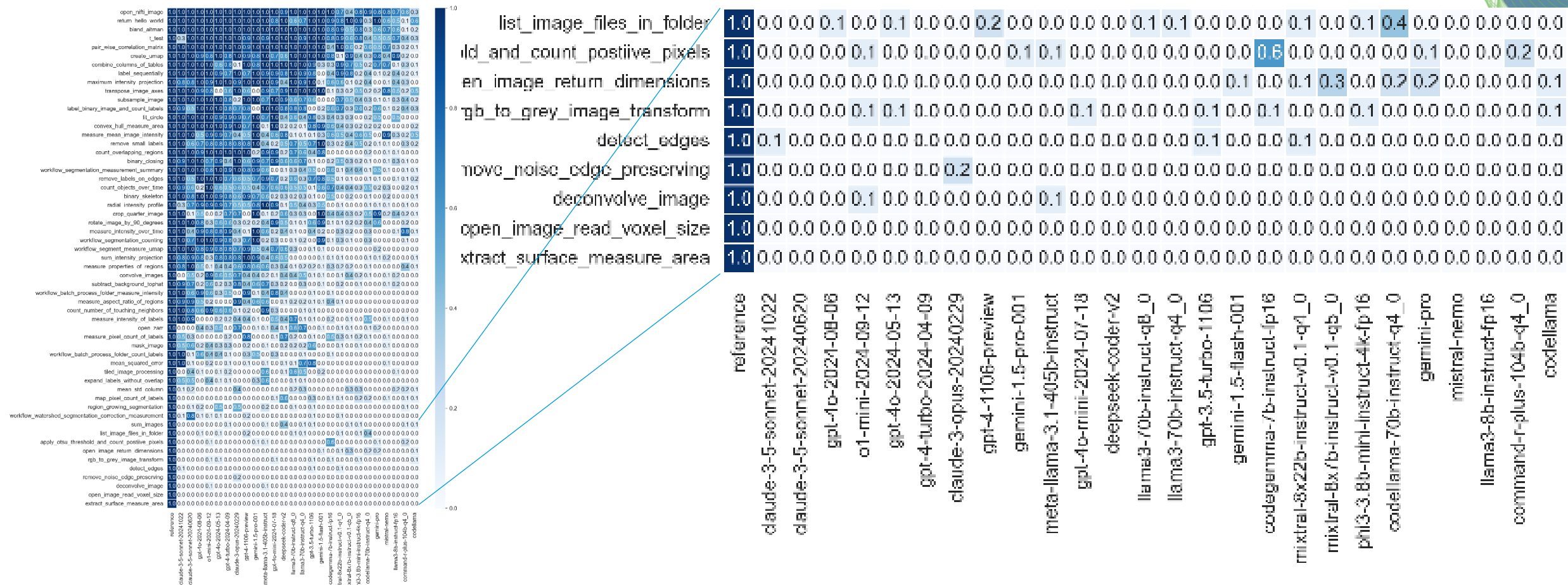
Results in more detail

Summary: 57 use-cases (yet), 26 LLMs (yet), n=10



Results in more detail

Summary: 57 use-cases (yet), 26 LLMs (yet), n=10



Results in even more detail

Gpt4 improvement over 10 months: 7%



gpt 4o

520%

[illegible]

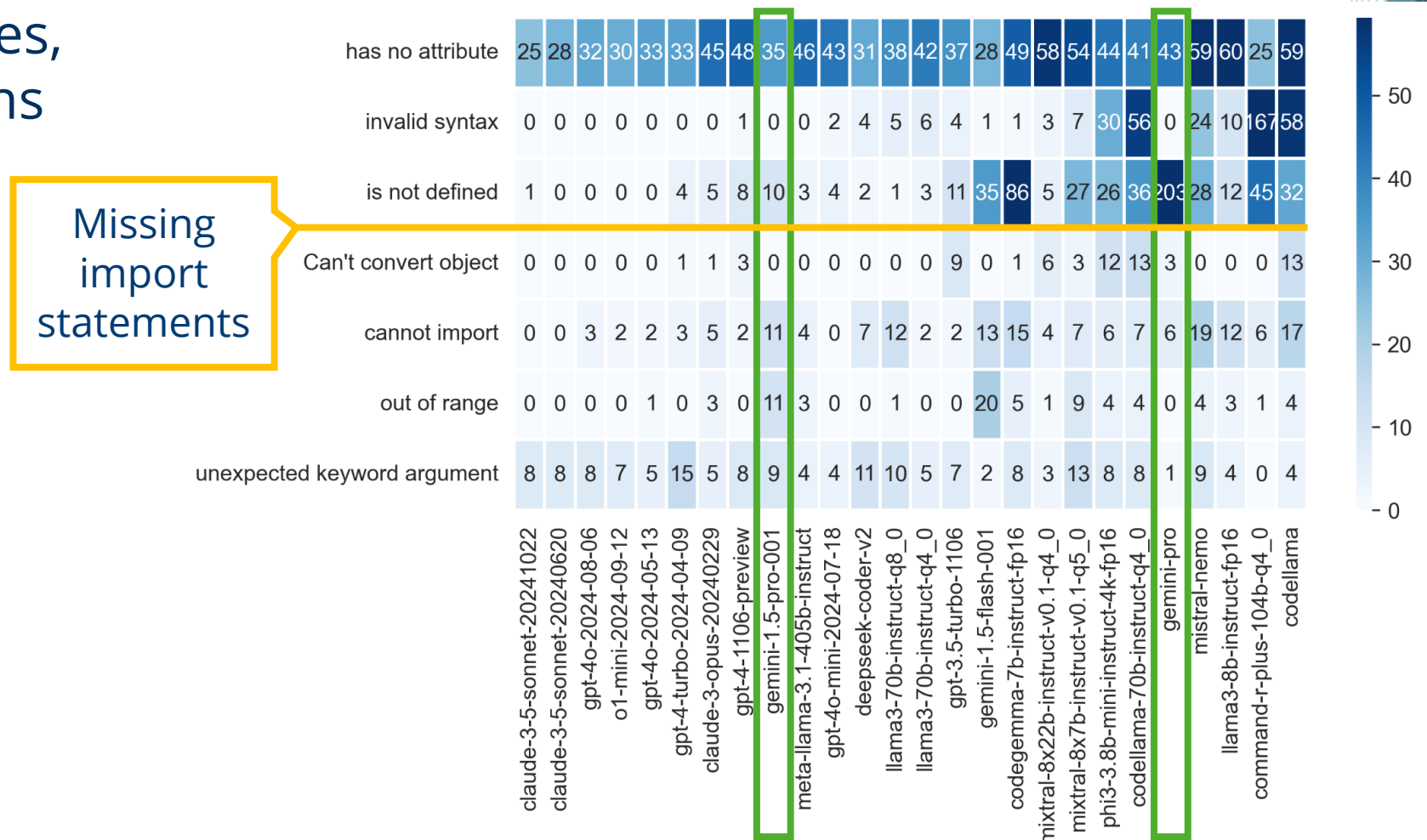
gpt 4preview
(November 2023) 46%

ScaDS.AI
DRESDEN LEIPZIG

Robert Haase
@haesleinhuepf
April 11th 2025

Results: strengths and weaknesses of LLMs

Common error messages,
e.g. for different versions
of the **gemini** model
(Google)



Results: strengths and weaknesses of LLMs

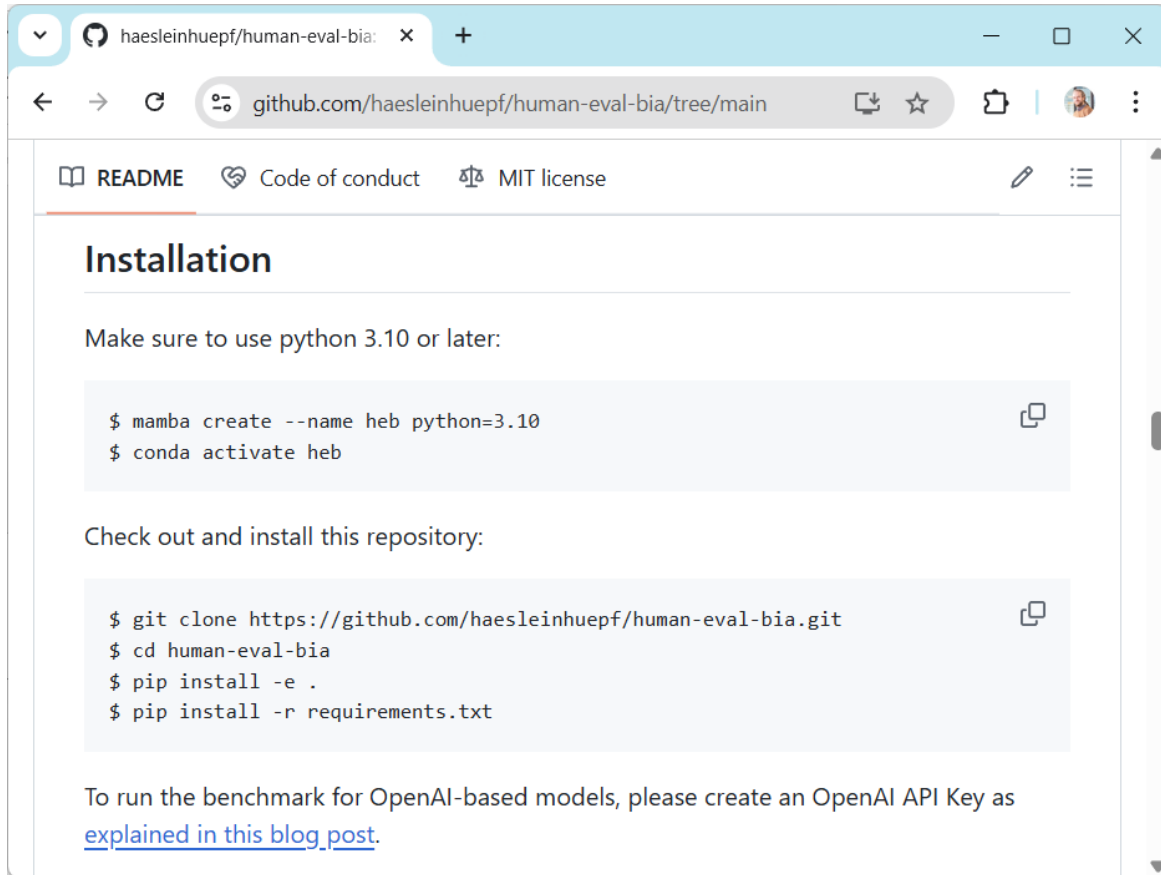
- LLMs use different Python libraries than we Bio-image Analysts do.
- What can we teach LLMs?
- What can we learn from this?

numpy	220	497	487	453	435	442	434	453	398	322	475	477	480	460	447	360	384	298	478	392	450	426	165	403	432	412	454
scipy	70	118	170	104	112	118	123	131	141	65	126	141	89	156	144	76	57	76	168	82	138	118	31	133	155	82	114
skimage	220	124	115	110	108	102	129	125	132	149	119	87	149	98	85	115	91	154	118	102	129	151	116	60	68	131	96
cv2	0	56	45	51	31	66	63	44	57	52	40	100	112	85	107	144	107	43	90	76	107	120	82	137	192	31	137
pandas	60	95	100	98	97	99	100	99	97	68	100	100	101	100	100	90	88	74	98	72	99	81	52	97	98	89	95
pyclesperanto_prototype	40	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
vedo	20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
umap	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	20	16	20	20	20	19	20
dask	10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0
nibabel	10	13	12	18	15	16	17	10	20	10	18	20	11	20	20	11	11	14	11	10	11	10	10	10	11	17	12
SimpleITK	0	7	10	2	5	7	2	10	1	4	2	0	0	0	0	8	10	2	7	10	0	9	7	8	0	0	1
trimesh	0	0	0	0	0	0	0	0	0	0	0	0	1	2	0	0	2	1	1	0	0	0	1	0	0	0	0
itk	0	7	10	2	5	6	2	10	0	3	2	0	0	0	0	9	10	4	8	10	0	9	7	6	2	0	1
	reference	claude-3-5-sonnet-2024-10-22	claude-3-5-sonnet-2024-06-20	gpt-4o-2024-08-06	o1-mini-2024-09-12	gpt-4o-2024-05-13	gpt-4-turbo-2024-04-09	claude-3-opus-2024-02-29	gpt-4-1106-preview	gemin-1.5-pro-001	meta-llama-3.1-405b-instruct	gpt-4o-mini-2024-07-18	deepseek-coder-v2	llama3-70b-instruct-q8_0	llama3-70b-instruct-q4_0	gpt-3.5-turbo-1106	gemin-1.5-flash-001	codegemma-7b-instruct-fp16	mixtral-8x22b-instruct-v0.1-q4_0	mixtral-8x7b-instruct-v0.1-q5_0	phi3-3.8b-mini-instruct-4k-fp16	codellama-70b-instruct-q4_0	gemin-pro	mistral-nemo	llama3-8b-instruct-fp16	command-r-plus-104b-q4_0	codellama

Human
reference

Additional online resources

Step-by-step guide for running the benchmark



haesleinhuepf/human-eval-bia: x +

github.com/haesleinhuepf/human-eval-bia/tree/main

README Code of conduct MIT license

Installation

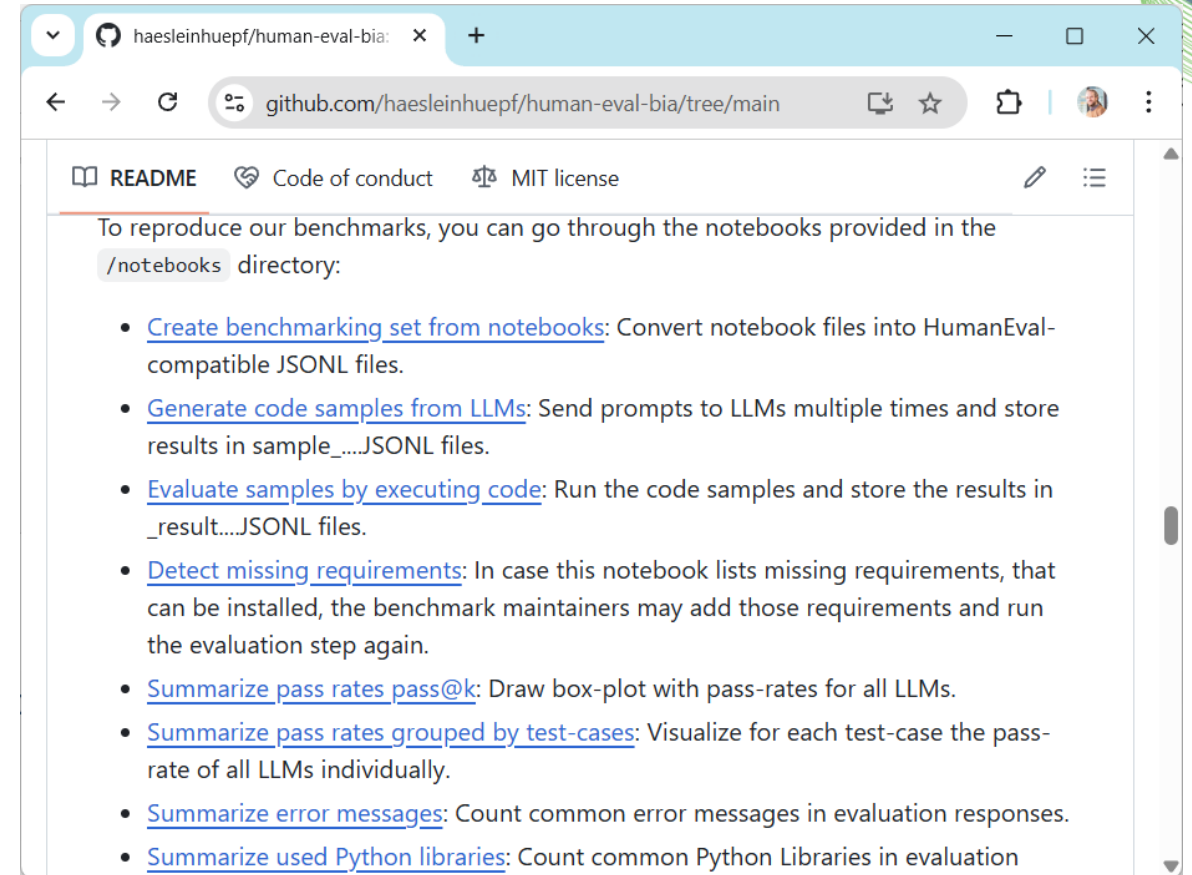
Make sure to use python 3.10 or later:

```
$ mamba create --name heb python=3.10
$ conda activate heb
```

Check out and install this repository:

```
$ git clone https://github.com/haesleinhuepf/human-eval-bia.git
$ cd human-eval-bia
$ pip install -e .
$ pip install -r requirements.txt
```

To run the benchmark for OpenAI-based models, please create an OpenAI API Key as [explained in this blog post](#).



haesleinhuepf/human-eval-bia: x +

github.com/haesleinhuepf/human-eval-bia/tree/main

README Code of conduct MIT license

To reproduce our benchmarks, you can go through the notebooks provided in the `/notebooks` directory:

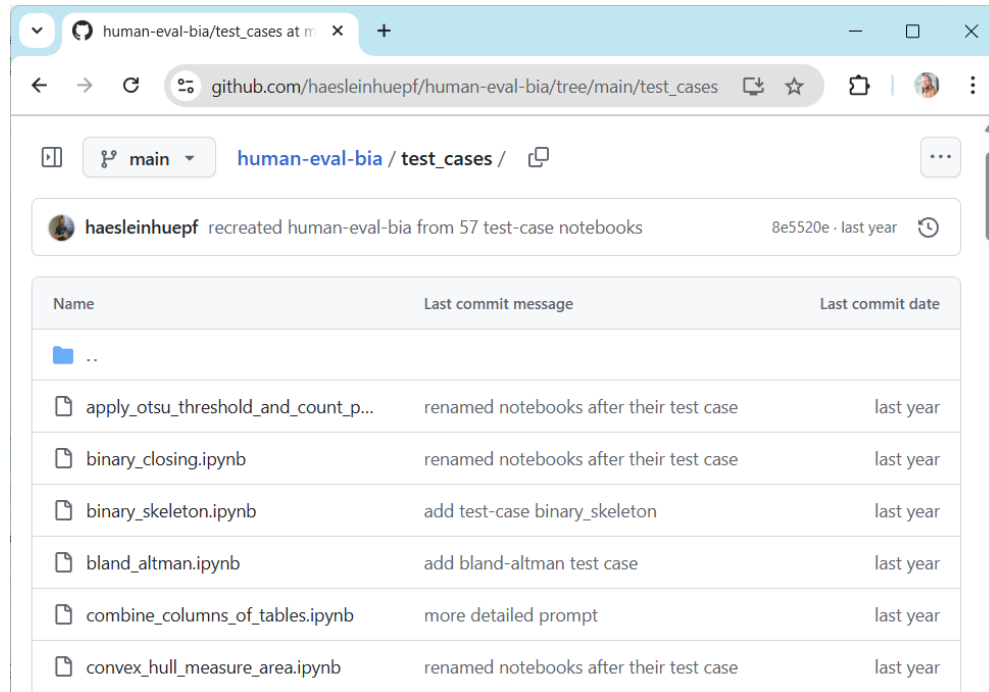
- [Create benchmarking set from notebooks](#): Convert notebook files into HumanEval-compatible JSONL files.
- [Generate code samples from LLMs](#): Send prompts to LLMs multiple times and store results in sample_...JSONL files.
- [Evaluate samples by executing code](#): Run the code samples and store the results in _result_...JSONL files.
- [Detect missing requirements](#): In case this notebook lists missing requirements, that can be installed, the benchmark maintainers may add those requirements and run the evaluation step again.
- [Summarize pass rates pass@k](#): Draw box-plot with pass-rates for all LLMs.
- [Summarize pass rates grouped by test-cases](#): Visualize for each test-case the pass-rate of all LLMs individually.
- [Summarize error messages](#): Count common error messages in evaluation responses.
- [Summarize used Python libraries](#): Count common Python Libraries in evaluation

Additional online resources

Complete set of use-cases online, open source

Advantage: Reproducibility

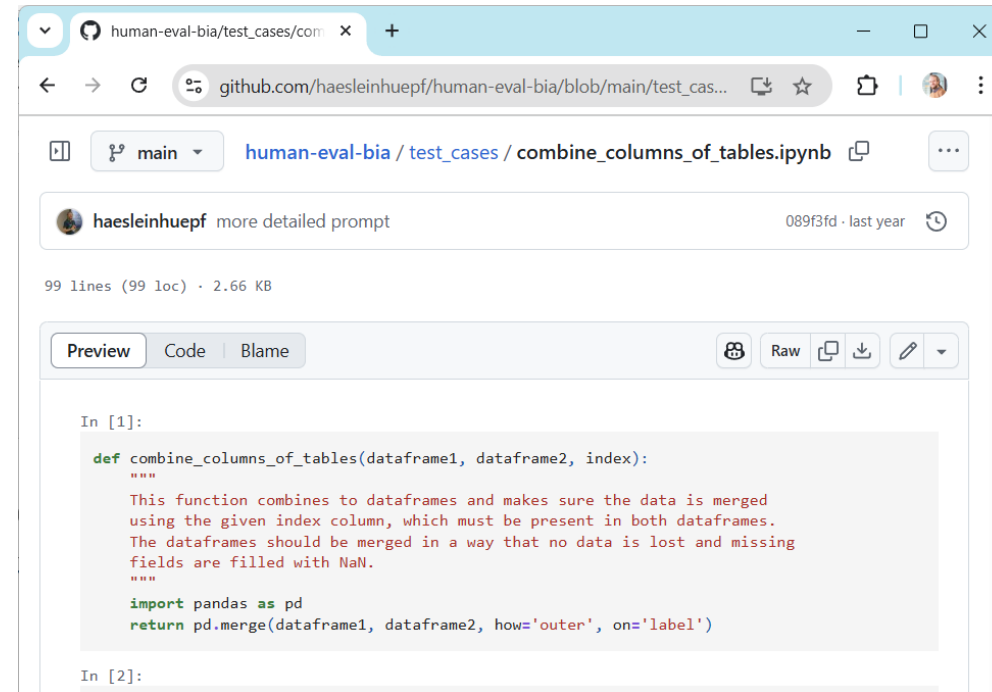
Disadvantage: Enables cheating in model training



human-eval-bia/test_cases at main · haesleinhuepf/human-eval-bia

haesleinhuepf recreated human-eval-bia from 57 test-case notebooks 8e5520e · last year

Name	Last commit message	Last commit date
..		
apply_otsu_threshold_and_count_p...	renamed notebooks after their test case	last year
binary_closing.ipynb	renamed notebooks after their test case	last year
binary_skeleton.ipynb	add test-case binary_skeleton	last year
bland_altman.ipynb	add bland-altman test case	last year
combine_columns_of_tables.ipynb	more detailed prompt	last year
convex_hull_measure_area.ipynb	renamed notebooks after their test case	last year



human-eval-bia/test_cases/combine_columns_of_tables.ipynb

haesleinhuepf more detailed prompt 089f3fd · last year

99 lines (99 loc) · 2.66 KB

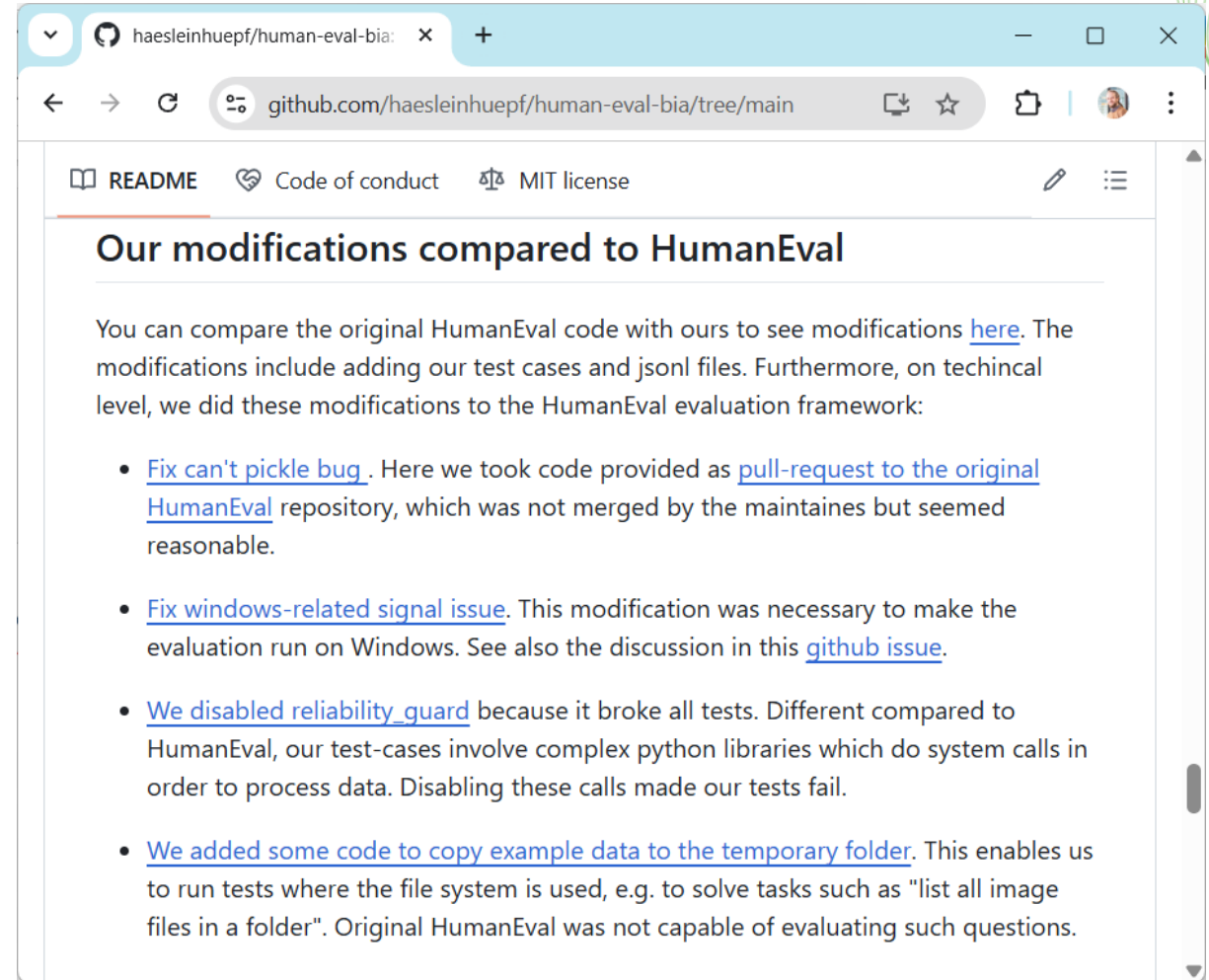
Preview Code Blame

```
In [1]:
def combine_columns_of_tables(dataframe1, dataframe2, index):
    """
    This function combines two dataframes and makes sure the data is merged
    using the given index column, which must be present in both dataframes.
    The dataframes should be merged in a way that no data is lost and missing
    fields are filled with NaN.
    """
    import pandas as pd
    return pd.merge(dataframe1, dataframe2, how='outer', on='label')

In [2]:
```

HumanEval versus HumanEvalBIA

- HumanEval does not
 - work with image data
 - work with files
 - does not work on Windows
 - allow system calls
- The HumanEval authors did not publish
 - all test-cases (just a selection)
 - the code for creating the samples.

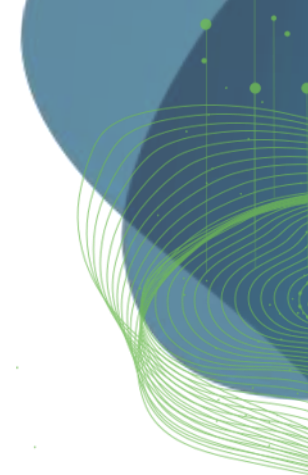


Critics

- Benchmark only measures functional correctness of code.
- Functional correctness in relation to executability would be interesting, too. -> How often does code run, but produce wrong results?
- Potential extensions:
 - Measure performance (speed, memory consumption)
 - Code quality



Critics

- 57 test-cases are not many. (20 more in the making)
 - Most test-cases written by a single person -> introduces bias
 - Open benchmarks may become part of training data of future models; rendering the benchmark useless for benchmarking
- 

Summary

- First benchmark of its kind in the bio-imaging / microscopy community
- Shows that commercial LLMs reach functional correctness of 50-60%
- Open-weight LLMs achieve < 50%
- May be useful for developing better models

Questions?