

CURIE: Evaluating LLMs On Multitask Scientific Long Context Understanding and Reasoning

Hao Cui^{*1}, Zahra Shamsi^{*1}, Gowoon Cheon^{*1}, Xuejian Ma¹, Shutong Li¹, Maria Tikhanovskaya², Peter Norgaard¹, Nayantara Mudur², Martyna Plomecka³, Paul Raccuglia¹, Yasaman Bahri¹, Victor V. Albert^{4,5}, Pranesh Srinivasan¹, Haining Pan⁶, Philippe Faist⁷, Brian Rohr⁸, Michael J. Statt⁸, Dan Morris¹, Drew Purves¹, Elise Kleeman¹, Ruth Alcantara¹, Matthew Abraham¹, Muqthar Mohammad¹, Ean Phing VanLee¹, Chenfei Jiang¹, Elizabeth Dorfman¹, Eun-Ah Kim⁹, Michael Brenner^{1,2}, Viren Jain¹, Sameera Ponda¹, Subhashini Venugopalan^{*^1}

¹Google, ²Harvard, ³University of Zurich, ⁴NIST, ⁵UMD College Park, ⁶Rutgers, ⁷FU Berlin, ⁸Modelyst, ⁹Cornell
[vsubhashini}@google.com](mailto:{vsubhashini}@google.com)



ICLR 2025

Can LLMs assist scientists in some workflows?

MIGRATION STUDIES VOLUME 3 • NUMBER 1 • 2015 • 89–110

89

Modeling internal migration flows in sub-Saharan Africa using census microdata

Andres J. Garcia^{†,‡,*}, Deepa K. Pindolia^{†,§},
Kenneth K. Lopiano^{*,††} and Andrew J. Tatem^{‡‡,§§,***}

[†]Department of Geography, University of Florida, Gainesville, FL, USA; [‡]Emerging Pathogens Institute, University of Florida, Gainesville, FL, USA; [§]Clinton Health Access Initiative, Boston MA, USA; ^{††}Department of Statistics, University of Florida, Gainesville, FL, USA; ^{‡‡}Statistical and Applied Mathematical Sciences Institute, Research Triangle Park, NC, USA; ^{§§}Department of Geography and Environment, University of Southampton, Highfield, Southampton, UK;

^{***}Fogarty International Center, National Institutes of Health, Bethesda, MD, USA; ^{***}Flowminder Foundation, Karolinska Institute, SE-171 77 Stockholm, Sweden.

*Corresponding author. Email: andygarcia@gmail.com

Abstract

Globalization and the expansion of transport networks has transformed migration into a major policy issue because of its effects on a range of phenomena, including resource flows in economics, urbanization, as well as the epidemiology of infectious diseases. Quantifying and modeling human migration can contribute towards a better understanding of the nature of migration and help develop evidence-based interventions for disease control policy, economic development, and resource allocation. In this study we paired census microdata from 10 countries in sub-Saharan Africa with additional spatial datasets to develop models for the internal migration flows in each country, including key drivers that reflect the changing social, demographic, economic, and environmental landscapes. We assessed how well these gravity-type spatial interaction models can both explain and predict migration. Results show that the models can explain up to 87 percent of internal migration, can predict future within-country migration with correlations of up to 0.91, and can also predict migration in other countries with correlations of up to 0.72. Findings show that such models are useful tools for understanding migration as well as predicting flows in regions where data are sparse, and can contribute towards strategic economic development, planning, and disease control targeting.

1. Introduction

Human population movements are an important component in a wide range of diverse

Can we reproduce the analysis in this study?

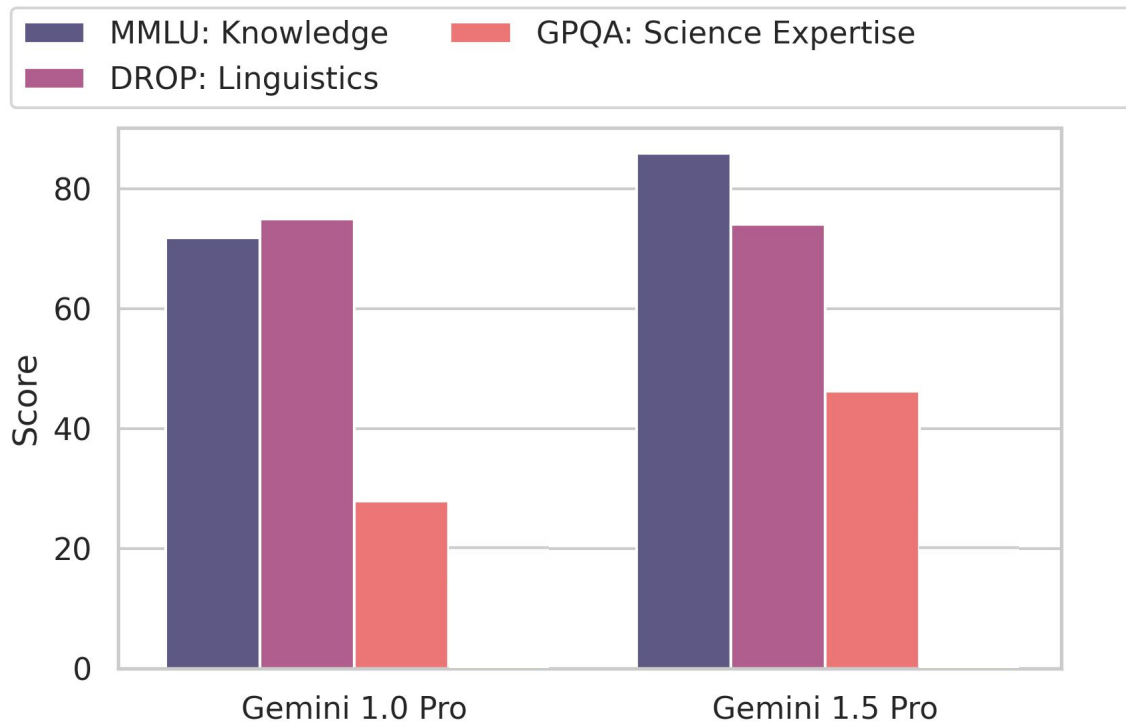
Can I apply the same techniques for Europe?

Can we measure scientific problem-solving ability?

This requires

- Knowledge of the domain
- Long-context capabilities
 - to understand context of the problem
- Reasoning ability
 - to apply the knowledge in the context of a given problem

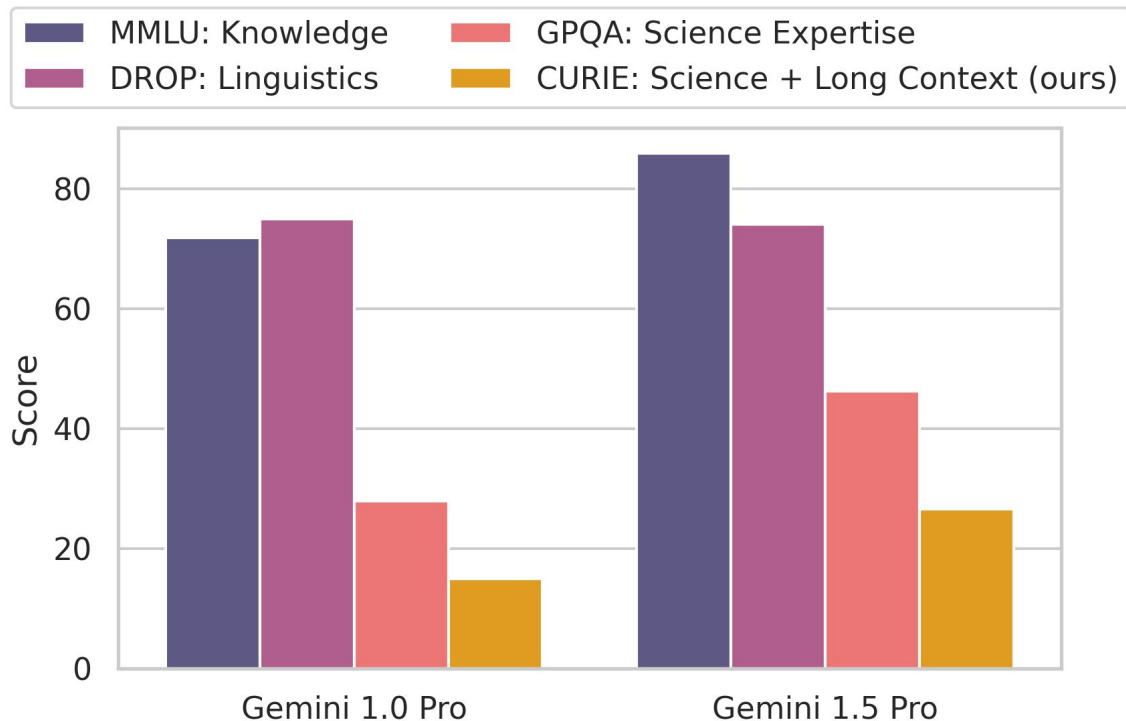
Existing benchmarks test for knowledge, linguistics



Language models demonstrate knowledge across a wide range of domains, as seen with perf. on MMLU, GPQA, DROP and other benchmarks.

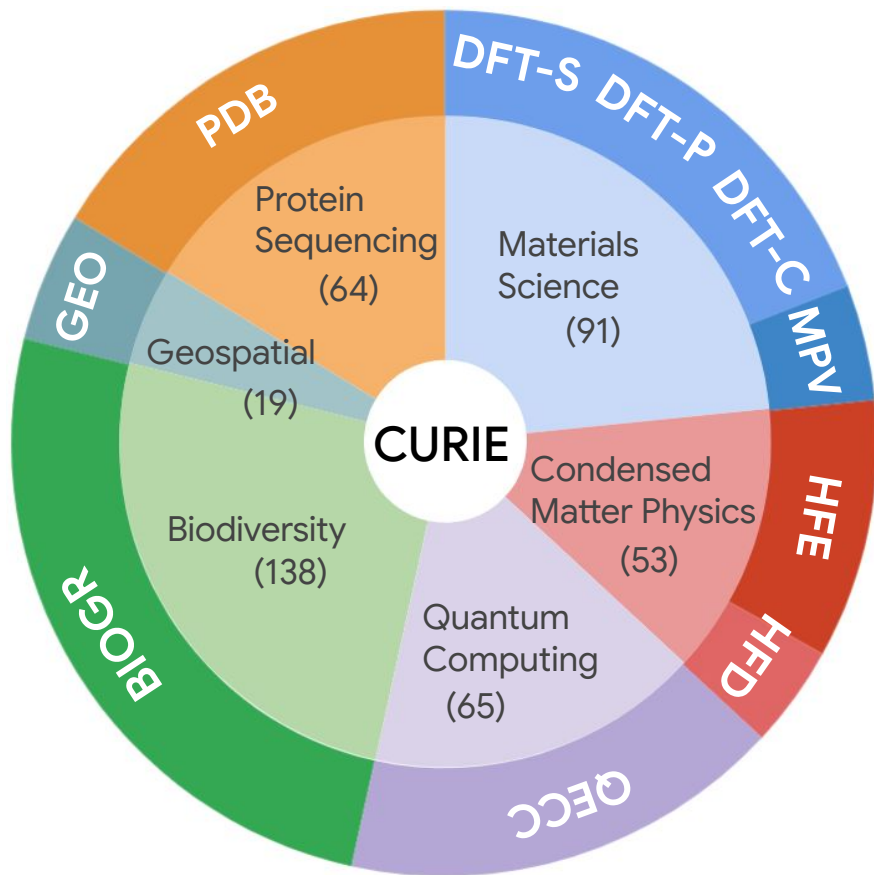
CURIE: Test scientific problem solving

(scientific long-Context Understanding Reasoning and Information Extraction benchmark)

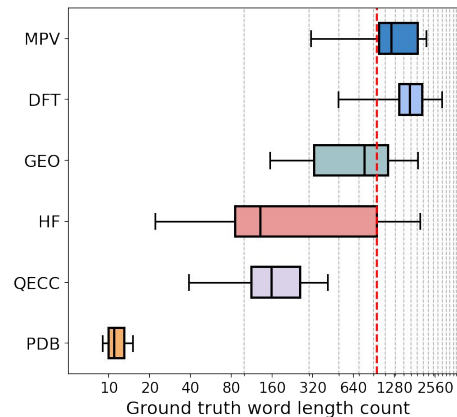
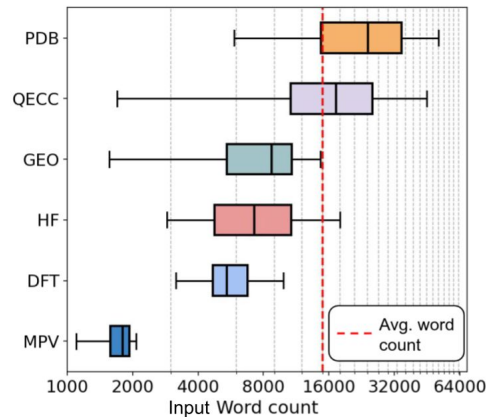
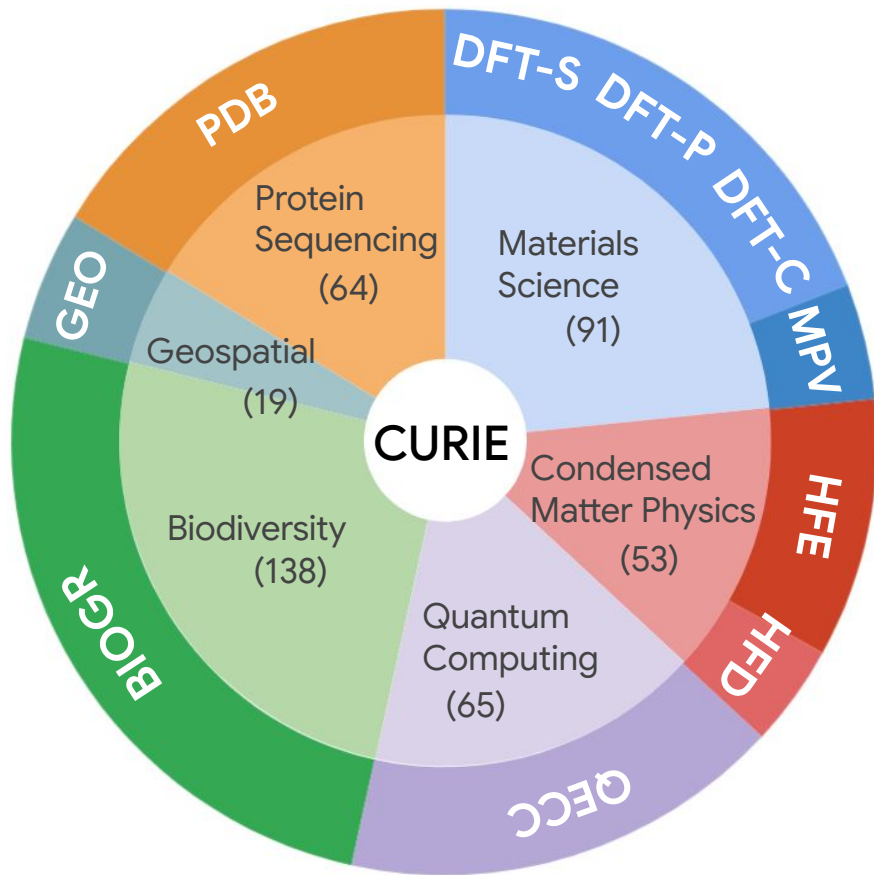


Benchmarks help drive progress.

Curated 10 tasks from experts in 6 domains



Avg. 15k words in the input, and 960 words in output



Example: Material Science

Given a paper we want to reproduce the DFT calculations done in this paper.

Task	Domain	# Qs	Brief Description
DFT-S	Material Science	74	Extracts input material structures for DFT calculations.
DFT-P	Material Science	74	Extract parameters for DFT calculations.
DFT-C	Material Science	74	Write functional code for DFT computations.
MPV	Material Science	17	Identify all instances of materials, their properties, and descriptors.

Coexistence of Co doping and strain on arsenene and antimonene: tunable magnetism and half-metallic behavior

Yungang Zhou, Geng Cheng and Jing Li

Effectively modulating the magnetism of two-dimensional (2D) systems is critical for the application of magnetic nanostructures in quantum information devices. In this work, by employing density functional theory calculations, we found the coexistence of Co doping and strain can effectively control the spin

Density Functional Theory (DFT)

DFT-S: Identify input structures.

DFT-P: Identify DFT calculations and params.

DFT-C: Write python code for DFT calculations.

```
"common_name": "arsenene",
"scientific_name": "NaN",
"type": "surface",
"composition": "As2",
"crystal_or_isolated": "surface",
"vacuum": "[0,0,15]",
"supercell": "[4,4,1]",
"cell_size": "NaN",
```

```
"software": "vasp",
"functional": "PBE",
"k-points": "[8,8,1]",
"pseudopotentials": "NaN",
"basis_set": "NaN",
"energy_cutoff": 500.0,
"force_convergence": 0.01,
"energy_convergence": "NaN",
```

```
def get_strained_structures(atoms: Atoms) -> list[Atom]
    strains = np.linspace(0.96, 1.08, 7)
    return_list = []
    for strain in strains:
        strained_atoms = deepcopy(atoms)
        atoms.cell *= strain
        atoms.positions *= strain
        return_list.append(strained_atoms)
```


CURIE: 10 tasks requiring different capabilities

Task	Domain	# Qs	Brief Description	Capability	Output Format	Primary Eval. metric
DFT-S	Material Science	74	Extracts input material structures for DFT calculations.	entity recognition, concept tracking	JSON	LLMSim-F1
DFT-P	Material Science	74	Extract parameters for DFT calculations.	concept extraction, tracking, aggregation	JSON	LLMSim-F1
DFT-C	Material Science	74	Write functional code for DFT computations.	concept aggregation, coding	TEXT	ROUGE-L
MPV	Material Science	17	Identify all instances of materials, their properties, and descriptors.	entity recognition, concept extraction, tracking	JSON	LLMSim-F1
QECC	Quantum Computing	65	Create a YAML file with the Error Correction Code's properties.	concept aggregation, summarization	YAML	ROUGE-L

Different kinds of outputs: dicts, equations, text etc.

Task	Domain	# Qs	Brief Description	Capability	Output Format	Primary Eval. metric
HFD	Condensed Matter Physics	64	Derive the Hartree-Fock mean-field Hamiltonian for a quantum many-body system.	concept extraction, algebraic manipulation, reasoning	TEXT	ROUGE-L
HFE	Condensed Matter Physics	19	Extract the most general mean-field Hamiltonian.	concept extraction	TEXT (latex equation)	ROUGE-L
GEO	Geospecial	15	Extract information for all geospatial datasets used along with the spatial and temporal extents.	concept extraction, aggregation	JSON	ROUGE-L
BIOGR	Biodiversity	38	Determine the latitude, longitude bounding box encompassing the region in the map image.	visual comprehension, reasoning	JSON (lat. lon. co-ordinates)	Intersection-over-Union (IoU)
PDB	Protein Sequencing	138	Reconstruct a protein's amino acid sequence from the 3D structure.	tracking, aggregation reasoning	Code or TEXT (seq.)	Identity ratio (IDr)

Evaluation metrics

Programmatic

Doesn't require an LLM e.g. ROUGE-L, IoU

LLM-based

Uses an LLM as a proxy to rate or measure semantic closeness

LMScore: Coarse evaluation of outputs

$$LMScore = \sum_{t=0}^2 p(x_t) \times w_t$$

$$x_t \in \{\text{bad}, \text{ok}, \text{good}\}$$

$$w_t \in \{0, 0.5, 1\}$$

LLMSim: LLM eval for optimal match b/w list of dicts

D_G A set of ground truth dictionaries

```
[  
  {"material": "Indium Nitride", "property": "band gap"},  
  {"material": "Silicon", "property": "power conversion efficiency"},  
  {"material": "Zinc Oxide", "property": "Direct band gap"},  
]
```

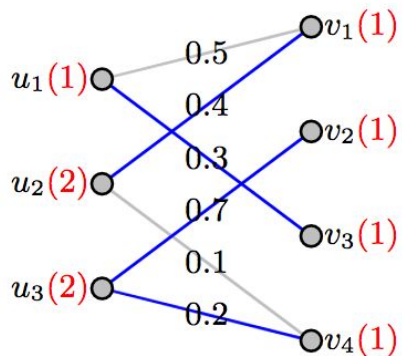
D_P A set of predicted dictionaries

```
[  
  {"material": "ZnO", "property": "Exciton binding energy"},  
  {"material": "Indium nitride", "property": "band gap"},  
  {"material": "Si", "property": "power conversion efficiency\nPCE"},  
  {"material": "ZnO", "property": "band gap"},  
]
```

LLMSim: LLM eval for optimal match b/w list of dicts

$\text{LLMSim} = M(D_P, D_g)$ *Match each ground truth record to a predicted record*

$= \begin{cases} \text{None, if no match in values} \\ D_p \in D_P : \arg \max s(f_i, D_p, D_g) \end{cases}$ *Select predicted record most similar to ground truth*



LLMSim: LLM eval for optimal match b/w list of dicts

$$\begin{aligned} \text{LLMSim} &= M(D_P, D_g) \quad \text{Match each ground truth record to a predicted record} \\ &= \begin{cases} \text{None, if no match in values} \\ D_p \in D_P : \arg \max s(f_i, D_p, D_g) \end{cases} \quad \text{Select predicted record most similar to ground truth} \end{aligned}$$

$$Pr = \frac{|(D_p, D_g) \in M|}{|D_P|}, Re = \frac{|(D_p, D_g) \in M|}{|D_G|} \quad \text{Compute Precision and Recall}$$

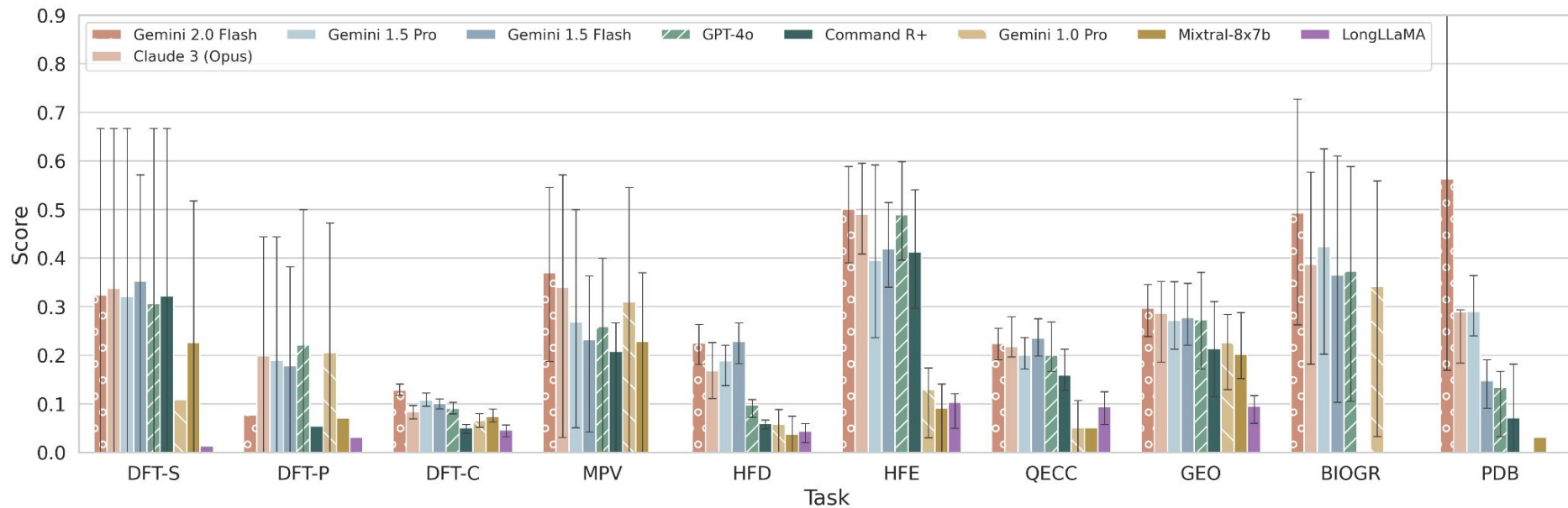
LLMSim: LLM eval for optimal match b/w list of dicts

$$\begin{aligned} \text{LLMSim} &= M(D_P, D_g) \quad \text{Match each ground truth record to a predicted record} \\ &= \begin{cases} \text{None, if no match in values} \\ D_p \in D_P : \arg \max s(f_i, D_p, D_g) \end{cases} \quad \text{Select predicted record most similar to ground truth} \end{aligned}$$

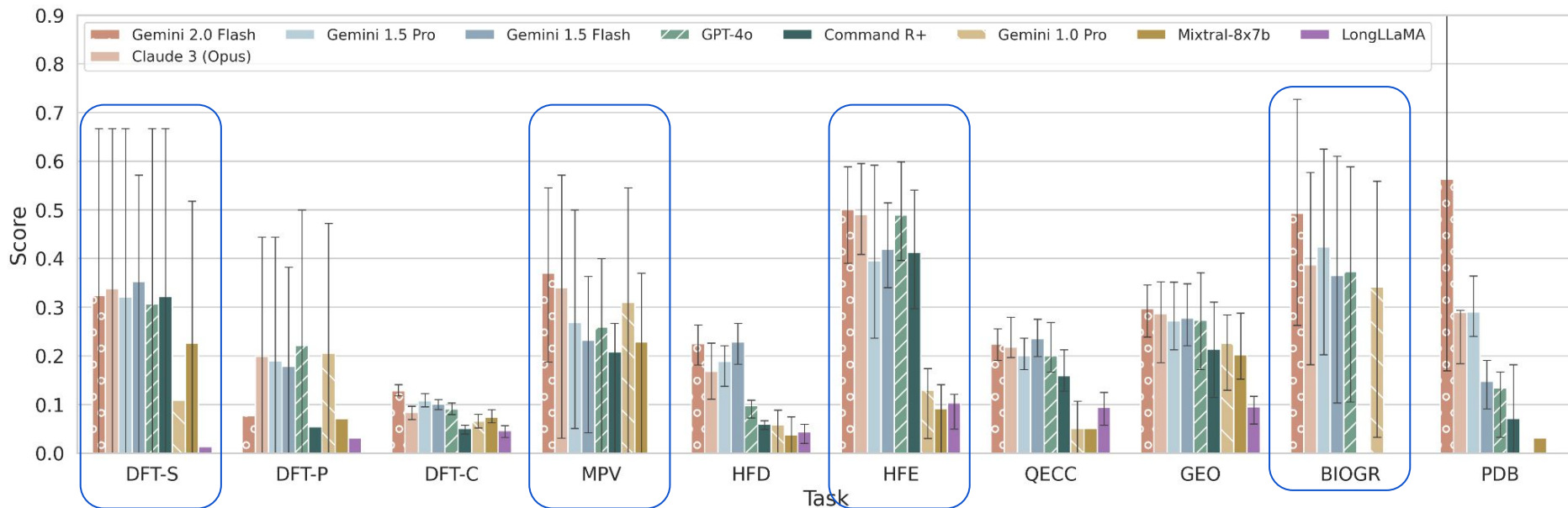
$$Pr = \frac{|(D_p, D_g) \in M|}{|D_P|}, Re = \frac{|(D_p, D_g) \in M|}{|D_G|} \quad \text{Compute Precision and Recall}$$

$$f_1 = \frac{2 \times Pr \times Re}{Pr + Re}, F1_{macro} = \frac{\sum_1^N f_1}{N} \quad \text{Compute f1 score and avg. F1}$$

Analysis across tasks

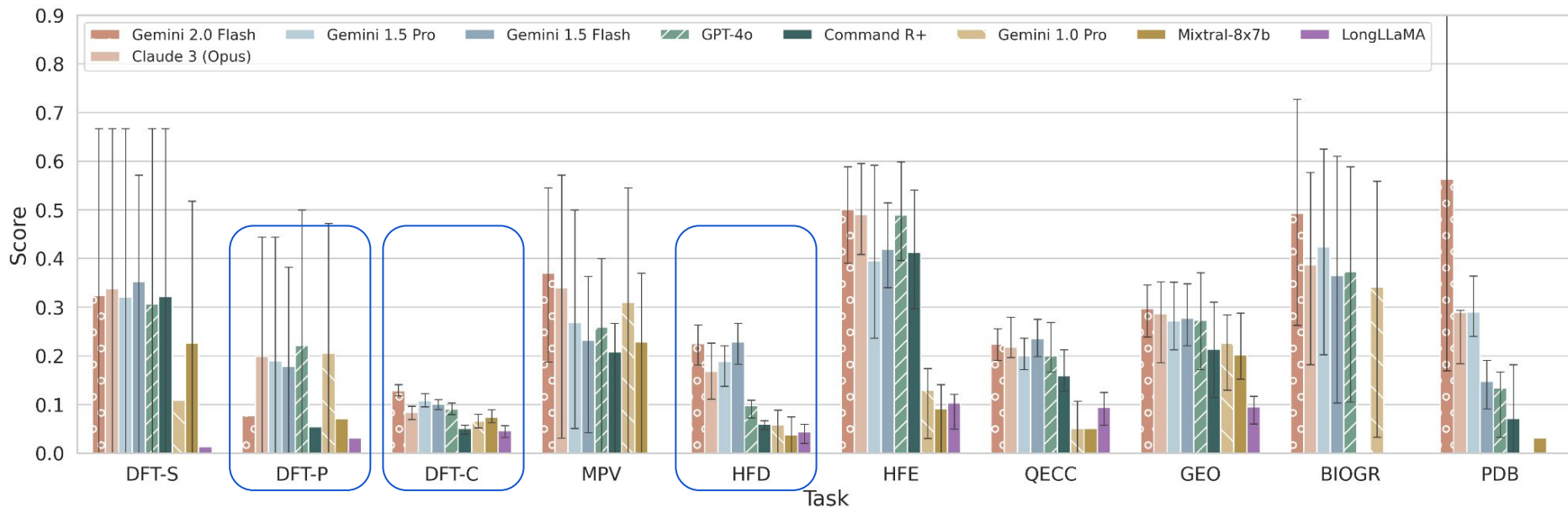


Frontier models do better on extraction tasks



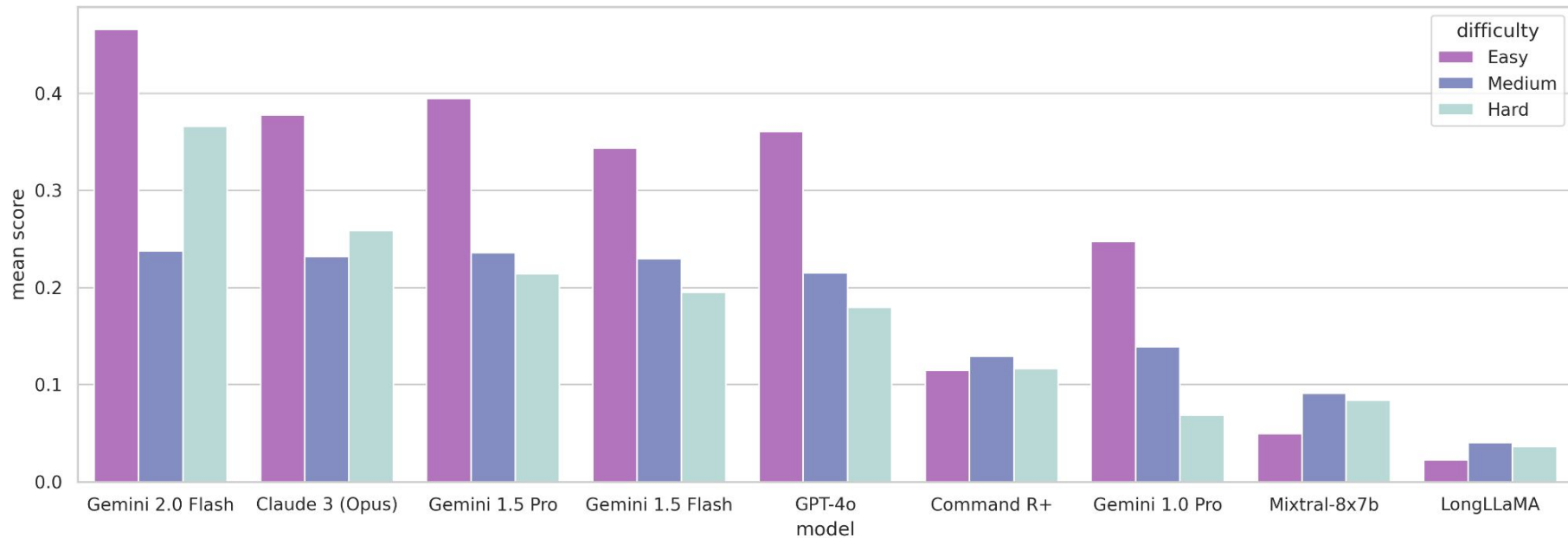
- Extraction tasks (DFT-S, MPV, HFE) and geo-referencing (BIOGR) are easier.

Reasoning - derivation, aggregation see lower perf.



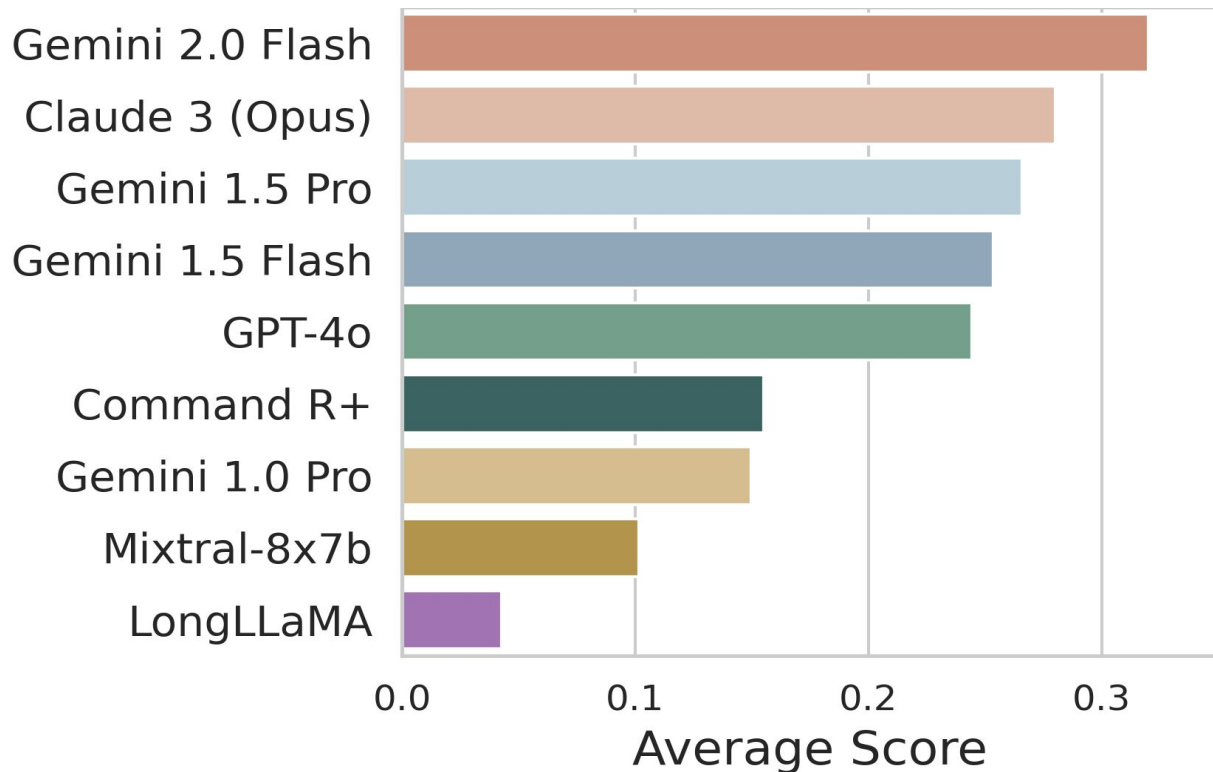
- Extraction tasks (DFT-S, MPV, HFE) and geo-referencing (BIOGR) better perf.
- Reasoning e.g. derivation (HFD), aggregation and coding (DFT-P, DFT-C) harder.

Sliced by difficulty, models do better on easy examples

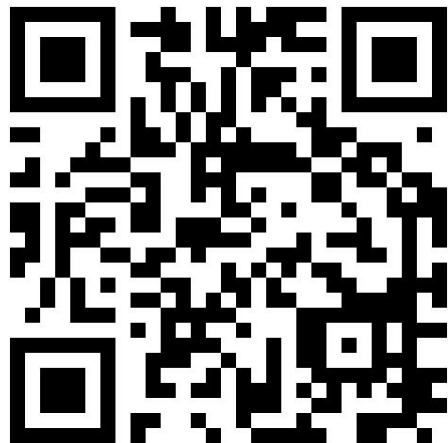


Experts marked each example as easy, difficult or hard, often based on how spread-out the information required to answer the question is.

Highest score 32% – much room for improvement



Data and code on GitHub

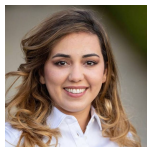


github.com/google/curie

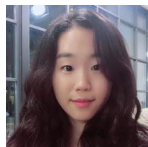


arxiv.org/abs/2503.13517

Thanks!



Zahra
Shamsi



Gowoon
Cheon



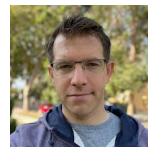
Jackson
Cui



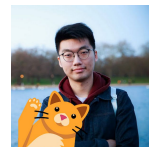
Subhashini
Venugopalan



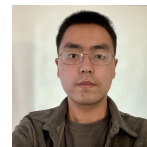
Sameera
Ponda



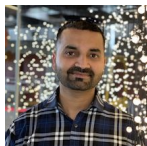
Peter
Norgaard



Shutong
Li



Xuejian
Ma



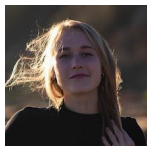
Matthew
Abraham



Nayantara
Mudur



Michael
Brenner



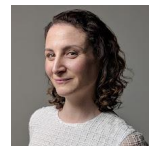
Maria
Tikhonovskaya



Martyna
Plomecka



Paul
Raccuglia



Lizzie
Dorfman



Yasaman
Bahri



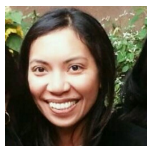
Dan
Morris



Drew
Purves



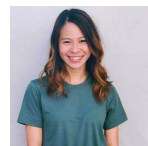
Elise
Kleeman



Ruth
Alcantara



Eun-Ah
Kim



Phing
Lee



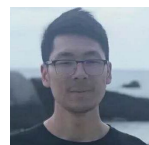
Chenfei
Jiang



Viren
Jain



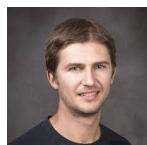
Muqthar
Mohammad



Haining
Pan



Philippe
Faist



Victor
Albert



Brian
Rohr



Michael
Statt