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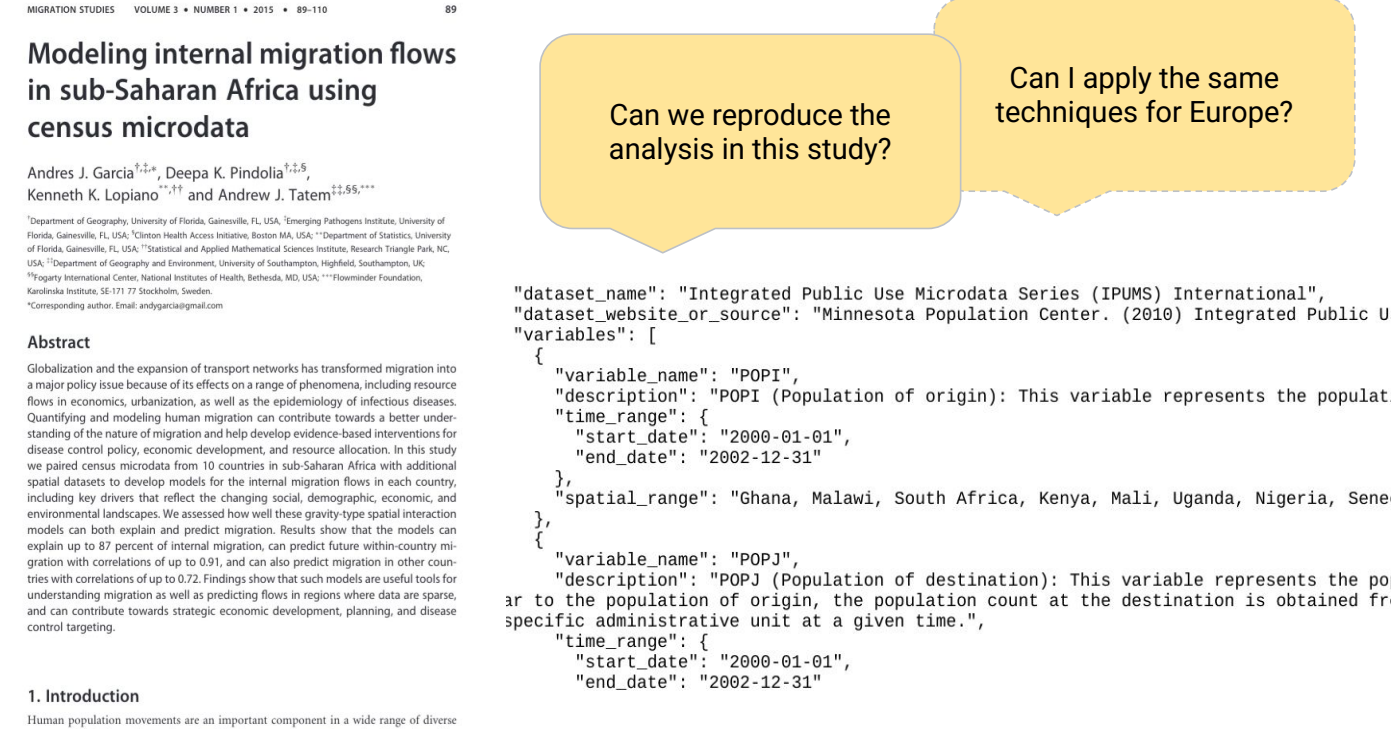
Google Research

Hao Cui*[1], Zahra Shamsi*[1], Gowoon Cheon*[1], Xuejian Ma[1], Shutong Li[1], Maria Tikhanovskaya[2], Peter Norgaard[1], Nayantara Mudur[2], Martyna Plomecka[3], Paul Raccuglia[1], Yasaman Bahri[1], Victor V. Albert[4,5], Pranesh Srinivasan[1], Haining Pan[6], Philippe Faist[7], Brian Rohr[8], Michael J. Statt[8], Dan Morris[1], Drew Purves[1], Elise Kleeman[1], Ruth Alcantara[1], Matthew Abraham[1], Muqthar Mohammad[1], Ean Phing VanLee[1], Chenfei Jiang[1], Elizabeth Dorfman[1], Eun-Ah Kim[9], Michael Brenner[1,2], Viren Jain[1], Sameera Ponda[1], Subhashini Venugopalan*^[1]

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

## Can LLMs assist Scientists with some workflows?

### Can we measure problem solving ability?
- Extract details of the data.
- Identify and extract the processes and methodology.
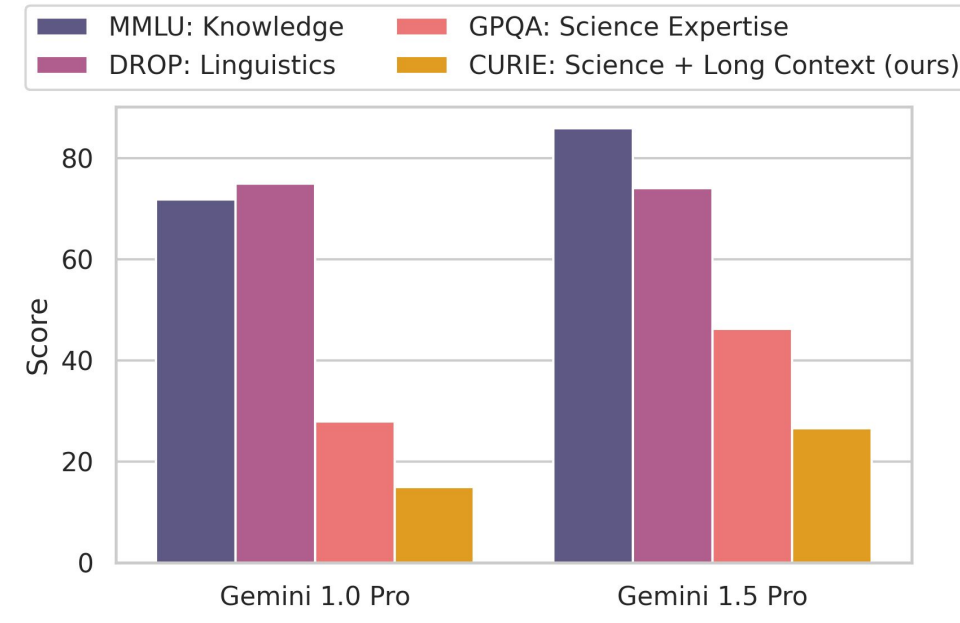- Write code to solve problem or reproduce study.

### This requires
- Knowledge of the domain.
- Processing long-context info.
- Reasoning ability to apply knowledge in the context of a problem.

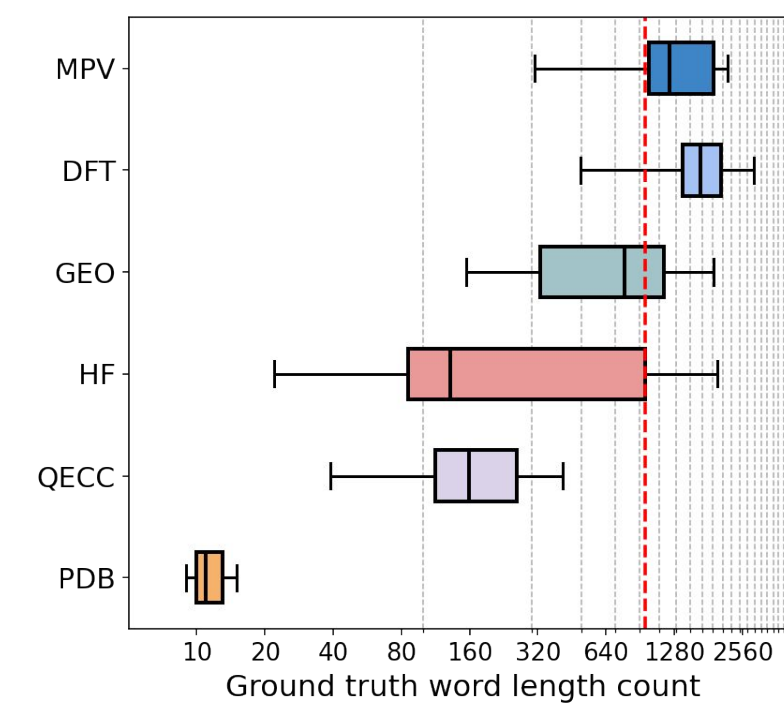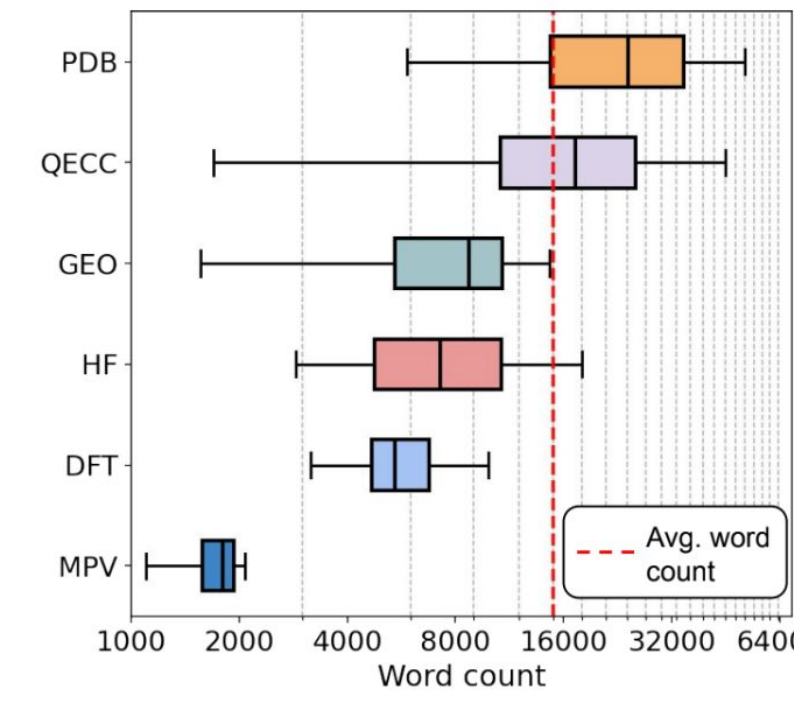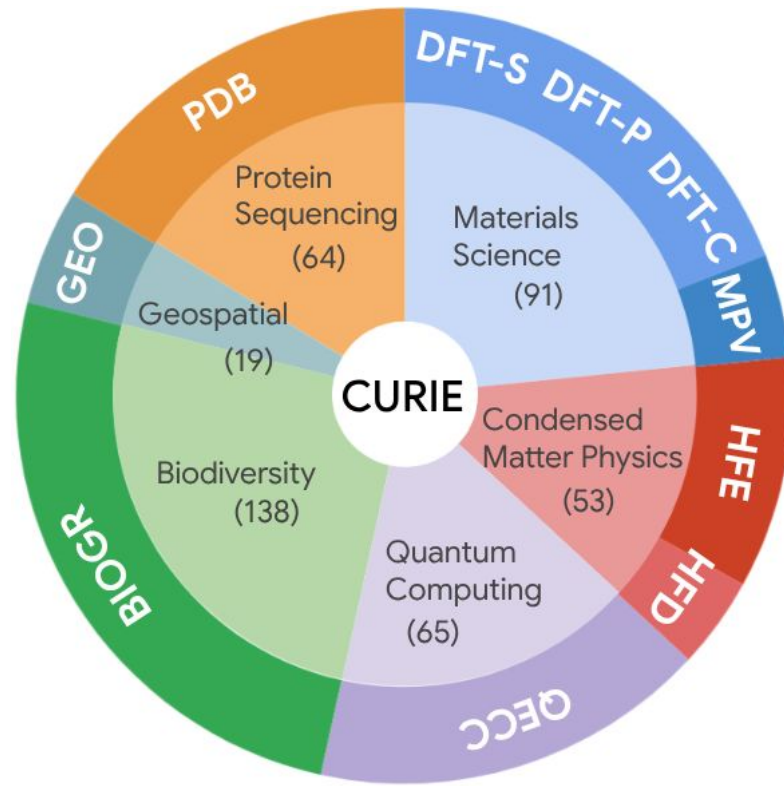### CURIE: Tests scientific problem solving
- Existing benchmarks test for knowledge, linguistics.
- CURIE: scientific long-**C**ontext **U**nderstanding **R**easoning and **I**nformation **E**xtraction benchmark

## The CURIE benchmark and dataset
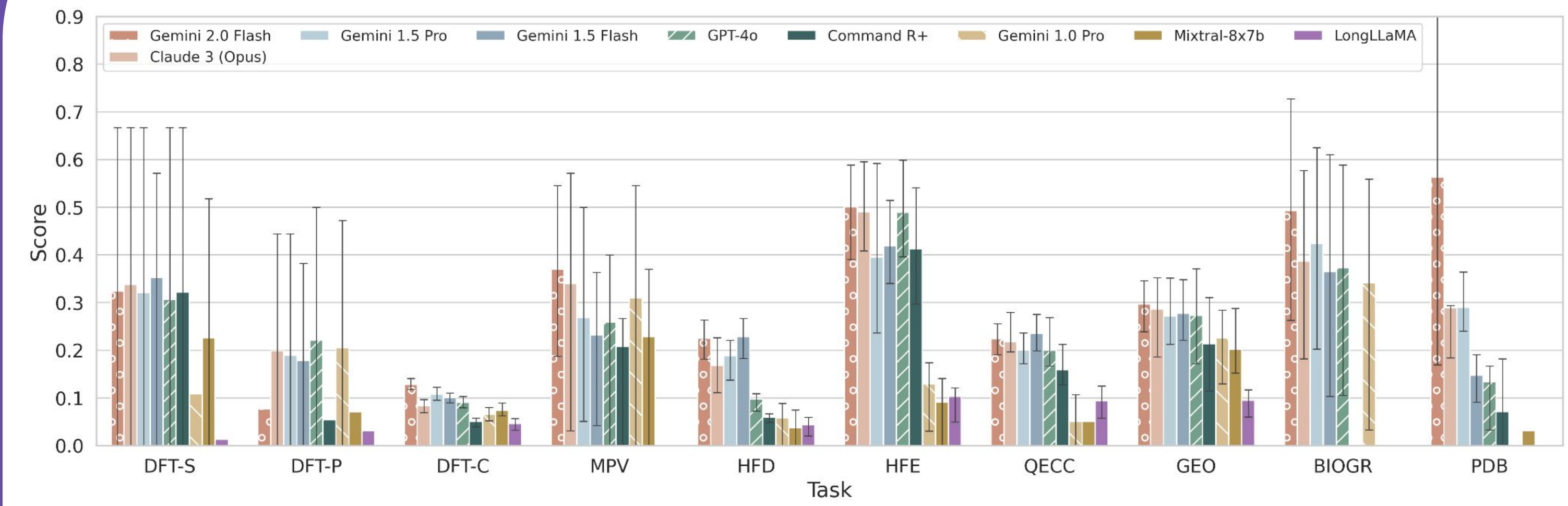
### 580 examples, 429 documents, 10 tasks, 6 domains
- Tasks require expertise + reasoning + long-context understanding
- Avg. input query length ~15k words, gold response length is ~1k

We collaborated closely with domain experts throughout benchmark development. This included:
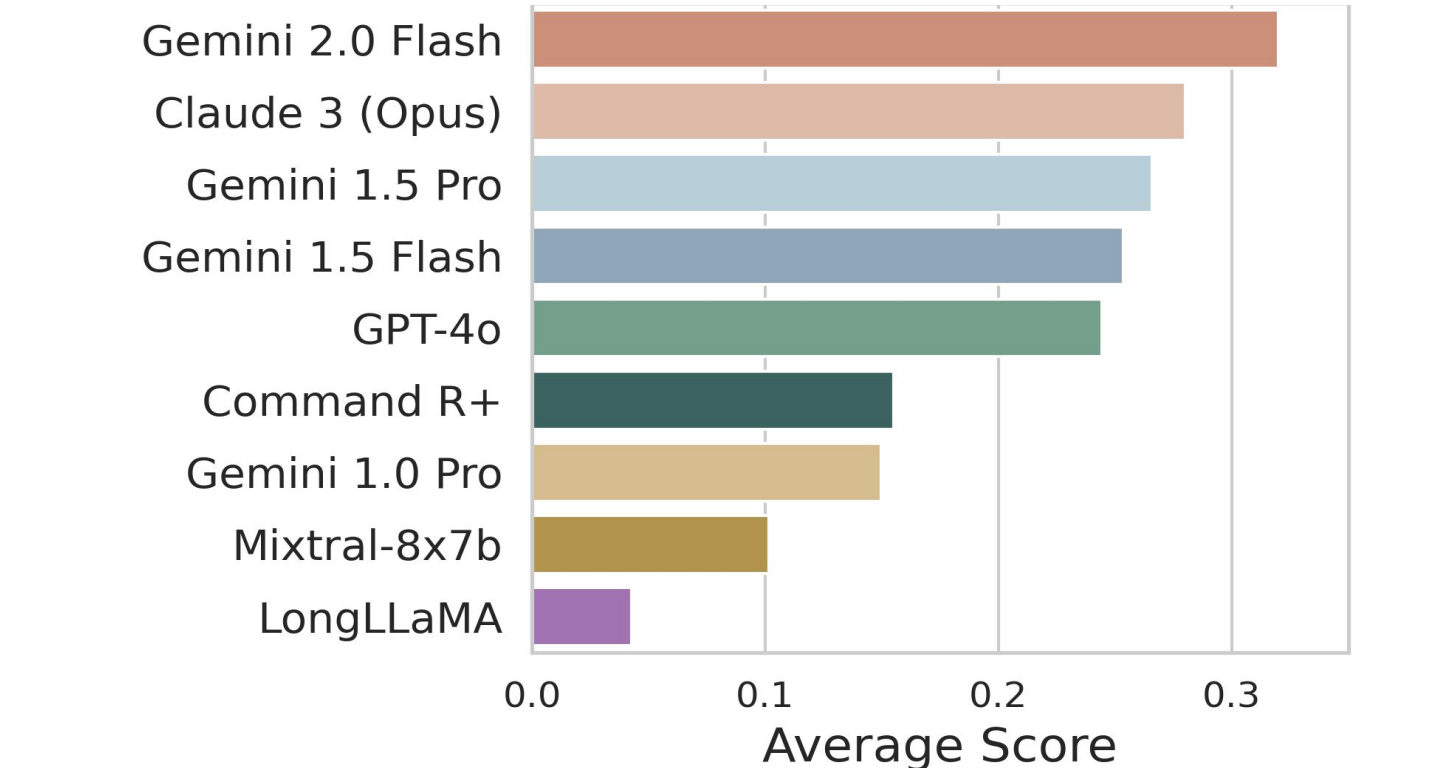
- Defining and identifying tasks reflecting realistic scientific workflows.
- Sourcing relevant papers from the domain.
- Creating accurate, nuanced, and comprehensive ground truth answers.
- Rating task difficulty based on salient features.
- Identifying and verifying evaluation metrics against expert judgments of model responses.
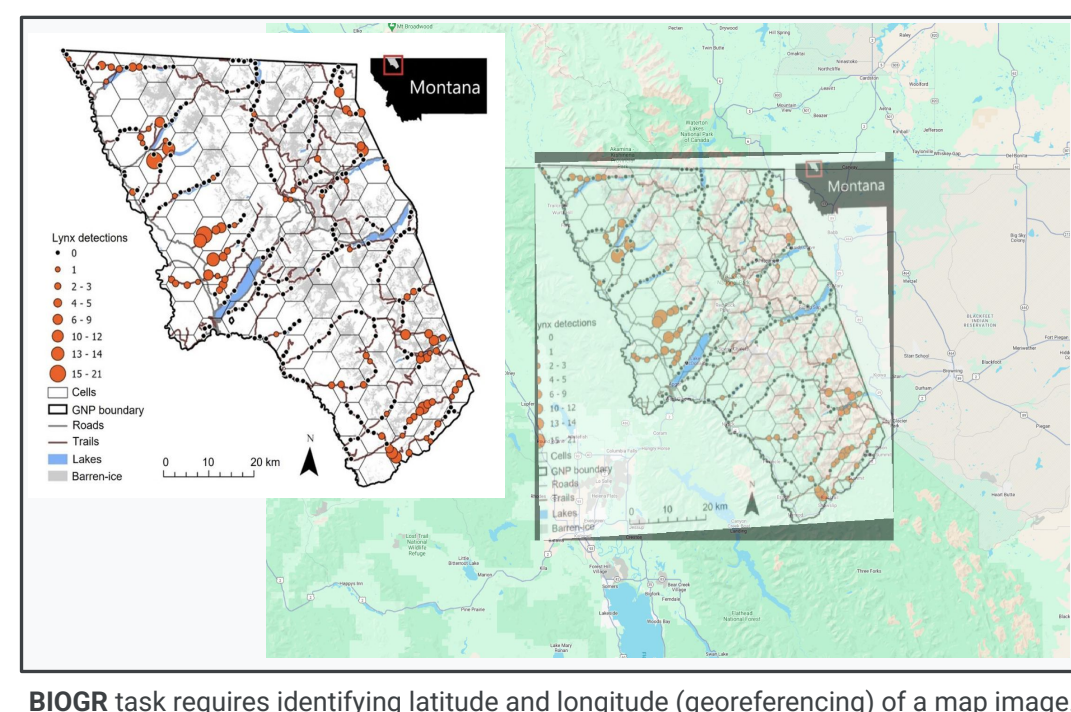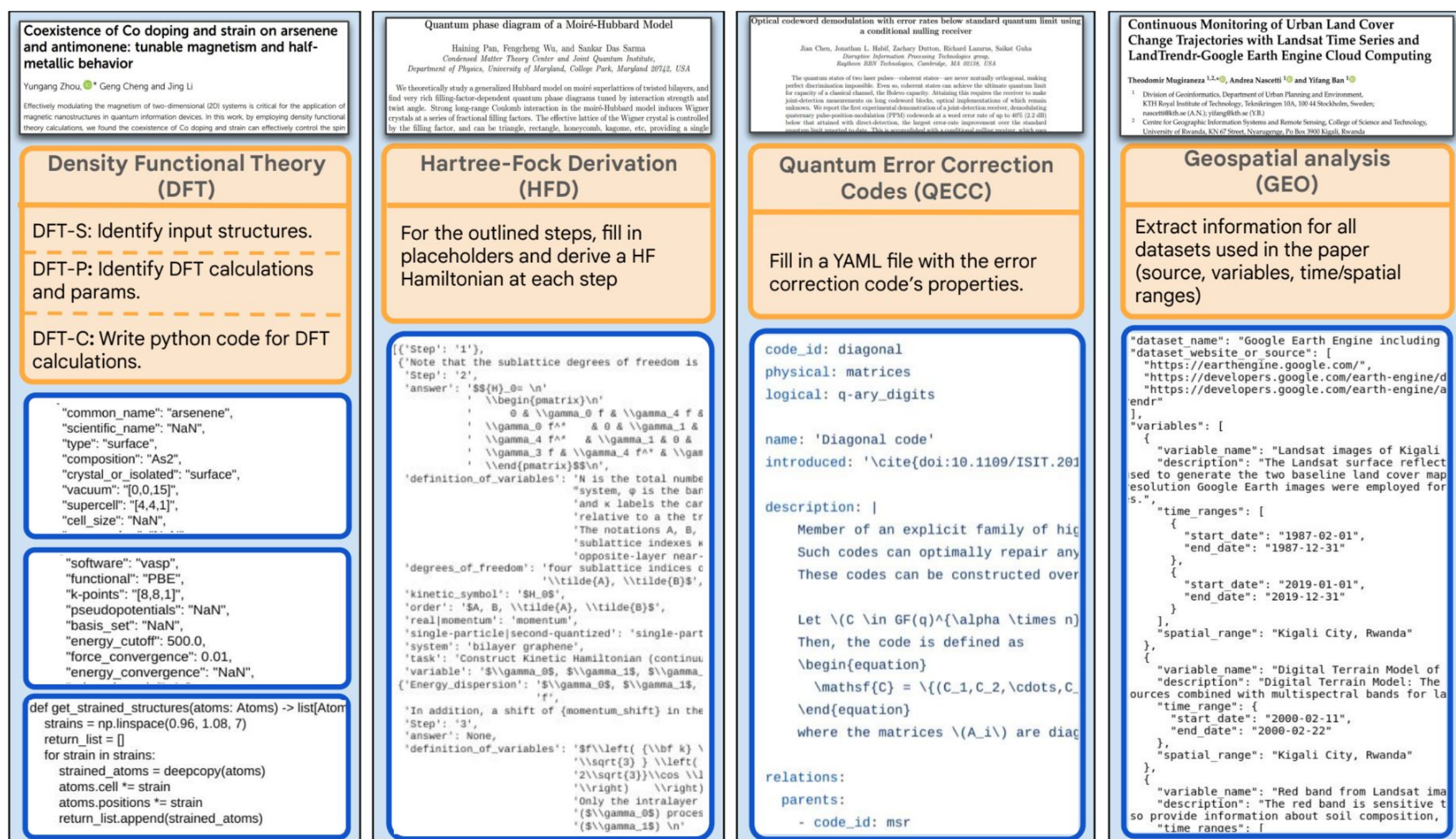
## Results

Per task normalized scores of various LLMs on the CURIE benchmark to measure performance on 10 long-context tasks requiring expertise across six scientific disciplines. Higher is better.

- Highest at 32%
- Gemini Flash 2 and Claude-3 do better. They understand the purpose of the extraction tasks and group responses.
- Exhaustive retrieval e.g. DFT-S, MPV, GEO are challenging for all.
- Flash 2 generated code to solve PDB ~50% of the time and was correct. When enumerating sequence it was similar to other models.
- Experts noted that summaries generated by models were succinct while including a multitude of details hard to comb out e.g. in QECC, and easy to remove afterwards.
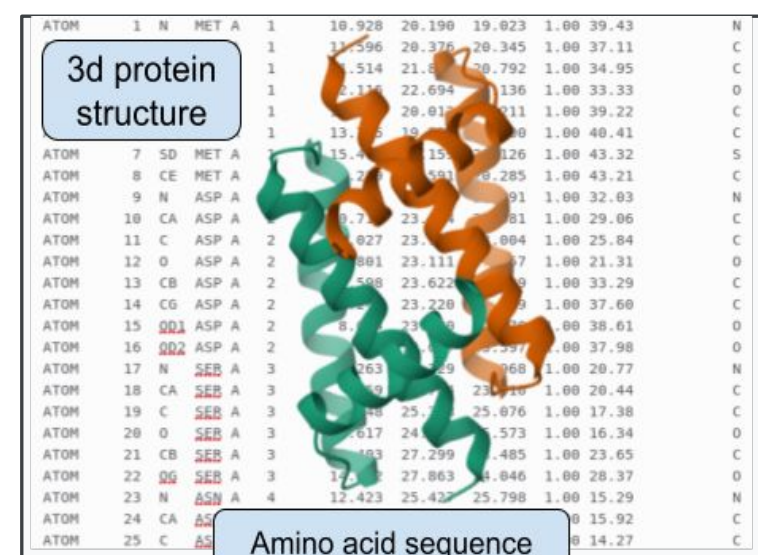
Average normalized performance of long-context LLMs across the 10 tasks from six scientific domains in CURIE.

## Illustrative examples of tasks in CURIE

### Density Functional Theory (DFT)
- DFT-S: Identify input structures.
- DFT-P: Identify input DFT calculations and params.
- DFT-C: Write python code for DFT calculations.

### Hartree-Fock Derivation (HFD)
For the outlined steps, fill in placeholders and derive a HF Hamiltonian at each step.

### Quantum Error Correction Codes (QECC)
Fill in a YAML file with the error correction code's properties.

### Geospatial analysis (GEO)
Extract information for all datasets used in the paper (source, variables, time/spatial ranges)

BIOGR task requires identifying latitude and longitude (georeferencing) of a map image.

PDB task requires reconstructing a protein's amino acid sequence from the 3D structure.

## Brief description of tasks and capabilities measured

| 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 |
| 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 |
| QECC | Quantum Computing | 65 | Create a YAML file with the Error Correction Code's properties. | concept aggregation, summarization | YAML | 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, aggregation | JSON (lat. lon. co-ordinates) | Intersection-over-Union (IoU) |
| PDB | Protein Sequencing | 138 | Reconstruct a protein's amino acid sequence form the 3D structure. | tracking, aggregation, reasoning | Code or TEXT (seq.) | Identity ratio (IDr) |

Tasks in CURIE are varied and have ground truth annotations in mixed and heterogeneous form e.g. as JSONs, latex equations, YAML files, or free-form text. We use programmatic metrics for a predominant number of tasks and propose LLM- based evaluation for others.

## Examples

### DFT-P

### DFT-S

### HFE

### MPV

### QECC

### PDB