

Can Large Language Models *Transform* Computational Social Science?

Caleb Ziems^{1*}, William Held^{2*}, Omar Shaikh^{1*}, Jiaao Chen^{2*}, Zhehao Zhang^{3*}, Diyi Yang¹

¹Stanford University, ²Georgia Institute of Technology, ³Dartmouth

* All substantially contributed to the implementation of this work



Overview

We evaluate **13 LLMs** on **25 CSS tasks** zero-shot and draw a road map 🚗🗺️ for social scientists who want to use them.

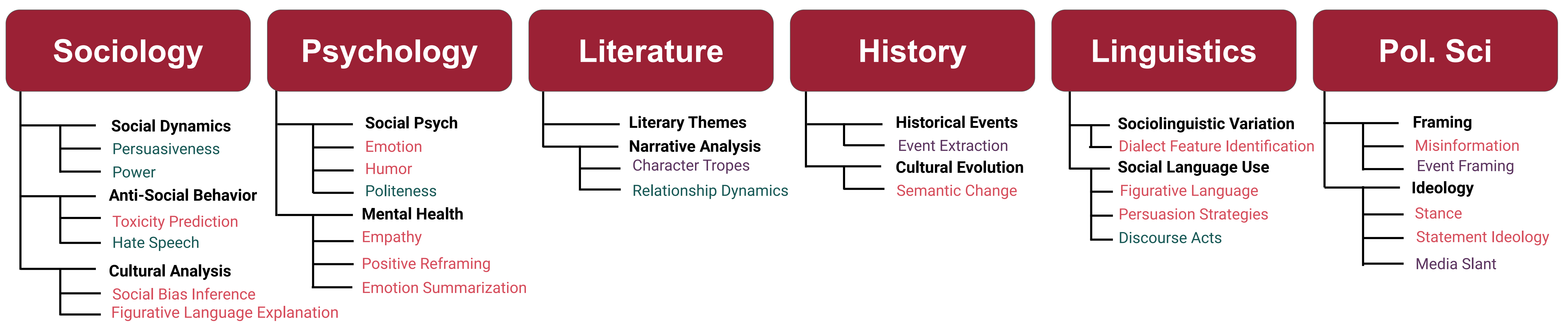
White House Ousts Top Climate Change Official

Which of the following describes the above news headline?

A: Misinformation ↩️

B: Trustworthy ↩️

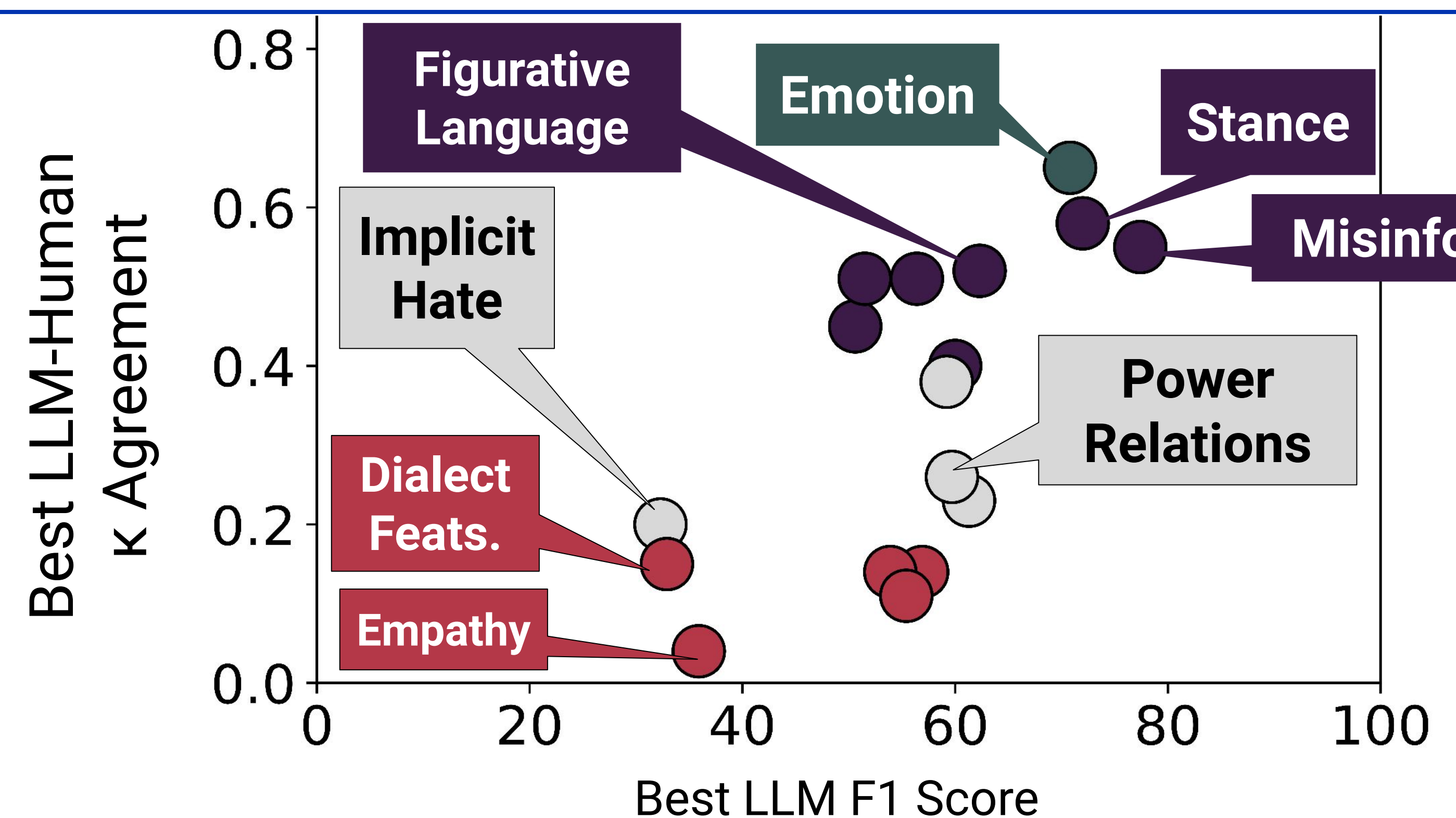
Constraint: Answer with only the option above that is most accurate and nothing else.



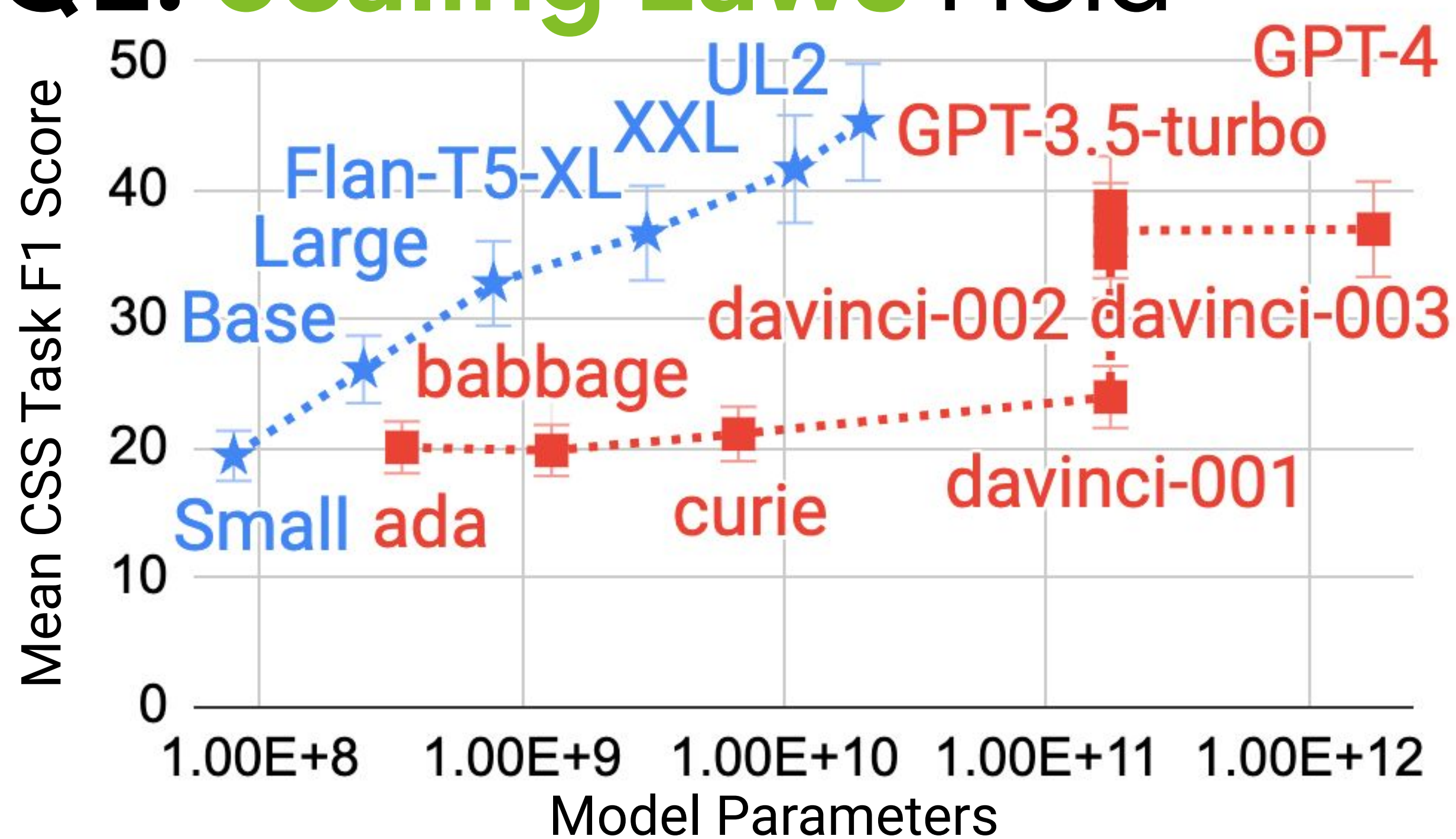
RQ1: LLMs can **Augment** but **not Replace** Human Annotation

Augment? Yes ✅ for ~half of tasks with $\kappa > 0.4$ (fair or good agreement with humans)

Replace? No ❌ especially not for expert taxonomies (Dialect Feats.) or parsing tasks



RQ2: **Scaling Laws** Hold

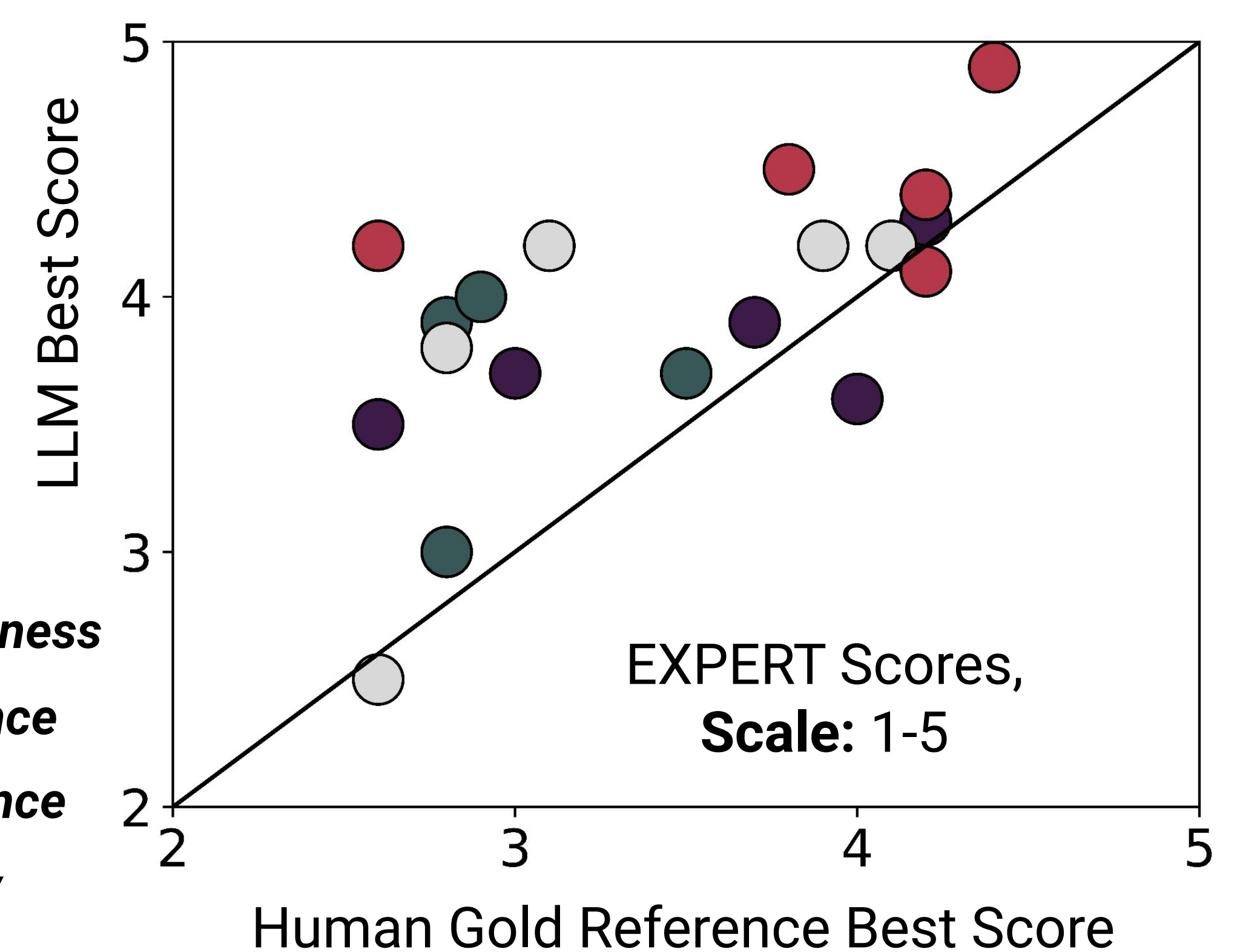


RQ4:

LLMs Can Also **Generate New Labels**

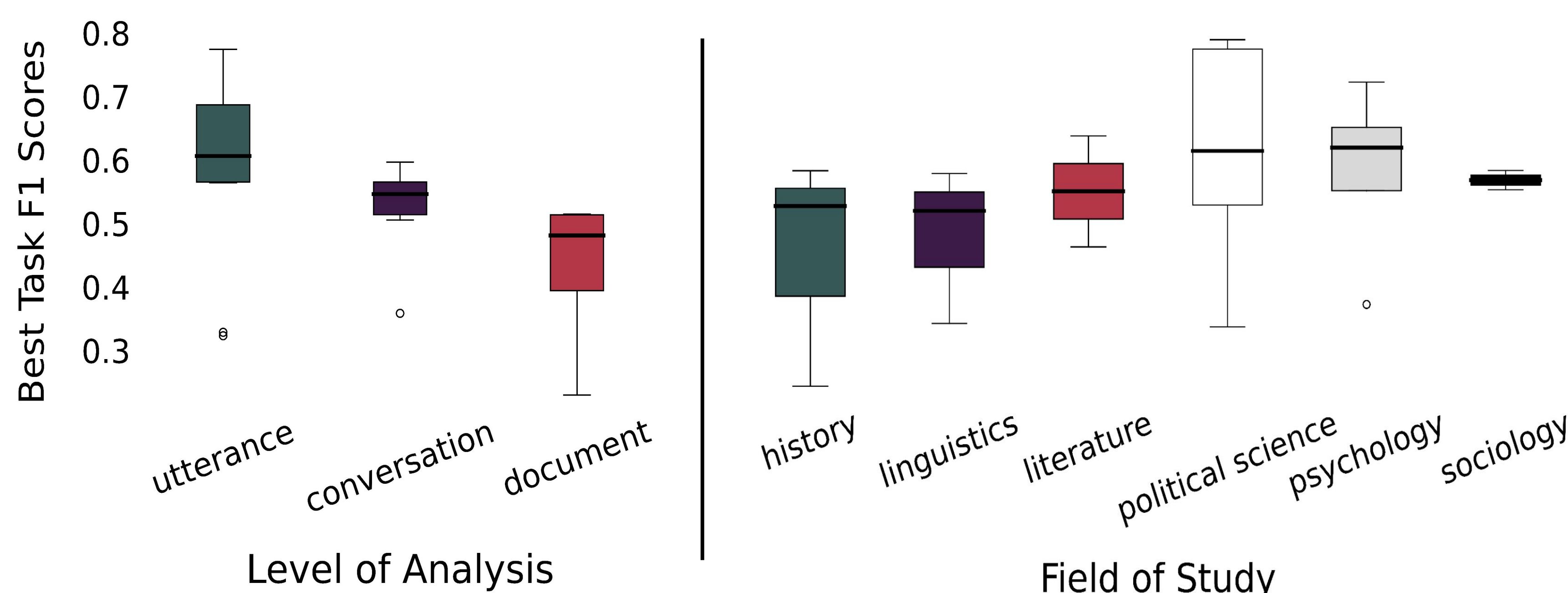
that **beat** reference levels of:

- Faithfulness
- Relevance
- Coherence
- Fluency



RQ3:

Performance Varies by **Complexity**, **not Field of Study**



Recommends:

1. LLM + Human co-annotation
2. Open Source models for classification
3. Reinvest in experts for reproducible CSS