# Hidden Non-Determinism in LLM APIs
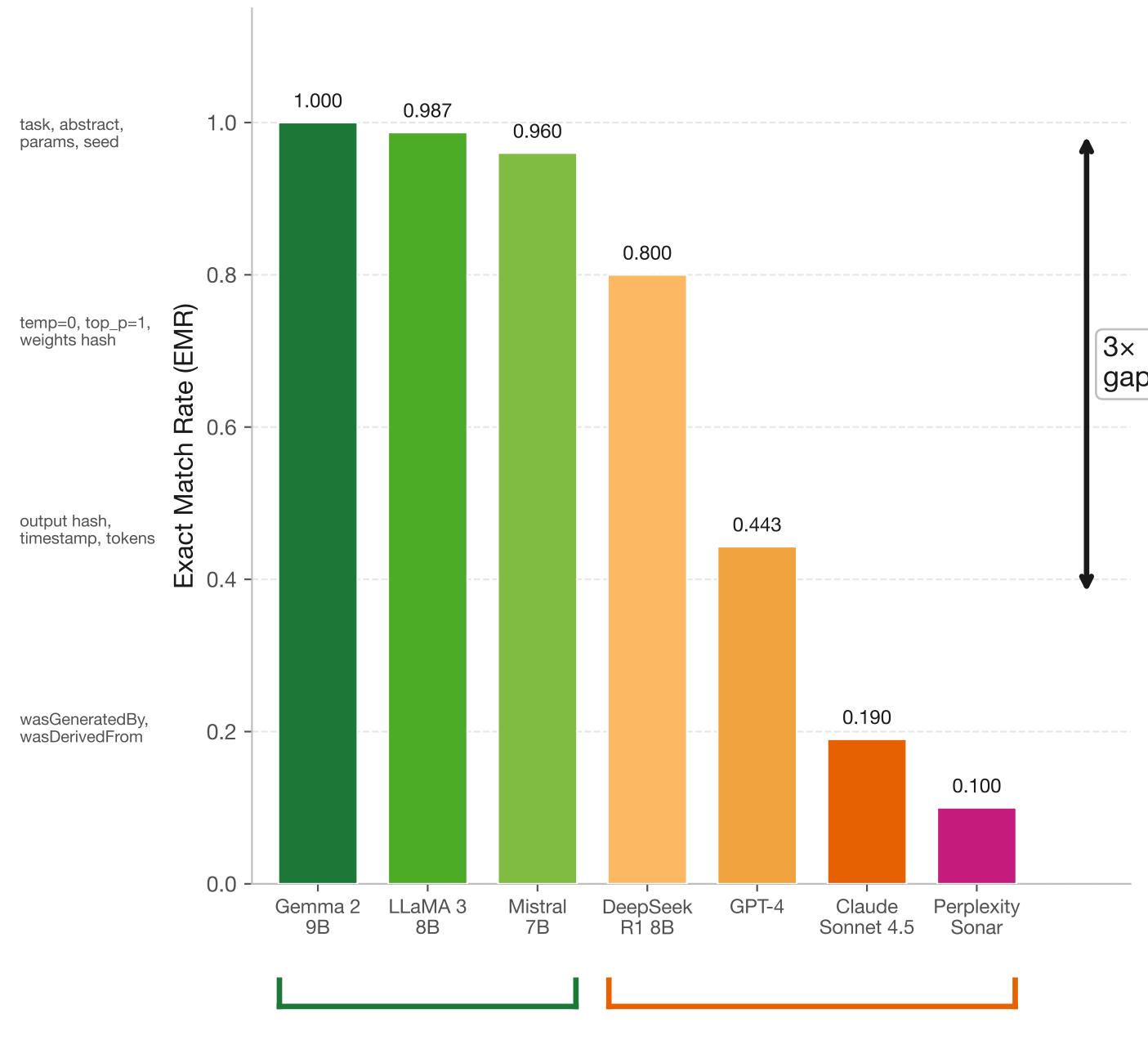
## A Lightweight Provenance Protocol for Reproducible Generative AI Research

### Protocol Pipeline

**Prompt Card** — task, abstract, params, seed

**Model Execution** — temp=0, top_p=1, weights hash

**Run Card** — output hash, timestamp, tokens

**PROV Graph** — wasGeneratedBy, wasDerivedFrom

**Reproducibility Verification**

### Greedy-Decoding Reproducibility (Extraction Task)



| Model | Exact Match Rate (EMR) |
|---|---|
| Gemma 2 9B | 1.000 |
| LLaMA 3 8B | 0.987 |
| Mistral 7B | 0.960 |
| DeepSeek R1 8B | 0.800 |
| GPT-4 | 0.443 |
| Claude Sonnet 4.5 | 0.190 |
| Perplexity Sonar | 0.100 |

3× gap

LOCAL — avg EMR = 0.982

API — avg EMR = 0.383

**github.com/Roverlucas/genai-reproducibility-protocol**

### Key Results

**3,804** total runs

**7** LLM models

**4** tasks

**4** API providers

**<1%** overhead

Local models reproduce near-perfectly; API outputs vary across identical calls.