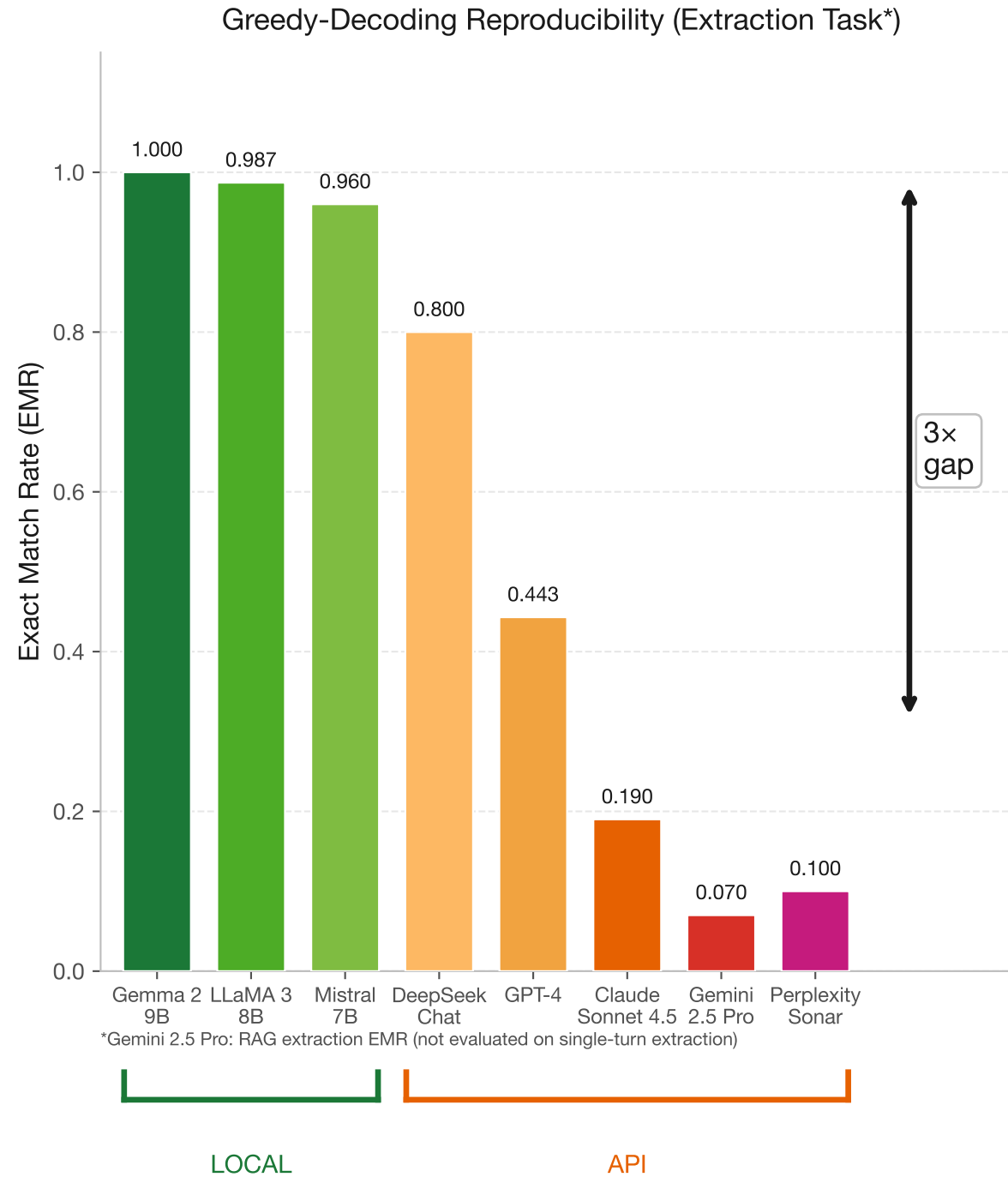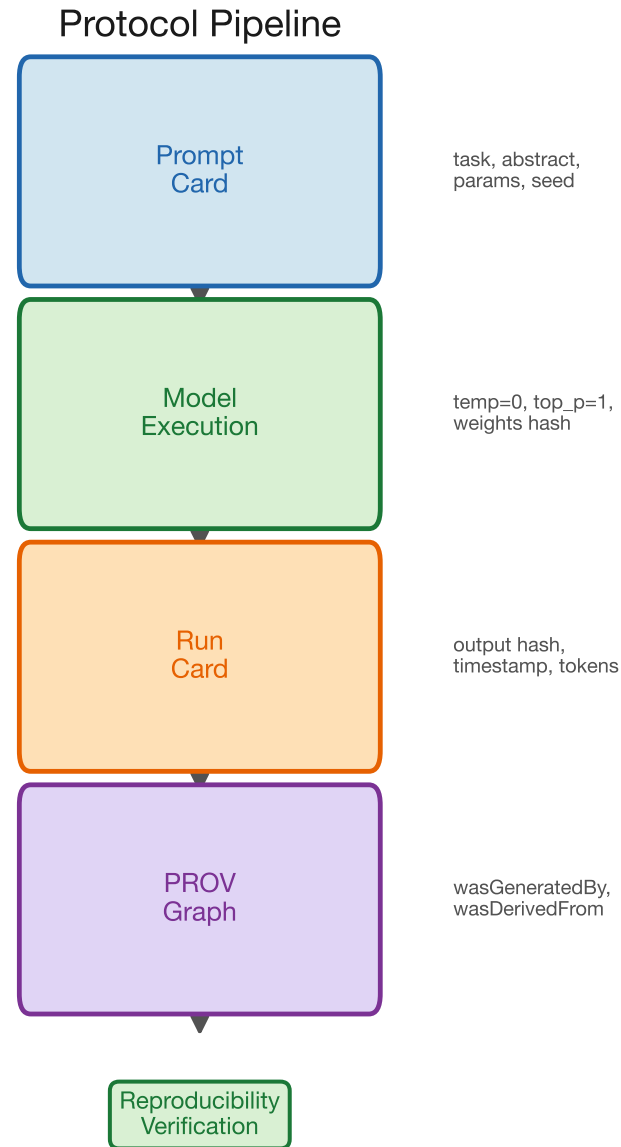# Hidden Non-Determinism in LLM APIs

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*)

Exact Match Rate (EMR)

| Model | EMR |
|---|---|
| Gemma 2 9B | 1.000 |
| LLaMA 3 8B | 0.987 |
| Mistral 7B | 0.960 |
| DeepSeek Chat | 0.800 |
| GPT-4 | 0.443 |
| Claude Sonnet 4.5 | 0.190 |
| Gemini 2.5 Pro | 0.070 |
| Perplexity Sonar | 0.100 |

3× gap

*Gemini 2.5 Pro: RAG extraction EMR (not evaluated on single-turn extraction)

LOCAL — avg EMR = 0.982
API — avg EMR = 0.321

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

## Key Results

- **4,104** total runs
- **9** model deployments
- **4** tasks
- **7** providers
- **<1%** overhead

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