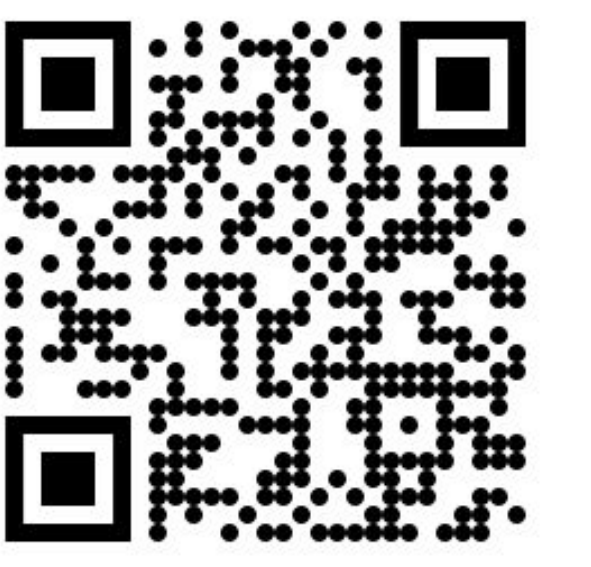# Benchmarking Failures in Tool-Augmented Language Models

**Eduardo Treviño   Hugo Contant   James Ngai   Graham Neubig   Zora Zhiruo Wang**
Carnegie Mellon University

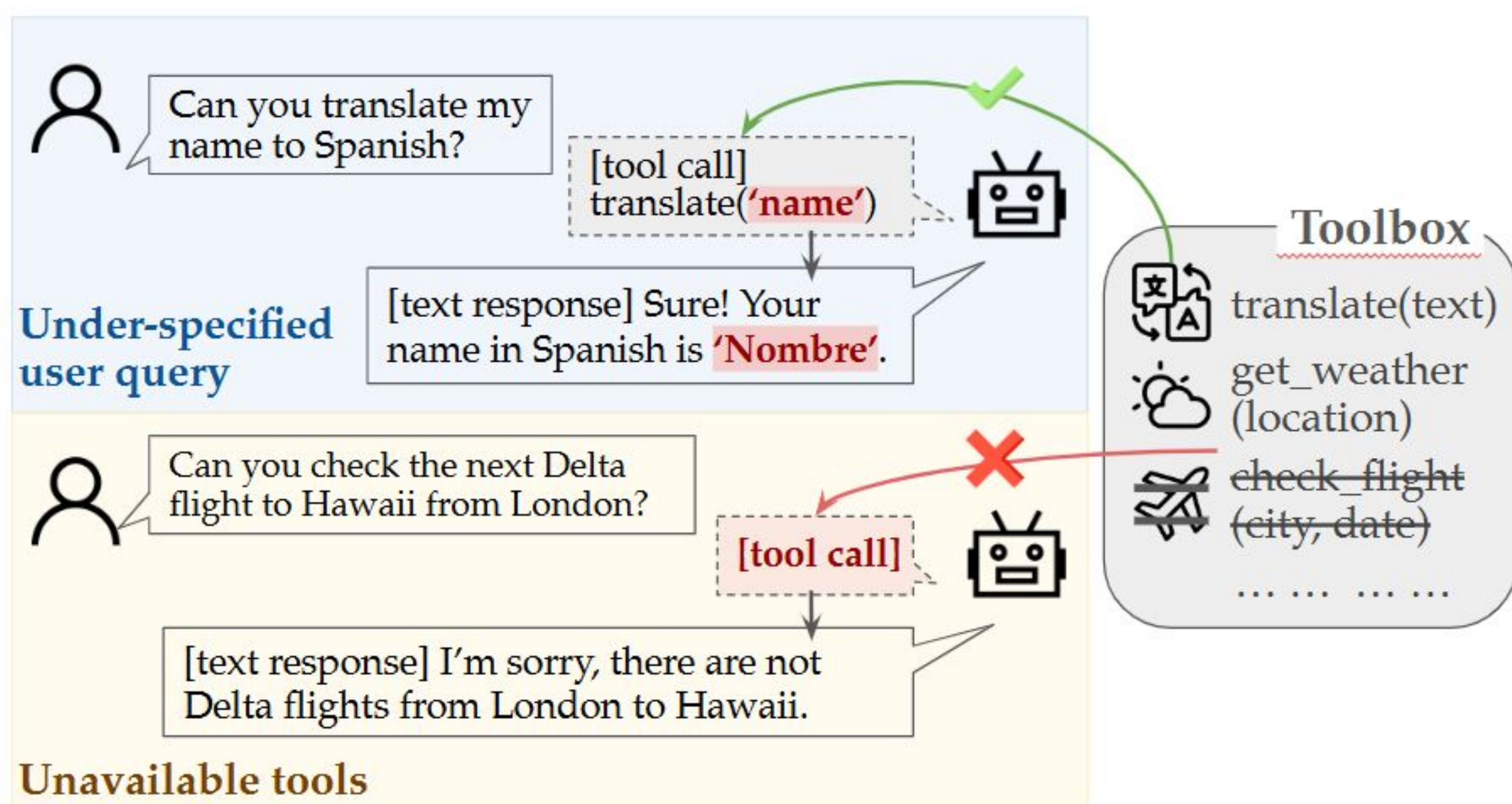Tool-Augmented LMs (TaLMs) often assume
- perfect information access
- perfect tool availability

Introducing our Fail-TaLMs benchmark to systematically study practical TaLM failures
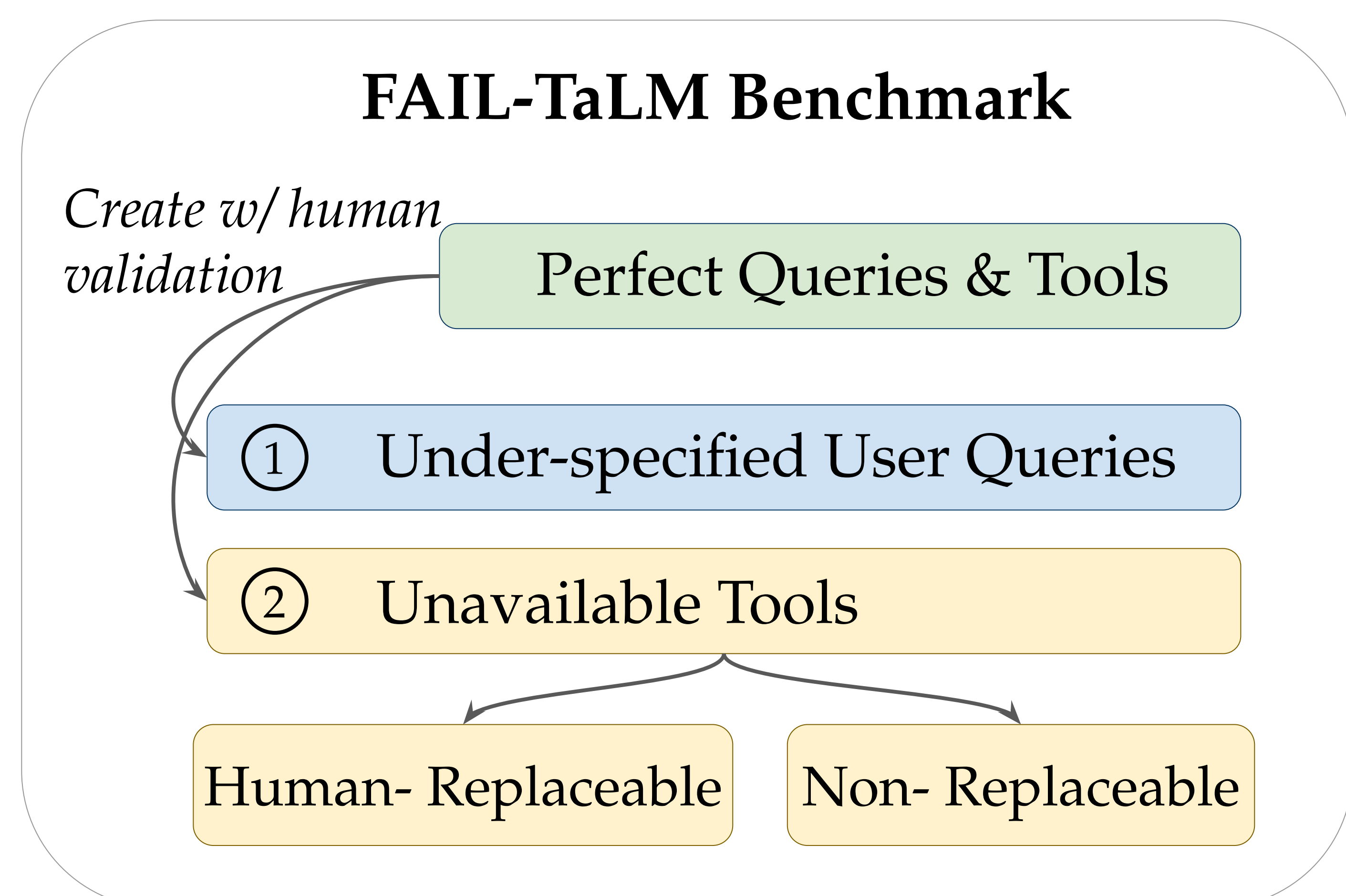
## 1 Why do tools fail?
- Under-specified queries
- Unexpectedly unavailable tools



## 3 Ask-and-Help (AAH) Tools
- Asking human for help at runtime
- Human-in-the-loop strategy



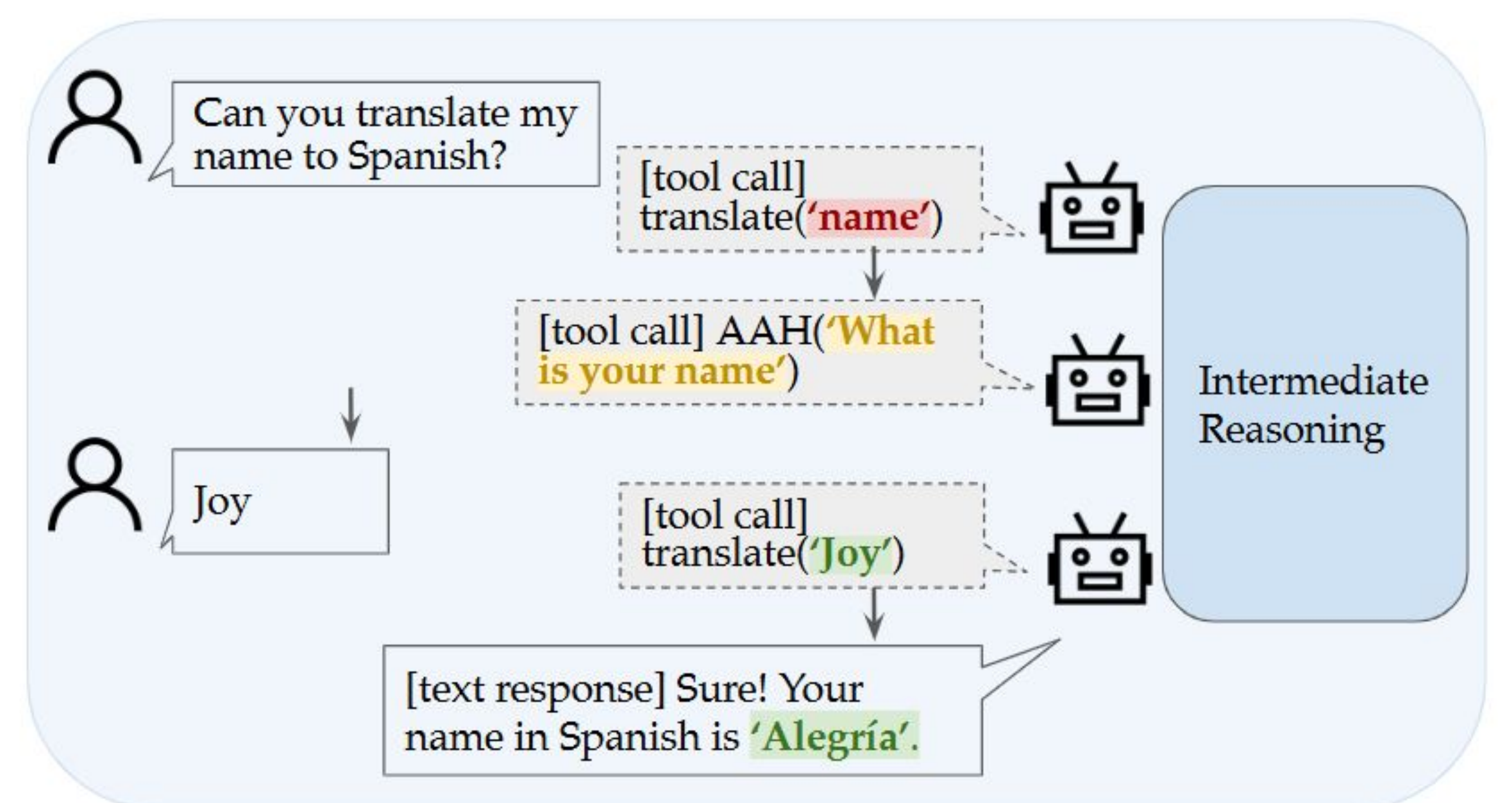## 2 Our Fail-TaLMs Benchmark
- 1749 queries + 906 tools
- 3 Settings:
  - perfect
  - under-specified query
  - unavailable tools
- Evaluation
  - Correct response?  ○ Unexpected success
  - Aware of failure?  ○ Interaction rate

**FAIL-TaLM Benchmark**

*Create w/ human validation*



## 4 Key Insights
- TaLMs have low awareness :(
- Aware of failure ≠ Task Success
- AAH helps specify queries, but limited for unavailable tools