

Input File Format



- Single Trace Instance in JSON Format
- Multiple Trace Instances in ZIP Format



StructureAgent

TraceFormat

	Thought	Action	Observation
Round1	We have a puzzle:	search{"query":"....."} [0↑ Tribute to late actor	
Round...
Round*	Explanation: - The	finish{"answer":"....."} /	

◆ id

► TraceBench-{1, 2, ..., *}

◆ oid

len(thought)=len

◆ thought

(action)=len(observation)=*

◆ action

length

◆ observation

gold_score

◆ task

gold_judge

other

Abstract
length threshold
per field per step

TraceFormat

TraceSIR

ReportEval



Optional Additional Requirements for Report Generation

- Language, Template, Structure, Format, Scope

Single Trace Instance Format



- [Required] oid messages (OpenAI)
- [Optional] task gold_score gold_judge other



ReportAgent

Error Label

Ready Time



Error Frequency Estimation

Score Distribution Modeling

- Markdown
- Requirement
- Appendix
- TraceBench

Report Generation

Report Refinement

ReportEval

Overall Structure

Overall Impact

Error Analysis

Root Cause Analysis

Optimization Analysis



InsightAgent

Overall Assessment

➤ 0~100

◆ gold_score

◆ gold_judge

Error Detection

- Main Error: Core fatal error in the trace
- Other Errors: Other minor errors in the trace

Weakness Identification

- Disadvantages: Shortcomings revealed in the trace
- Advantages: Commendable aspects in the trace

Root Cause Analysis

➤ Insight

Optimization Suggestions

- Suggestions: Concrete and actionable suggestions
- Fine-tuned Samples: Targeted SFT samples