

# What is AI reasoning?

## Position: Reasoning is a Learnable Rule-Based Process

Rachel Lawrence<sup>1\*</sup>, Jacqueline Maasch<sup>2\*</sup>

\*Equal contribution <sup>1</sup>Microsoft Research, Cambridge, UK <sup>2</sup>Cornell Tech, New York, NY

<https://bit.ly/ai-reasoning>



## Core Positions

**Thesis 1 Define, then measure.** Ambiguous, overloaded, or absent definitions for AI reasoning have led to confusion and mismeasurement. *Operational definitions* should be stated for the reasoning phenomena under investigation, and the *construct validity* of reasoning evaluation should be justified.

**Thesis 2 Reasoning is a learnable rule-based process.** Reasoning is a process of exact rule application, not an output. Learnable rules unambiguously map reasoning inputs to outputs and can encompass theorems, axioms, functions, policies, decision boundaries, laws, etc., including rules pertaining to stochasticity, uncertainty, and approximation.

**Thesis 3 Rule-based reasoning is valid.** The validity of a reasoning process arises from exact rule application, independent of rule selection.

## Operationalizing Reasoning

**Def. 1 (Reasoning, informal).** The process of selecting and applying sequences of rules that act on prior beliefs and current evidence to obtain principled belief updates in evolving states.

**Def. 2 (Reasoner).** A goal-oriented decision-maker that implements a reasoning process.

**Def. 3 (Reasoning, formal).** Let  $\mathcal{S}_t := \langle \mathcal{B}_t, \mathcal{E}_t, \mathcal{R}_t \rangle$  denote the reasoner's state at time step  $t$ , where  $\mathcal{B}_t$  denotes current belief,  $\mathcal{E}_t$  denotes aggregated evidence up to time  $t$ , and  $\mathcal{R}_t$  denotes the current set of established rules. Reasoning is the iterated application over steps  $t$  of rules  $r \in \mathcal{R}_{t-1}$  to prior beliefs  $\mathcal{B}_{t-1}$  and current evidence  $\mathcal{E}_t$ , by which we dynamically update states  $\mathcal{S}_t$ , and where every output  $\mathcal{B}_t$  for  $t > 0$  is the result of a rule application  $r(\mathcal{B}_{t-1}, \mathcal{E}_t)$  to the contents of state  $\mathcal{S}_{t-1}$ .

**Rules can be learned and revised.** Rules  $\mathcal{R}$  are operators whose operands are information: (1) extrinsically obtained information (**evidence**) and (2) intrinsically generated information (**beliefs**).

**Local rules**  $\mathcal{R}^L$  update beliefs. **Meta rules**  $\mathcal{R}^M$  update the rule set (Algorithm 1).

## Validity & Soundness

**Def. 4 (Validity).** A transition from state  $\mathcal{S}_t$  to  $\mathcal{S}_{t+1}$  is *valid* if and only if it arises from the application of a rule  $r \in \mathcal{R}_t$  to components of state  $\mathcal{S}_t$ .

**Def. 5 (Soundness).** A valid transition  $\mathcal{S}_t \rightarrow \mathcal{S}_{t+1}$  is *sound* if and only if all premises encoded by  $\mathcal{B}$ ,  $\mathcal{E}$ , and  $\mathcal{R}$  are true w.r.t. external evaluation.

---

### Algorithm 1 Valid reasoning as exact rule application.

---

**Input.** Initial rules  $\mathcal{R}_0$ , beliefs  $\mathcal{B}_0$ , evidence stream  $\{\mathcal{E}_i\}_{i=0}^T$ , stopping rule  $\mathfrak{s}_{\text{stop}}$ .

$\mathcal{R}, \mathcal{B}, \mathcal{E} \leftarrow \mathcal{R}_0, \mathcal{B}_0, \mathcal{E}_0$ ;  $\mathcal{S} \leftarrow (\mathcal{R}, \mathcal{B}, \mathcal{E})$ ;  $t \leftarrow 0$

**while not**  $\mathfrak{s}_{\text{stop}}(\mathcal{S})$  **do**

$\mathcal{E}' \leftarrow \mathcal{E}_{t+1}$  {Update extrinsic evidence.}

$r^L \leftarrow \mathfrak{s}_L(\mathcal{R}, \mathcal{B}, \mathcal{E}')$  {Select local rule.}

$\mathcal{B}' \leftarrow r^L(\mathcal{B}, \mathcal{E}')$  {Apply local rule, update beliefs.}

$r^M \leftarrow \mathfrak{s}_M(\mathcal{R}, \mathcal{B}', \mathcal{E}')$  {Select meta rule.}

$\mathcal{R}' \leftarrow r^M(\mathcal{R}, \mathcal{B}', \mathcal{E}')$  {Apply meta rule, update rules.}

$\mathcal{S}' \leftarrow (\mathcal{R}', \mathcal{B}', \mathcal{E}')$

$\mathcal{T}.\text{append}(\text{tr}(\mathcal{S}, r^L, r^M, \mathcal{S}'))$  {Update auditable trace  $\mathcal{T}$ .}

$\mathcal{R}, \mathcal{B}, \mathcal{E}, \mathcal{S} \leftarrow \mathcal{R}', \mathcal{B}', \mathcal{E}', \mathcal{S}'$

$t += 1$

Return  $\mathcal{B}, \mathcal{T}$  {Return belief set and reasoning trace.}

---