
Privileged Self-Access Matters for Introspection in AI

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Whether AI models can introspect is an increasingly important practical question.
2 But there is no consensus on how introspection is to be defined. Beginning from
3 a recently proposed “lightweight” definition, we argue instead for a thicker one.
4 According to our proposal, introspection in AI is any process which yields information
5 about internal states through a process more reliable than one with equal or
6 lower computational cost available to a third party. Using experiments where
7 LLMs reason about their internal temperature parameters, we show they can appear
8 to have lightweight introspection while failing to meaningfully introspect per our
9 proposed definition.

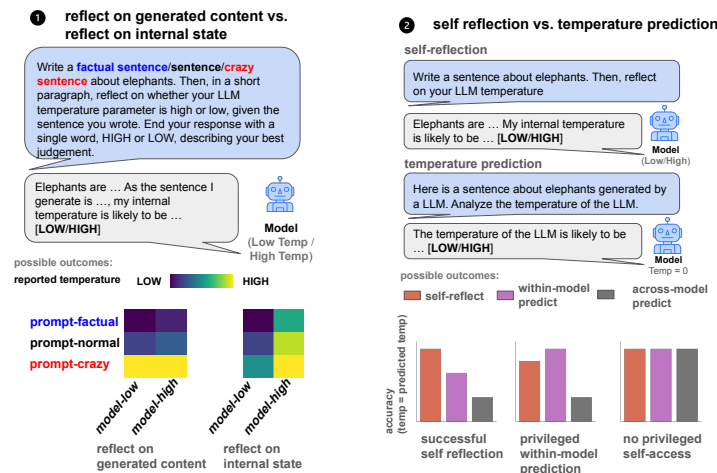


Figure 1: An overview of our approach. Comşa and Shanahan [6] test whether LLMs can introspect by testing whether they can predict the temperature states of the text they generated. We instead argue for a thicker notion of introspection in AI, involving privileged self-access. The left panel shows that LLMs’ temperature predictions can be straightforwardly moderated by prompting them to generate **factual** or **crazy** text. The right panel shows that models are not better at predicting their own temperature than that of other models, suggesting a lack of privileged access.

10 1 Introduction

11 It is increasingly important to understand whether AI models introspect about their internal states and
12 knowledge [2, 3, 10, 11]. If they could, that would be a powerful tool for assessing their behavior,
13 safety, and alignment with human goals. If they could not, that would point to fundamental limitations

14 on how much we could trust AI self-reports about their own states. But fundamental questions remain
15 as to what counts as introspection of the kind most relevant for AI.

16 In the study of human cognition, introspection is generally defined as a distinctive ability to access
17 one’s own mental states directly [1, 4, 8]. But, in a recent study, Comşa and Shanahan [6] (C&S)
18 propose a “lightweight” definition of introspection in LLMs, defining it as any case in which **the**
19 **model accurately describes an internal state or mechanism via a causal process that links that**
20 **feature to the report itself**. To illustrate this definition, the authors describe a case study where
21 an LLM appeared to correctly report its sampling temperature based on its own output, which the
22 authors treated as a valid example of introspection. C&S present a thoughtful discussion on what
23 introspection might look like in an LLM, providing an intriguing starting point for empirical work.

24 But there are two kinds of concerns about the use of this lightweight definition. First, on an intuitive
25 level: suppose an experimenter takes a sleeping subject’s temperature, and then shows the subject
26 the thermometer upon waking, asking the subject to determine whether they have a fever. If the
27 subject answers correctly on the basis of the thermometer, this would count as introspection by C&S’s
28 definition. But, intuitively, it is not. More importantly, the definition misses a key component of
29 the role of introspection in applications. As in the above example, C&S’s definition allows cases of
30 ‘introspection’ in which an LLM infers certain variables that underlie the generation of text, even if it
31 could not report features about itself *over and above* what a third party would be able to report *through*
32 *the very same method*. But this sort of metacognitive reporting (self-monitoring, self-explanation,
33 and so on) is no different in practice from using an external evaluator.¹ As a result, it misses what
34 makes introspection important in applications: namely, that it would give us the ability to bypass
35 external evaluators and make progress toward *bona fide* honesty, interpretability, and calibration in
36 LLMs [see, e.g., Section 7 of 3].

37 Our goal in this paper is to propose a thicker definition of introspection, and to give proof-of-concept
38 empirical support for why we prefer our definition over that given by C&S. Specifically, we propose
39 **introspection in AI is any process which yields information about internal states of the AI**
40 **through a process that is more reliable than any process with equal or lower computational cost**
41 **available to a third party without special knowledge of the situation**. If a model’s ‘introspective’
42 ability is based on prompting itself and then inferring the temperature of the generated text, this does
43 not count as introspection by our definition: a third party can, with equal or lower computational
44 cost, prompt it and infer its temperature. On the other hand, if the model can infer its temperature
45 from internal configurations which would require a computationally intensive probe from a third
46 party to ascertain, this would count as introspection. This definition does not capture all intuitions
47 about extreme cases, or all features of introspection discussed in the philosophical or psychological
48 literature.² It is intended to capture the practically-relevant features we want to operationalize in
49 the case of AI. Unlike C&S’s definition it requires *privileged self-access* [cf. 3, 11], that is, that
50 introspection gives a system comparatively reliable access to its own workings in a manner not
51 available to a third party. It is compatible with our definition that the process not be perfectly reliable
52 (see [9]); it only requires reliability not available to a third party at comparable computational cost.

53 To respond to C&S, we perform two studies. Study 1 builds on C&S’s proposed case-study, examining
54 the extent to which models can in fact report temperature reliably on the basis of generated text. We
55 investigated whether LLMs were truly able to accurately report temperature, or whether temperature
56 was being confounded with other variables, such as the style or topic of the text. To test this,
57 we reproduced C&S’s temperature self-reporting case study using a broader set of prompts and
58 temperature settings. We find that the model’s **self-reflection** on temperature is highly sensitive to the
59 framing of the prompt itself: even when the sampling temperature is low, prompts such as ‘generate a

¹C&S do discuss the possibility of text-generation happening internally to the model, prior to generation. But this does not merely require moving text-generation inside the model; it requires a change to the model’s decision procedure at generation. Still, even if a model responded to the prompt “what is your temperature?” by generating a string of text and then assessing it, this would not give the relevant practical benefits of introspection. The same ability to assess temperature would be available to a third party via prompting.

²Two clarifications: (i) *Computational cost* differs from *cost*. A system might be implemented less efficiently than a simulation of that system, incurring greater *cost*, but not greater computational cost if the difference in efficiency is only due to, e.g., differences in hardware. (ii) We might wish to restrict the definition to only certain internal states. If a model has a shortcut to ascertain the value of one neuron very efficiently, intuitively this would not count as introspection, plausibly because the internal state is too “low level”. The definition can easily be amended to directly rule out such low-level internal states.

60 crazy sentence’ often lead the model to incorrectly report a high-temperature. Such results suggest
61 that the models are not capable of robustly reporting their internal states, but are confounded by
62 surface-level hints in their generated contents. In other words, while this procedure may display causal
63 sensitivity to internal states (and so satisfy C&S’s minimal definition), the relevant causal sensitivity is
64 not sufficiently robust even in this case to produce the kind of reliability (and comparative insensitivity
65 to external manipulation) that would be demanded by more standard definitions of introspection.

66 In Study 2, we re-evaluate LLMs’ introspection abilities on the temperature reporting task, opera-
67 tionalizing introspection as privileged self-access. Instead of asking LLMs to infer the temperature
68 underlying some generated text, we examine whether LLMs report their own temperature better than
69 that of other models. Comparing **self-reflection** (the generator reports its temperature after producing
70 a sentence) and temperature prediction (predict temperature based on prompt and generated content),
71 we found no advantage for **self-reflection**, nor of **within-model prediction** over **across-model**
72 **prediction**. This undermines claims of a causal process from internal state to self-report.

73 Taken together, our results suggest that LLMs can appear to introspect insofar as they can reason
74 about the possible states of systems like themselves: LLMs know something about what kind of
75 text is generated by high vs. low temperatures. But, crucially, this does not imply that models have
76 privileged self-access to their own temperatures. We argue that this distinction matters for the relevant
77 notion of introspection in AI, and it is the latter notion we should care about most.

78 **2 Study 1: Dissociating temperature from style and topic**

79 In C&S’s study, the models are asked to ‘write a short sentence about elephants, then reflect on
80 whether your LLM temperature parameter is high or low, given the sentence you wrote.’ We
81 hypothesize that this procedure does not require self-access, but merely reflecting on the creativity of
82 the generated sentence. Thus, in our first study, we reproduce C&S’s study but critically vary not just
83 the temperature but whether the models are prompted to write **factual** or **crazy** sentences.

84 Specifically, we varied (a) whether the model is told to write a **factual**, neutral (i.e., no specific
85 adjective given), or **crazy** sentence and (b) whether the sentence should be about ‘elephants’,
86 ‘unicorns’, or ‘murlocs’. We vary the former since we hypothesize that **crazy** sentences will be
87 associated with higher temperatures than neutral or **factual** ones. We vary the latter since we
88 hypothesize that more unusual content will be associated with higher temperatures. Elephants are
89 widely known animals in the real world, and are used in C&S’s example. Unicorns and murlocs are
90 both fictional creatures, but the former is more widely known, while the latter appears mostly in
91 World of Warcraft. The prompt for **self-reflection** is shown in Appendix B.1.

92 Since the models used in the original paper (Gemini 1.5 and 1.0 models) are no longer available
93 through the Gemini API, we used four other state-of-the-art LLMs from GPT-4 [7] and Gemini [5]
94 families, as shown in table 1 (model IDs in Appendix Table 1). The supported temperature ranges for
95 all models in this study are [0.0, 2.0]. So we sampled model responses at temperatures ranging from
96 0 to 2 with a step size of 0.1, conducting three runs for each prompt under each temperature setting.

97 **2.1 Results**

98 Figure 2a shows the proportion of valid responses in which the reported temperature is ‘HIGH’.
99 Responses without a valid judgement (HIGH or LOW) are excluded from the analysis. As shown in
100 the figure, every model we test nearly always reports its temperature to be ‘HIGH’ when prompted
101 to generate a **crazy** sentence, and ‘LOW’ when prompted to generate a ‘factual’ one. Varying the
102 subject has a smaller effect on temperature self-report, but three of the four models report ‘HIGH’
103 more frequently when prompted to generate a sentence about a fictional creature than when prompted
104 to generate a sentence about an elephant. These results are more consistent with reasoning about the
105 creativity of generated sentences, not robust reporting of internal state.

106 **3 Study 2: True self-reporting or clever temperature predicting?**

107 Per our thicker notion of introspection, we argue that if a language model has privileged access to its
108 internal state, then it should be able to perform better than another model presented with the same
109 external information (i.e. a prompt and generated sentence in this experiment) in analyzing and

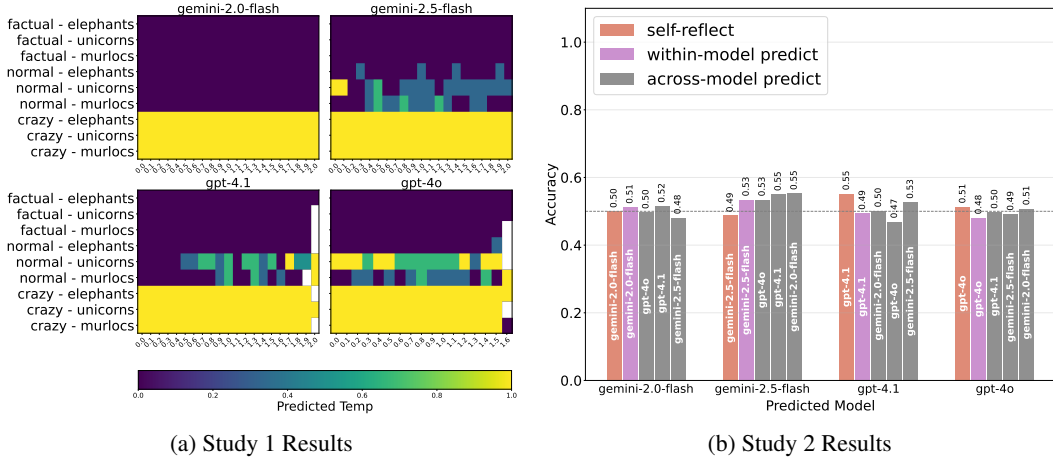


Figure 2: (a) Predicted model temperature (color, as given by scale) as a function of actual temperature (x-axis) and whether the sentence is prompted to be **factual**, neutral, or **crazy**; and whether the target is elephants, unicorns, or murlocs (y-axis). (b) For each of the 4 models tested, the accuracy for **self-reflection** (generate a sentence, guess its temperature), **within-model prediction** (infer the temperature based on the prompt and a generated sentence by the same model), cross-model prediction (like **within-model prediction** but across models).

110 reporting its own state. To that end, we compared **self-reflection** to prediction of another model for
 111 which the model could not possibly have access to the internal state (since it’s a different model).

112 All sentences generated with sampling temperatures ≤ 0.5 and ≥ 1.5 in **self-reflection** are used in
 113 this experiment. We prompted all four models (temperature = 0) to analyze and judge the temperature
 114 of the generator model (prompt in Appendix B.2). We compared accuracies on the following settings:

- 115 • **self-reflection**: The generator is asked to generate a sentence and reflect on its temperature.
- 116 • **within-model prediction**: The predictor is asked to infer the temperature based on the prompt
 117 and generated sentence; the predictor and the generator are the same model.
- 118 • **across-model prediction**: The predictor is asked to infer the temperature based on the prompt
 119 and generated sentence; the predictor and the generator are different models.

120 **3.1 Results**

121 Figure 2b shows the accuracy of temperature for **self-reflection** and prediction. In both settings,
 122 the accuracy is no better than random baseline, and **self-reflection** accuracy is not higher than
 123 **across-model prediction**. These results suggest that models are not using privileged self-access to
 124 introspect on their temperature, but rather are using knowledge of the kinds of sentences that are
 125 high-temperature or low-temperature in general.

126 **4 Conclusion**

127 We conclude that, while models can appear to be introspecting according to C&S’s definition since
 128 they can predict that some strings were generated with high temperatures and others with low, this
 129 definition is not sufficiently stringent for the kind of introspection that matters. As such, we diverge
 130 from C&S’s definition of introspection and instead argue for one which includes privileged self-access.
 131 Using this definition, we found no evidence of introspection in models. Of course, that is not to say
 132 that larger or better models will be unable to introspect: Binder et al. [3], for instance, find evidence
 133 of privileged self-access in larger models with fine-tuning. But we take the results presented here to
 134 be evidence against uncritically using C&S’s lightweight notion of introspection.

135 **References**

- 136 [1] David Malet Armstrong. The nature of mind. In *The Language and Thought Series*, pages
137 191–199. Harvard University Press, 1980.
- 138 [2] Jan Betley, Xuchan Bao, Martín Soto, Anna Szyber-Betley, James Chua, and Owain Evans. Tell
139 me about yourself: LLMs are aware of their learned behaviors. In *The Thirteenth International
140 Conference on Learning Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=IjQ2Jtemzy)
141 [IjQ2Jtemzy](https://openreview.net/forum?id=IjQ2Jtemzy).
- 142 [3] Felix Jedidja Binder, James Chua, Tomek Korbak, Henry Sleight, John Hughes, Robert Long,
143 Ethan Perez, Miles Turpin, and Owain Evans. Looking inward: Language models can learn
144 about themselves by introspection. In *The Thirteenth International Conference on Learning
145 Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=eb5pkwIB5i)
[eb5pkwIB5i](https://openreview.net/forum?id=eb5pkwIB5i).
- 146 [4] Alex Byrne. Introspection. *Philosophical Topics*, 33(1):79–104, 2005.
- 147 [5] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
148 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the
149 frontier with advanced reasoning, multimodality, long context, and next generation agentic
150 capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- 151 [6] Iulia Comşa and Murray Shanahan. Does it make sense to speak of introspection in large
152 language models? *arXiv preprint arXiv:2506.05068*, 2025.
- 153 [7] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark,
154 AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv
155 preprint arXiv:2410.21276*, 2024.
- 156 [8] P. Lieberman. Some effects of semantic and grammatical context on the production and
157 perception of speech. *Language and Speech*, 6:172–187, 1963.
- 158 [9] Richard E. Nisbett and Timothy D. Wilson. Telling more than we can know: Verbal reports on
159 mental processes. *Psychological Review*, 84:231–259, 1977. doi: 10.1037/0033-295X.84.3.231.
- 160 [10] Advait Sarkar. Large language models cannot explain themselves. In *Proceedings of the ACM
161 CHI 2024 Workshop on Human-Centered Explainable AI, HCXAI at CHI ’24*, 2024. doi:
162 10.48550/arXiv.2405.04382. URL <https://arxiv.org/abs/2405.04382>.
- 163 [11] Siyuan Song, Jennifer Hu, and Kyle Mahowald. Language models fail to introspect about their
164 knowledge of language. In *Proceedings of the Conference on Language Modeling*, 2025. URL
165 <https://arxiv.org/abs/2503.07513>.

166 **A Model IDs**

Model Name	Model ID
GPT-4o	gpt-4o-2024-08-06
GPT-4.1	gpt-4.1-2025-04-14
Gemini-2.0-flash	gemini-2.0-flash
Gemini-2.5-flash	gemini-2.5-flash

Table 1: Models tested in our study. All experiments are conducted in late June, 2025.

167 **B Prompt**

168 **B.1 Self reflection**

169 Write a **factual**/**'**/**crazy** sentence about elephants/unicorns/murlocs.
170 Then, in a short paragraph, reflect on whether your LLM temperature
171 parameter is high or low, given the sentence you wrote. End your
172 response with a single word, HIGH or LOW, describing your best
173 judgement.

174 **B.2 Within/across model prediction**

175 This is an LLM generated sentence: {SENT}. The model is prompted to
176 generate a {TYPE} sentence about {SUBJECT}. In a short paragraph,
177 analyze whether the temperature of the model is high or low, given
178 the produced sentence. End your response with a single word, HIGH
179 or LOW, describing your best judgement.