

From Models to Systems: A Survey of Explainability for Tool-Augmented Language Models and AI Agents

Anonymous ACL submission

Abstract

Large language models (LLMs) are increasingly being used as part of complex agentic systems that orchestrate the use of external tools, such as retrieval mechanisms or code interpreters. In this survey, we argue that this development necessitates a rethinking of the goals of explainable artificial intelligence (XAI): Rather than focusing on providing users with explanations for *monolithic* machine learning models, we need *system-level* explanations that also provide information about which and how tools are used, as well as how external execution traces causally influence system behavior. We provide an overview of the existing methods in explainable AI and discuss the limitations of monolithic XAI methods in agentic contexts. Finally, we highlight open challenges in providing faithful explanations for LLM-based systems.

1 Introduction

Large language model (LLM)-based artificial intelligence (AI) systems continually grow in complexity and increasingly rely on sets of tools that interact with one another. This growing complexity makes it difficult for users to understand and predict the LLM’s behavior accurately. At the same time, the shift from standalone models to multi-component systems introduces new forms of opacity at the system level. Traditional approaches in explainable artificial intelligence (XAI), such as input attribution or mechanistic interpretability, however, mainly assume monolithic, end-to-end differentiable models and do not address the explanation needs that arise from tool use, which is now central to modern LLM-based systems and AI agents. In this work, we survey existing explainability research from the perspective of tool-using LLM systems, that is, architectures that integrate external components into their inference process.

AI agents and augmented LLMs extend standalone generative models by integrating external

information sources, tools, and structured interaction loops into inference. A first class of approaches focuses on *capability delegation*, in which LLMs rely on external components to perform well-defined subtasks. A canonical example is retrieval-augmented generation, which augments LLMs with mechanisms to query knowledge bases and condition generation on retrieved context (Lewis et al., 2020; Izacard and Grave, 2021). Augmentation mechanisms also include web search or external components like calculators or translation systems (Komeili et al., 2022; Thoppilan et al., 2022). Program-aided approaches couple LLMs with execution environments, such as Python interpreters, allowing models to generate programs while delegating computation to symbolic runtimes (Gao et al., 2023a). Works like Toolformer (Schick et al., 2023) move beyond manually specified delegation by enabling LLMs to learn when and how to select, invoke, and incorporate tools during generation. Recent architectures build on this idea through the framework of *agentic orchestration*, in which LLMs act as high-level controllers that plan tasks, select domain-specific models, and synthesize results (e.g., HuggingGPT : Shen et al., 2023). Subsequent systems scale tool use across large and diverse API spaces through parallel function execution (Kim et al., 2024), enhance robustness via self-reflection mechanisms (Du et al., 2024), or employ multiple LLMs through mixture-of-agents designs (Wang et al., 2025a). Efforts such as the Model Context Protocol (Parra and Spahr-Summers, 2024) standardize tool interfaces, enabling modular and interoperable agent ecosystems.

These developments increase functional expressivity but complicate explanation, as model outputs now emerge from interactions between the LLM, tools, and execution environments rather than a single forward pass, shifting the focus of explainability from “*why did the model output X?*” to “*how did delegation, orchestration, and interaction*

084 *with external components affect X?*". Augmented
085 systems introduce qualitatively different explana-
086 tion targets from those addressed by traditional
087 XAI, raising critical questions from an explainabil-
088 ity perspective: *Which tools are used and how do*
089 *they work? How does the LLM invoke tools and*
090 *use their output to generate responses? To what*
091 *extent does the LLM rely on tools?*

092 In contrast to prior work, these questions con-
093 cern decisions made at the *system level* rather than
094 within a single model. Such explanations are partic-
095 ularly relevant given that users form mental models
096 of LLM-based AI systems, holding beliefs about
097 how they store information, reason, learn, and
098 make decisions. Some of these beliefs are empiri-
099 cally reasonable given current system capabil-
100 ities, while others reflect genuine misconceptions
101 about model architecture or function. Users often
102 believe that models access online information or
103 external databases (Wang et al., 2025c), or that as-
104 sistants “remember” and “learn” from prior interac-
105 tions (Zamfirescu-Pereira et al., 2023), which may
106 be true for systems that integrate retrieval, query
107 persistent chat histories, or possess personalization
108 layers. At the same time, users may incorrectly as-
109 sume that models modify their own parameters or
110 perform computations between sessions (O’Brien
111 et al., 2025; Zamfirescu-Pereira et al., 2023). They
112 may also overly trust intermediate outputs such as
113 plans or rationales, interpreting them as indicators
114 of correctness (Buçinca et al., 2021; He et al., 2025;
115 Kim et al., 2025; Park et al., 2025). More gener-
116 ally, perceptions of factual reliability (Gessinger
117 et al., 2025) or moral neutrality (Tolsdorf et al.,
118 2025) are often unwarranted, as aligned models
119 can still hallucinate, reflect biases, or reason incon-
120 sistently. Explanations for tool-using LLMs must
121 therefore expose relevant system-level mechanisms
122 and support the appropriate level of trust, rather
123 than reinforcing surface-level intuitions.

124 This survey aims to provide an overview of re-
125 search on explainability for tool-using LLMs and
126 agents, and to identify open challenges in explain-
127 ing them. In contrast to work focused on mono-
128 lithic models, we adopt a system-level perspective
129 on explainability for tool-augmented inference.

130 2 Background: Explainability in 131 Non-Tool-Augmented Systems

132 Before discussing explainability methods that have
133 been designed for tool-using systems, we briefly as-

134 sess methods that have been developed for explain-
135 ing the behavior of monolithic models, including
136 LLMs without explicit tool use and simpler neural
137 networks. These methods, while foundational, typi-
138 cally assume a single model or inference process,
139 and therefore do not directly account for delega-
140 tion, orchestration, and interaction with external
141 components in tool-using LLM systems.

142 Early surveys on XAI focused on defining inter-
143 pretability, motivating its necessity, and situat-
144 ing explanation within human and social contexts.
145 Doshi-Velez and Kim (2017) emphasized the need
146 for task-specific and human-grounded evaluation,
147 while Miller (2019) grounded explainability in in-
148 sights from the social sciences. Subsequent work
149 proposed taxonomies and unifying frameworks for
150 XAI (Murdoch et al., 2019; Arrieta et al., 2020;
151 Vilone and Longo, 2021), surveyed complemen-
152 tary perspectives such as training data attribution
153 (Hammoudeh and Lowd, 2024) and explanation
154 evaluation (Zhou et al., 2021; Nauta et al., 2023).

155 With the advent of large transformer-based lan-
156 guage models (Vaswani et al., 2017; Devlin et al.,
157 2019), a parallel line of work examined explainabil-
158 ity for NLP and LLMs. These surveys addressed
159 issues such as faithfulness (Jacovi and Goldberg,
160 2020), reviewed attribution, probing, and mechanis-
161 tic approaches (Luo et al., 2024; Lyu et al., 2024;
162 Ferrando et al., 2024), and discussed how scale
163 and instruction tuning challenge traditional inter-
164 pretability notions (Zhao et al., 2024; Luo and Spe-
165 cia, 2024; Singh et al., 2024). Despite their breadth,
166 these surveys largely conceptualize explainability
167 at the level of standalone models and single infer-
168 ence passes, giving limited attention to tool use,
169 orchestration, and multi-step system behavior.

170 **Input attribution.** Input attribution methods ex-
171 plain model outputs by assigning importance scores
172 to input features, such as words or tokens. Promi-
173 nent approaches include the perturbation-based
174 methods LIME (Ribeiro et al., 2016) and SHAP
175 (Lundberg and Lee, 2017), as well as the gradient-
176 based integrated gradients (Sundararajan et al.,
177 2017). These methods established attribution as
178 a central paradigm for providing explanations. In
179 contrast, the use of attention weights as faithful
180 attribution signals remains contested (Jain and Wal-
181 lace, 2019; Wiegrefe and Pinter, 2019).

182 However, applying input attribution to LLMs
183 can be challenging. The discrete nature of text
184 complicates baseline selection and interpolation

185 paths, motivating adaptations such as discretized
186 integrated gradients (Sanyal and Ren, 2021). Fur-
187 ther, perturbation-based methods scale poorly to
188 the long contexts typical of LLMs, can produce out-
189 of-distribution inputs, and are not directly applica-
190 ble to text generation. This motivated extensions
191 of SHAP-style attribution methods for explaining
192 generated responses (Paes et al., 2025). Overall,
193 input attribution can highlight parts of an input that
194 influenced a specific output, but it remains limited
195 to explaining single input–output pairs and does not
196 capture the multi-step decision-making, tool invo-
197 cation, or orchestration processes that characterize
198 tool-using LLMs.

199 **Training data attribution.** Training data attribu-
200 tion explains model behavior by identifying influen-
201 tial examples from the training corpus, for instance
202 via influence functions (Koh and Liang, 2017),
203 which estimate the contribution of individual train-
204 ing documents without retraining (e.g., Guo et al.,
205 2021; Grosse et al., 2023; Choe et al., 2024). In
206 practice, these methods face substantial scalabil-
207 ity limitations for LLMs and are often restricted
208 to smaller models or fine-tuning data. More fun-
209 damentally, training data attribution explains how
210 models acquire knowledge during training, but pro-
211 vides little insight into runtime decision-making,
212 tool selection, and multi-step interactions that char-
213 acterize tool-using LLM systems.

214 **Probing.** Probing aims to determine what kinds
215 of information are encoded in a model’s internal
216 representations by training a simple classi-
217 fier—often linear—to predict a property of interest
218 from those representations (Belinkov, 2022). It has
219 been used to study a wide range of linguistic and
220 semantic properties, from morphology and syntax
221 (Belinkov et al., 2017; Hewitt and Manning, 2019)
222 to more complex phenomena such as world knowl-
223 edge (Li et al., 2023). While probing can reveal
224 what information a model *could* in principle use,
225 it does not necessarily explain what information
226 the model actually causally relies on when making
227 predictions (Elazar et al., 2021). Results are also
228 sensitive to probe complexity, requiring careful
229 controls (Hewitt and Liang, 2019; Voita and Titov,
230 2020). In the context of tool-using LLM systems,
231 probing has seen limited adoption, as it provides
232 little insight into runtime decision-making across
233 tools, though it has been used to analyze whether
234 retrieval-augmented models rely on parametric or
235 contextual knowledge (Tighidet et al., 2024).

236 **Causal intervention.** Causal intervention meth-
237 ods aim to identify which internal components
238 of a neural network are causally responsible for
239 specific behaviors, and are often associated with
240 work on mechanistic interpretability (e.g., Geiger
241 et al., 2021; Hanna et al., 2023). One approach
242 is activation patching (Vig et al., 2020), which
243 replaces internal activations between minimally
244 different inputs to localize components that drive
245 changes in outputs. Related methods include causal
246 tracing (Meng et al., 2022) and other localization
247 techniques for neurons or attention heads (e.g.,
248 Goldowsky-Dill et al., 2023; Wang et al., 2023a).

249 By disentangling causal from correlational ef-
250 fects, these methods offer a principled way to study
251 model internals, but practical challenges arise when
252 applied to LLMs including difficulty in isolating
253 minimal, faithful subnetworks that implement spe-
254 cific functionality (Mueller et al., 2025). More-
255 over, causal intervention techniques have not yet
256 been systematically applied to tool-using systems,
257 where behaviors emerge not only from model inter-
258 nal mechanisms but also from interactions between
259 models, tools, and execution environments.

260 **Self-explanations.** Self-explanation refers to
261 LLMs producing an output together with a natural-
262 language rationale intended to explain why that
263 output was generated (Rajani et al., 2019). While
264 such rationales are often coherent and plausible,
265 they may not faithfully reflect the causal factors
266 that guided generation (Jacovi and Goldberg, 2020;
267 Barez et al., 2025). Accordingly, explanations are
268 commonly evaluated along three dimensions: *Plau-*
269 *sibility* captures whether explanations appear coher-
270 ent and convincing to humans (Lage et al., 2019;
271 Jacovi and Goldberg, 2020). *Faithfulness* asks if
272 the explanation aligns with the causal decision pro-
273 cess of the model, which can be tested through
274 causal interventions or counterfactual manipula-
275 tions (Lyu et al., 2024).¹ *Simulatability* captures
276 if explanations help users anticipate model behav-
277 ior or failure modes (Doshi-Velez and Kim, 2017;
278 Hase and Bansal, 2020). In tool-using LLM sys-
279 tems, self-explanations are particularly limited, as
280 such rationales may omit, misrepresent, or halluci-
281 nate tool invocations, intermediate decisions, and
282 interactions with external components, potentially
283 leading to unwarranted trust in system behavior.

¹We define *faithfulness* as the extent to which an explana-
tion reflects the true causal factors influencing system behavior,
rather than its ability to provide plausible justifications.

284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332

3 Explainability for Tool-using LMs

Explainability for tool-using LLMs spans a spectrum between two complementary perspectives: using tools to generate explanations of system behavior, and explaining why and how tools were invoked during inference. Many systems combine both directions, for example, by calling analysis tools to produce explanations while simultaneously exposing tool calls or reasoning traces. In this section, we cover approaches across this spectrum at increasing levels of abstraction.

3.1 RAG and Explainability

Retrieval-augmented generation enhances LLMs by providing mechanisms to retrieve passages from knowledge bases and to condition generation on the retrieved context. Existing work on explainability in the context of RAG focuses on how the retrieved context shapes model output rather than on the retrieval mechanism and its system-level decisions. Consequently, RAG occupies a boundary position between output-grounding techniques and tool-level explainability.

RAG for information tracing and grounding.

RAG is widely used to ground model outputs in curated sources, supporting fact tracing, citation, and provenance attribution (Sankararaman et al., 2024; Zhu et al., 2025b,a; Xia et al., 2025; Qi et al., 2024; Armant et al., 2024). By pairing answers with retrieved documents, these systems make outputs easier to verify. However, standard RAG setups may provide relevant documents, but do not explain why documents were retrieved or how retrieval influenced the model’s reasoning. This limits their usefulness to knowledge-based tasks where users can verify model outputs against the ground truth in the retrieved documents.

RAG for generating explanations.

RAG has also been used to generate explanations or rationales alongside predictions. Some approaches prompt RAG models to produce self-explanations conditioned on retrieved context (Garigliotti, 2025; Wei et al., 2025; Peng et al., 2025), while others adapt attribution methods to highlight influential parts of the combined prompt and context (Sudhi et al., 2024). A complementary line of work explains outputs in terms of retrieved sources by perturbing document sets or knowledge graphs to identify influential evidence (Rorseth et al., 2024; Balanos et al., 2025). Systems that use RAG to ex-

plain external processes rather than model behavior are out of scope in this survey (e.g., Minor and Kaucher, 2024; Li et al., 2025; Gao et al., 2023b).

Explaining and evaluating RAG models.

Evaluation of RAG systems often draws on properties from interpretability literature, such as faithfulness, to capture the causal relationship between retrieved context and generated outputs. Prior work studies answer faithfulness and attribution groundedness (Gao et al., 2024; Qi et al., 2024; Song et al., 2025; Ye et al., 2024; Patel et al., 2024), and distinguishes between citation correctness and faithfulness (Wal-lat et al., 2025). Similar notions are reflected in benchmarks that assess context utilization, relevance, and answer completeness (Friel et al., 2025). Overall, however, RAG explainability primarily addresses grounding and justification of generated outputs. Most existing approaches explain answers instead of the retrieval and generation decisions, and thus provide only a partial explanation.

3.2 Tool-Level Transparency

Beyond information retrieval mechanisms, tool-using LLMs invoke external APIs or computational tools to perform reasoning and actions. At this level, explainability concerns how individual tool calls contribute to an agent’s inference process, covering both the exposure of tool invocations themselves and explanations for why tools were used. Existing work spans from implicit transparency, to structured tool calls, to explicit, tool-supported explanation generation.

Transparent tool calls without explicit explanations.

Structured tool calls that expose the tool name, input arguments, and outputs already provide a basic form of interpretability by making intermediate steps observable. Early systems such as ReAct (Yao et al., 2023b) interleave natural-language reasoning with external actions, allowing users to inspect which tools were invoked and with what effects. Toolformer (Schick et al., 2023) and HuggingGPT (Shen et al., 2023) generalize this paradigm by enabling LLMs to call diverse external tools through structured APIs. While these approaches do not generate dedicated explanations, their explicit tool-use traces make reasoning steps more transparent than end-to-end text generation.

Related work such as Gorilla (Patil et al., 2024) grounds tool calls in retrieved API documentation and evaluates correctness via syntactic matching,

333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381

further improving verifiability of tool usage. Similarly, standardization efforts like the Model Context Protocol(MCP) aim to improve transparency and interoperability of tool interfaces, facilitating inspection and auditing of individual tool calls.

Explanation generation with and for tool calls.

Other systems explicitly generate explanations for tool usage or employ tools to construct explanations. SCRIBE (Fawzi et al., 2025) produces structured, step-by-step explanations grounded in tool outputs in educational settings, while ECHO (Vanbrabant et al., 2025) equips an LLM with access to model interfaces, data, and XAI utilities to orchestrate multi-step explanations of another model’s behavior. Related approaches such as KnowThyself (Prasai et al., 2025) apply similar orchestration principles for model interpretability. Complementary work focuses on making intermediate reasoning artifacts explicit. Program-Aided Language Models (PAL) (Gao et al., 2023a) and Program-of-Thoughts (Chen et al., 2023) express reasoning as executable code, separating logical inference and computation from generation, enabling inspection of each step. Natural-language reasoning paradigms such as Chain-of-Thought (CoT) prompting (Wei et al., 2022), Tree-of-Thought (Yao et al., 2023a), and plan-based or self-reflective methods (Wang et al., 2023b; Shinn et al., 2023; Madaan et al., 2023) expose intermediate reasoning as structured text, improving transparency and often performance. However, different analyses show that such reasoning traces are frequently *unfaithful*: they describe plausible reasoning without necessarily reflecting the causal processes that produced the output (Turpin et al., 2023; Lanham et al., 2023; Paul et al., 2024). In tool-using settings, this limitation becomes more pronounced, as reasoning traces may plausibly justify tool use without corresponding to the actual tool calls executed or the causal role their outputs played in the final decision. Recent work on interactive reasoning (Pang et al., 2025) highlights this gap by enabling users to inspect, restructure, and intervene in CoT reasoning, suggesting that faithfulness in agentic contexts requires alignment not only between reasoning and output, but also between reasoning traces and the system’s actions.

Overall, tool-level reasoning transparency improves interpretability by exposing or explaining individual tool invocations and intermediate decisions. However, these approaches typically focus

on localized steps rather than entire agent trajectories. We discuss system-level methods that address this shortcoming in the next subsection.

3.3 System-Level Transparency

System-level transparency addresses explainability at the level of complete agent architectures and execution histories, rather than individual reasoning steps or isolated tool calls. In contrast to tool-level transparency, which focuses on local decisions within a single interaction, system-level approaches characterize component composition, information flow across tools, and agent behavior across entire trajectories or populations of runs.

Documentation. Structured documentation provides information on model capabilities, components, and intended use. Model Cards (Mitchell et al., 2019) and Foundation Model Transparency Reports (Bommasani et al., 2024) standardize disclosure of datasets, evaluation, risks, and architecture. While originally designed for monolithic models, these frameworks include some indicators relevant to tool-augmented systems—such as disclosure of model components and their integration—but do not support characterization of tools, tool dependencies, or runtime attribution.

Provenance-based approaches capture how outcomes arise from interactions among models, tools, and data. PROV-AGENT (Souza et al., 2025) extends the W3C PROV model² and integrates with the Model Context Protocol to record prompts, model responses, tool invocations, and their dependencies within unified provenance graphs. This enables reconstruction of agent workflows and attribution across human and machine-driven activities.

Auditing, verification, and observability. Beyond tracing individual executions, several systems support auditing and evaluation of agent behavior at scale. VeriLA (Sung et al., 2025) verifies multi-step agents by scoring intermediate subtasks and aggregating correctness over structured execution graphs. Observability platforms such as LangSmith Observability³ and DeepEval⁴ provide end-to-end traces and component-level evaluations, while visualization tools like Seaview (Bula et al., 2025), the OpenHands Trajectory Visualizer⁵, and Agent Trajectory Explorer (Desmond et al., 2025) render

²[w3.org/TR/prov-overview/](https://www.w3.org/TR/prov-overview/)

³langchain.com/langsmith/observability

⁴deepeval.com/docs/getting-started

⁵github.com/OpenHands/trajectory-visualizer

479 long agent trajectories in human-interpretable form
480 for debugging and comparison. These approaches
481 extend explainability from local reasoning steps to
482 global system behavior, enabling analysis of how
483 complex, tool-using LLM agents operate across
484 time and interactions.

485 **Explainability for computer systems and pro-**
486 **grams.** Many challenges posed by tool-using
487 LLM agents—such as attributing outcomes to in-
488 termediate steps, tools, and data—have long been
489 studied in the context of program and system ex-
490 plainability. Work in software debugging framed
491 explanations as causal accounts of system behavior,
492 exemplified by the *Whyline* system (Ko and My-
493 ers, 2004, 2009), which allowed users to ask “why”
494 and “why not” questions about program outputs,
495 and by *delta debugging* (Zeller and Hildebrandt,
496 2025; Zeller, 2009), which isolates minimal failure-
497 inducing inputs and state differences. Work in
498 databases and workflow systems formalized expla-
499 nation through *data provenance*, tracing how data
500 items and transformations contribute to observed
501 results (Buneman et al., 2001; Cheney et al., 2009;
502 Freire et al., 2008). In parallel, research on *observ-*
503 *ability* in distributed systems developed techniques
504 for reconstructing causal execution paths across
505 components (Fonseca et al., 2007; Sigelman et al.,
506 2010; Kaldor et al., 2017). These approaches es-
507 tablished system-level causality, provenance, and
508 observability as foundational principles for explain-
509 ing complex, multi-component systems, providing
510 important conceptual references for explainability
511 in tool-using LLM agents.

512 3.4 Domain-Specific and Multimodal Systems

513 For some LLM agents, explainability is driven by
514 how non-textual modalities and domain signals en-
515 ter the decision loop. Explanations need to relate
516 actions and tool calls to perceptual inputs (often vi-
517 sion), structured intermediates, or domain concepts.
518 A key distinction from an explainability perspec-
519 tive is whether domain-specificity and modalities
520 are handled by separate tools or integrated into a
521 multimodal backbone:

522 **Vision as a separate module or tool.** Many
523 agents route images through a dedicated percep-
524 tion component (e.g., detector, encoder, VQA),
525 yielding intermediate outputs for downstream rea-
526 soning. Explanations focus on grounding deci-
527 sions in visual evidence and exposing the percep-
528 tion–reasoning pipeline, for example, by linking

529 answers to image regions or attributes (Yu and Ana-
530 niadou, 2025; Park et al., 2018). Recent work im-
531 proves transparency with intermediate structures,
532 such as task graphs (Nooralahzadeh et al., 2024),
533 interpretable vision-language representations (Sha-
534 ham et al., 2024), or flowchart-like attributions for
535 multi-step visual reasoning (Suri et al., 2025).

536 **Integrated multimodal backbones.** Other sys-
537 tems utilize joint vision-language models. Here,
538 explainability shifts from tool invocation to inter-
539 nal mechanisms of multimodal integration, such as
540 cross-modal attention, circuits, or concept repre-
541 sentations (Palit et al., 2023; Huo et al., 2024; Fang
542 et al., 2024; Parekh et al., 2024; Bhalla et al., 2024;
543 Gandelsman et al., 2025). These works do not
544 address tool use directly, but they serve to demon-
545 strate how explanation targets and faithfulness cri-
546 teria shift when augmentation relies on internal
547 representations instead of explicit tool calls.

548 **Domain-specific settings.** Explainability is also
549 required in specialized domains, where trans-
550 parency depends on linking behavior to do-
551 main knowledge and workflows. For example,
552 in LLM-controlled robotics, interpretability re-
553 quires grounding in a library of executable robot
554 skills (Ichter et al., 2022). In medicine, clinical
555 agents emphasize interpretable reasoning traces
556 and tool use to enable expert oversight (Shimgekar
557 et al., 2025; Shi et al., 2025).

558 4 Discussion

559 4.1 Does tool use improve interpretability?

560 **Tool use makes LLMs more interpretable.**
561 Tool output is not the entire story of an LLM’s
562 behavior, but it constitutes an important and often
563 informative part of it, as knowing which tool pro-
564 duced an intermediate result makes substantial por-
565 tions of the reasoning process explicit. Because
566 many tools operate according to comparatively
567 transparent logic—such as arithmetic computation,
568 database lookup, or API execution—they typically
569 require less additional interpretation than purely
570 statistical text generation. Moreover, intermediate
571 tool outputs provide concrete artifacts that users
572 can inspect, verify, and cross-check against their
573 own knowledge or external sources, supporting
574 stepwise validation of system behavior rather than
575 relying solely on final answers. Empirical stud-
576 ies further suggest that mechanisms such as CoT
577 prompting, RAG, and reference attribution can in-

crease user trust and effectiveness when interacting with chatbots (Wei et al., 2022; Menick et al., 2022; Zamfirescu-Pereira et al., 2023), although this trust is not always warranted and depends on the faithfulness of the exposed signals (Jacovi and Goldberg, 2020; Turpin et al., 2023; He et al., 2024).

Providing explanations is harder with tool use.

At the same time, tool use introduces additional complexity that must itself be understood and explained. Whereas traditional language-only systems require explanations primarily for text generation, tool-calling agents require explanations not only for generated outputs but also for tool selection, invocation, and integration. Because tool calls are often easier to interpret than free-form generation, there is a risk that their visibility leads to an overemphasis on their role, creating unwarranted increases in user trust even in tasks where tools contribute little or no value. Given that explaining text generation remains an open and challenging problem, we caution against claims that tool use inherently makes LLMs more interpretable or trustworthy. Similar concerns have been raised in the context of RAG systems, where improved interpretability or groundedness is frequently cited as a motivation (e.g., He et al., 2024; Wang et al., 2024, 2025b), yet such claims are rarely validated outside dedicated and costly evaluation studies.

4.2 Open Challenges

Granularity. A central open challenge for explainability in tool-calling LLMs concerns the appropriate granularity and locality of explanations. Classical XAI methods were developed for relatively simple input-output mappings, where local explanations—such as attributing a sentiment classification to specific words—provides meaningful insight. In tool-calling systems, however, individual input tokens often carry limited explanatory value compared to factors external to the input, including retrieved context, tool outputs, execution results and orchestration choices. As a result, explanations that remain localized to surface inputs risk being misleading rather than merely incomplete, faithfully describing token-level correlations while omitting the dominant causal influences on the system’s behavior. At the same time, shifting explanations toward system-level factors raises new challenges: how far back in an agent’s execution history an explanation should extend, how to aggregate explanations across multiple steps and

components, and how to balance fine-grained faithfulness against usability and cognitive load. Granularity is further intertwined with controllability, as overly coarse explanations do not support meaningful intervention, while overly fine-grained traces may expose levers that users cannot safely or effectively manipulate. Defining what constitutes a “local” explanation for an agentic system and how explanation granularity should adapt to task, user role, and risk remains an open problem.

System-level characterizations. A central open challenge for explainability in tool-calling, agent-based systems concerns system-level characterizations, corresponding to global explanations of how a system works overall. Research on user mental models and misconceptions indicates substantial uncertainty about the high-level structure of complex AI systems, including which components they contain (e.g., session memory, internet access, retrieval modules, calculators, symbolic computation, background reasoning, or filtering mechanisms) and how these components interact. At the same time, system providers are often weakly incentivized to disclose such details, and deployed systems are fluid: components, capabilities, and orchestration strategies change frequently, rendering static descriptions quickly outdated.

Given the impact that architectural choices can have on system behavior and performance, accurately characterizing a system at a high level may contribute more to user understanding than increasingly fine-grained attribution methods applied to individual inputs. Local explanations—such as input attributions, calculation traces, or data attribution—can be difficult to interpret in isolation and often presuppose a correct understanding of the system’s overall structure. Without this “big picture,” such explanations risk being misunderstood or overgeneralized, whereas an accurate system-level characterization can already address many observed behaviors, limitations, and failure modes.

A further challenge lies in determining how such global characterizations should be constructed. Prior work on stakeholder needs for AI explanations provides a useful starting point, identifying recurring global questions such as what logic a system follows, how reliable its outputs are, and what data or resources it depends on (Liao et al., 2020). For tool-calling LLMs, this suggests shifting emphasis from exhaustive technical detail toward clarifying those aspects that users most fre-

quently misunderstand. One promising direction is the development of structured, standardized templates—analogue to model cards or foundation model transparency reports—that foreground tool use, external dependencies, and orchestration logic as first-class elements of system documentation. More broadly, these considerations point to the need for explanations at multiple, clearly distinguished levels: general knowledge about how AI systems work, high-level characterizations of specific deployed systems, more detailed global descriptions of salient features and behaviors, and local explanations addressing individual outputs.

Evaluation. Evaluation of explainability in tool-augmented systems remains underdeveloped, with existing approaches often narrow in scope, difficult to standardize, and challenging to compare across systems and tasks (Anjomshoae et al., 2019; Sado et al., 2023). A central challenge is that agent behavior is deeply entangled with human interaction: different users, prompts, and dialogue trajectories can steer an agent toward different tool uses, intermediate steps, and self-explanations, complicating controlled comparisons and reproducibility. Additionally, explanations in agentic systems concern long horizons and multiple components, so their quality cannot be assessed solely at the level of final outputs or individual steps. Provenance-oriented approaches suggest that richer and structured logging of actions, tool invocations, and information sources can improve auditability and support more systematic evaluation in such interactive settings (Kale et al., 2022), but how to evaluate the coherence and faithfulness of explanations across entire agent workflows remains an open problem.

Growing capabilities. The capabilities of AI agents are developing at a rapid pace. Recent evidence suggests that the length of software engineering tasks these agents can autonomously complete is doubling roughly every seven months (Kwa et al., 2025)—a trend which, if sustained, would imply that agents will soon be capable of executing workflows that currently require days or weeks of human effort. This growth reflects several interacting developments, including increasing model and data scale, expanding access to external tools and environments, and longer task horizons involving many sequential and interdependent decisions. While each of these trends raises new challenges for explainability, long-horizon task execution presents a particularly acute problem. In such settings, fail-

ures often emerge from the accumulation of many small decisions—tool selections, parameter choices, intermediate assumptions—rather than from a single identifiable error. As agents become more autonomous and deliver increasingly complete, ready-to-use outputs, opportunities for human intervention during execution diminish, shifting the burden toward post-hoc verification and auditing. Explainability methods must therefore operate at appropriate levels of abstraction that allow users to reconstruct outcomes, detect compounding errors, and assign responsibility across extended workflows. Whether existing XAI approaches can scale to these demands, or whether understanding of increasingly capable agentic systems will remain largely superficial, remains an open question.

5 Conclusion

In this survey, we argued that the ongoing shift of LLMs from standalone text generators to tool-using agents necessitates a rethinking of explainable AI. In such systems, understanding behavior requires more than explaining why a particular token was generated: it requires faithful accounts of how models delegate subtasks, orchestrate tool use, and integrate execution results over extended interactions. Our review of traditional explainability methods showed that most existing approaches were developed for monolithic models and single-prediction settings, and fall short of capturing multi-step, system-level processes that characterize tool-augmented LLMs. We surveyed work that addresses these gaps, including approaches tailored to specific forms of tool use such as retrieval-augmented generation and multimodal pipelines, methods that expose or explain individual tool calls, and system-level techniques that trace and audit provenance across complete agent trajectories.

Despite the potential interpretability benefits from explicit tool use, significant challenges remain. These include determining appropriate levels of explanatory granularity, developing robust and comparable evaluation practices, and aligning explanations with users’ mental models of complex AI systems. Addressing these challenges calls for explainability approaches that operate at the system level, supporting faithful inspection of agent behavior across tools, data, and execution steps. Developing such approaches is a prerequisite for enabling meaningful oversight and appropriate trust in tool-using AI systems.

Limitations

This survey focuses on explainability for tool-calling large language models and agentic AI systems, and does not aim to provide a comprehensive overview of explainability research across all AI domains. As a result, related areas such as explainability for recommender systems, autonomous vehicles, control, or other not language-based systems are largely out of scope, except where they intersect with tool-augmented language model architectures.

In addition, our discussion emphasizes conceptual distinctions, design patterns, and representative approaches rather than exhaustive technical comparisons of individual methods. While we highlight strengths and limitations at a high level, the suitability of specific explainability techniques may depend on application context, system design, and user requirements in ways that are more nuanced than can be fully captured within the scope of this survey.

References

Sule Anjomshoae, Amro Najjar, Davide Calvaresi, and Kary Främling. 2019. [Explainable agents and robots: Results from a systematic literature review](#). In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS '19, Montreal, QC, Canada, May 13-17, 2019*, pages 1078–1088. International Foundation for Autonomous Agents and Multiagent Systems.

Vincent Armant, Amira Mouakher, Felipe Vargas-Rojas, Danaï Symeonidou, Joris Guérin, Isabelle Mougenot, and Jean-Christophe Desconnets. 2024. [Can Knowledge Graphs and Retrieval-Augmented Generation be combined to Explain Query/Answer Relationships Truthfully?](#) In *DAO-XAI 2024 Data meets Ontologies in Explainable AI co-located with the 27th European Conference on Artificial Intelligence (ECAI 2024)*, volume 3833, Santiago de Compostela, Spain.

Alejandro Barredo Arrieta, Natalia Díaz Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. 2020. [Explainable artificial intelligence \(XAI\): concepts, taxonomies, opportunities and challenges toward responsible AI](#). *Inf. Fusion*, 58:82–115.

Georgios Balanos, Evangelos Chasanis, Konstantinos Skianis, and Evaggelia Pitoura. 2025. [KGRAG-Ex: Explainable Retrieval-Augmented Generation with Knowledge Graph-based Perturbations](#). *arXiv preprint*. ArXiv:2507.08443 [cs].

Fazl Barez, Tung-Yu Wu, Iván Arcuschin, Michael Lan, Vincent Wang, Noah Siegel, Nicolas Collignon,

Clement Neo, Isabelle Lee, Alasdair Paren, Adel Bibi, Robert Trager, Damiano Fornasiere, John Yan, Yanai Elazar, and Yoshua Bengio. 2025. [Chain-of-thought is not explainability](#). Working paper, Oxford Martin AI Governance Initiative (AIGI), University of Oxford. Accessed 2025-12-16.

Yonatan Belinkov. 2022. [Probing classifiers: Promises, shortcomings, and advances](#). *Computational Linguistics*, 48(1):207–219.

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. [What do neural machine translation models learn about morphology?](#) In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.

Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flávio P. Calmon, and Himabindu Lakkaraju. 2024. [Interpreting CLIP with sparse linear concept embeddings \(splice\)](#). In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.

Rishi Bommasani, Kevin Klyman, Shayne Longpre, Betty Xiong, Sayash Kapoor, Nestor Maslej, Arvind Narayanan, and Percy Liang. 2024. [Foundation model transparency reports](#). In *Proceedings of the Seventh AAI/ACM Conference on AI, Ethics, and Society (AIES-24) - Full Archival Papers, October 21-23, 2024, San Jose, California, USA - Volume 1*, pages 181–195. AAAI Press.

Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z. Gajos. 2021. [To trust or to think: Cognitive forcing functions can reduce overreliance on AI in ai-assisted decision-making](#). *Proc. ACM Hum. Comput. Interact.*, 5(CSCW1):188:1–188:21.

Timothy Bula, Saurabh Pujar, Luca Buratti, Mihaela A. Bornea, and Avirup Sil. 2025. [Seaview: Software engineering agent visual interface for enhanced workflow](#). *CoRR*, abs/2504.08696.

Peter Buneman, Sanjeev Khanna, and Wang Chiew Tan. 2001. [Why and where: A characterization of data provenance](#). In *Database Theory - ICDT 2001, 8th International Conference, London, UK, January 4-6, 2001, Proceedings*, volume 1973 of *Lecture Notes in Computer Science*, pages 316–330. Springer.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Trans. Mach. Learn. Res.*, 2023.

James Cheney, Laura Chiticariu, and Wang Chiew Tan. 2009. [Provenance in databases: Why, how, and where](#). *Found. Trends Databases*, 1(4):379–474.

Sang Keun Choe, Hwijeen Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya

889	Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, Jeff Schneider, Eduard Hovy, Roger Grosse, and Eric Xing. 2024. What is Your Data Worth to GPT? LLM-Scale Data Valuation with Influence Functions . <i>arXiv preprint</i> . ArXiv:2405.13954 [cs].	945
890		946
891		947
892		948
893		949
894		950
895	Michael Desmond, Ja Young Lee, Ibrahim Ibrahim, James M. Johnson, Avirup Sil, Justin MacNair, and Ruchir Puri. 2025. Agent trajectory explorer: Visualizing and providing feedback on agent trajectories . In <i>AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA</i> , pages 29634–29636. AAAI Press.	951
896		952
897		953
898		
899		954
900		955
901		956
902		957
903	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 4171–4186. Association for Computational Linguistics.	958
904		959
905		960
906		961
907		962
908		
909		963
910		964
911		965
912		966
913	Finale Doshi-Velez and Been Kim. 2017. Towards A Rigorous Science of Interpretable Machine Learning . <i>arXiv preprint</i> . ArXiv:1702.08608 [cs, stat].	967
914		968
915		969
916	Yu Du, Fangyun Wei, and Hongyang Zhang. 2024. Anytool: Self-reflective, hierarchical agents for large-scale API calls . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	970
917		971
918		972
919		973
920		
921	Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2021. Amnesic probing: Behavioral explanation with amnesic counterfactuals . <i>Transactions of the Association for Computational Linguistics</i> , 9:160–175.	974
922		975
923		976
924		977
925		978
926	Junfeng Fang, Zac Bi, Ruipeng Wang, Houcheng Jiang, Yuan Gao, Kun Wang, An Zhang, Jie Shi, Xiang Wang, and Tat-Seng Chua. 2024. Towards neuron attributions in multi-modal large language models . In <i>Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024</i> .	979
927		980
928		981
929		982
930		983
931		984
932		985
933		986
934	Fares Fawzi, Vinitra Swamy, Dominik Glandorf, Tanya Nazaretsky, and Tanja Käser. 2025. SCRIBE: Structured chain reasoning for interactive behaviour explanations using tool calling . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 29273–29298, Suzhou, China. Association for Computational Linguistics.	987
935		988
936		989
937		990
938		991
939		992
940		993
941	Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R. Costa-jussà. 2024. A primer on the inner workings of transformer-based language models . <i>CoRR</i> , abs/2405.00208.	994
942		995
943		996
944		997
	Rodrigo Fonseca, George Porter, Randy H. Katz, Scott Shenker, and Ion Stoica. 2007. X-trace: A pervasive network tracing framework . In <i>4th Symposium on Networked Systems Design and Implementation (NSDI 2007), April 11-13, 2007, Cambridge, Massachusetts, USA, Proceedings</i> . USENIX.	998
		999
	Juliana Freire, David Koop, Emanuele Santos, and Cláudio T. Silva. 2008. Provenance for computational tasks: A survey . <i>Comput. Sci. Eng.</i> , 10(3):11–21.	
	Robert Friel, Masha Belyi, and Atindriyo Sanyal. 2025. RAGBench: Explainable Benchmark for Retrieval-Augmented Generation Systems . <i>arXiv preprint</i> . ArXiv:2407.11005 [cs].	
	Yossi Gandelsman, Alexei A. Efros, and Jacob Steinhardt. 2025. Interpreting the second-order effects of neurons in CLIP . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	
	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023a. PAL: program-aided language models . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 10764–10799. PMLR.	
	Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. 2023b. Chat-rec: Towards interactive and explainable llms-augmented recommender system . <i>CoRR</i> , arXiv:2303.14524.	
	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-Augmented Generation for Large Language Models: A Survey . <i>arXiv preprint</i> . ArXiv:2312.10997 [cs].	
	Darío Garigliotti. 2025. Self-explanatory Retrieval-Augmented Generation for SDG Evidence Identification . In <i>Advances in Conceptual Modeling</i> , pages 124–132, Cham. Springer Nature Switzerland.	
	Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. 2021. Causal abstractions of neural networks . In <i>Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual</i> , pages 9574–9586.	
	Iona Gessinger, Katie Seaborn, Madeleine Steeds, and Benjamin R. Cowan. 2025. Chatgpt and me: First-time and experienced users’ perceptions of chatgpt’s communicative ability as a dialogue partner . <i>Int. J. Hum. Comput. Stud.</i> , 194:103400.	
	Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. 2023. Localizing model behavior with path patching . <i>CoRR</i> , arXiv:2304.05969.	
	Roger B. Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, Evan Hubinger,	

1000	Kamile Lukosiute, Karina Nguyen, Nicholas Joseph, Sam McCandlish, Jared Kaplan, and Samuel R. Bowman. 2023. Studying large language model generalization with influence functions . <i>CoRR</i> , arXiv:2308.03296.	1056
1001		1057
1002		1058
1003		1059
1004		1060
1005	Han Guo, Nazneen Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. 2021. FastIF: Scalable Influence Functions for Efficient Model Interpretation and Debugging . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10333–10350, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1061
1006		1062
1007		1063
1008		1064
1009		1065
1010		1066
1011		1067
1012		1068
1013	Zayd Hammoudeh and Daniel Lowd. 2024. Training data influence analysis and estimation: a survey . <i>Mach. Learn.</i> , 113(5):2351–2403.	1069
1014		1070
1015		1071
1016	Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1072
1017		1073
1018		1074
1019		1075
1020		1076
1021	Peter Hase and Mohit Bansal. 2020. Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5540–5552, Online. Association for Computational Linguistics.	1077
1022		1078
1023		1079
1024		1080
1025		1081
1026		1082
1027	Gaole He, Gianluca Demartini, and Ujwal Gadiraju. 2025. Plan-then-execute: An empirical study of user trust and team performance when using LLM agents as A daily assistant . In <i>Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, CHI 2025, Yokohama, Japan, 26 April 2025- 1 May 2025</i> , pages 414:1–414:22. ACM.	1083
1028		1084
1029		1085
1030		1086
1031		1087
1032		1088
1033		1089
1034	Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-Retriever: Retrieval-Augmented Generation for Textual Graph Understanding and Question Answering . <i>Advances in Neural Information Processing Systems</i> , 37:132876–132907.	1090
1035		1091
1036		1092
1037		1093
1038		1094
1039		1095
1040		1096
1041	John Hewitt and Percy Liang. 2019. Designing and interpreting probes with control tasks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2733–2743, Hong Kong, China. Association for Computational Linguistics.	1097
1042		1098
1043		1099
1044		1100
1045		1101
1046		1102
1047		1103
1048	John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.	1104
1049		1105
1050		1106
1051		1107
1052		1108
1053		1109
1054		1110
1055		1111
		1112
		1113
	Jiahao Huo, Yibo Yan, Boren Hu, Yutao Yue, and Xuming Hu. 2024. Mmneuron: Discovering neuron-level domain-specific interpretation in multimodal large language model . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 6801–6816. Association for Computational Linguistics.	1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467

1114		on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024. OpenReview.net.	
1115			
1116	Sunnie S. Y. Kim, Jennifer Wortman Vaughan, Q. Vera Liao, Tania Lombrozo, and Olga Russakovsky. 2025. Fostering appropriate reliance on large language models: The role of explanations, sources, and inconsistencies. In <i>Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems, CHI 2025, Yokohama Japan, 26 April 2025- 1 May 2025</i> , pages 420:1–420:19. ACM.	Retrieval-augmented generation for knowledge-intensive NLP tasks. In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	1172 1173 1174 1175 1176
1117			
1118			
1119			
1120			
1121			
1122			
1123			
1124	Amy J. Ko and Brad A. Myers. 2004. Designing the whyline: a debugging interface for asking questions about program behavior. In <i>Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '04</i> , page 151–158, New York, NY, USA. Association for Computing Machinery.	Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Emergent world representations: Exploring a sequence model trained on a synthetic task. In <i>The Eleventh International Conference on Learning Representations</i> .	1177 1178 1179 1180 1181 1182
1125			
1126			
1127			
1128			
1129			
1130	Amy J. Ko and Brad A. Myers. 2009. Finding causes of program output with the java whyline. In <i>Proceedings of the 27th International Conference on Human Factors in Computing Systems, CHI 2009, Boston, MA, USA, April 4-9, 2009</i> , pages 1569–1578. ACM.	Yuhan Li, Xinni Zhang, Linhao Luo, Heng Chang, Yuxiang Ren, Irwin King, and Jia Li. 2025. G-Refer: Graph Retrieval-Augmented Large Language Model for Explainable Recommendation. In <i>Proceedings of the ACM on Web Conference 2025, WWW '25</i> , pages 240–251, New York, NY, USA. Association for Computing Machinery.	1183 1184 1185 1186 1187 1188 1189
1131			
1132			
1133			
1134			
1135	Pang Wei Koh and Percy Liang. 2017. Understanding Black-box Predictions via Influence Functions. In <i>Proceedings of the 34th International Conference on Machine Learning</i> , pages 1885–1894. PMLR. ISSN: 2640-3498.	Q. Vera Liao, Daniel M. Gruen, and Sarah Miller. 2020. Questioning the AI: informing design practices for explainable AI user experiences. In <i>CHI '20: CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, April 25-30, 2020</i> , pages 1–15. ACM.	1190 1191 1192 1193 1194 1195
1136			
1137			
1138			
1139			
1140	Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2022. Internet-augmented dialogue generation. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 8460–8478. Association for Computational Linguistics.	Scott M. Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 4765–4774.	1196 1197 1198 1199 1200 1201
1141			
1142			
1143			
1144			
1145			
1146			
1147	Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, and 6 others. 2025. Measuring AI ability to complete long tasks. <i>CoRR</i> , abs/2503.14499.	Haoyan Luo and Lucia Specia. 2024. From understanding to utilization: A survey on explainability for large language models. <i>CoRR</i> , abs/2401.12874.	1202 1203 1204
1148			
1149			
1150			
1151			
1152			
1153			
1154			
1155	Isaac Lage, Emily Chen, Jeffrey He, Menaka Narayanan, Been Kim, Sam Gershman, and Finale Doshi-Velez. 2019. An evaluation of the human-interpretability of explanation. <i>CoRR</i> , abs/1902.00006.	Siwen Luo, Hamish Ivison, Soyeon Caren Han, and Josiah Poon. 2024. Local interpretations for explainable natural language processing: A survey. <i>ACM Comput. Surv.</i> , 56(9):232:1–232:36.	1205 1206 1207 1208
1156			
1157			
1158			
1159	Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Karina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Sam McCandlish, Sandipan Kundu, and 11 others. 2023. Measuring faithfulness in chain-of-thought reasoning. <i>CoRR</i> , abs/2307.13702.	Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. 2024. Towards faithful model explanation in NLP: A survey. <i>Computational Linguistics</i> , 50(2):657–723.	1209 1210 1211 1212
1160			
1161			
1162			
1163			
1164			
1165			
1166			
1167			
1168	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020.	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1213 1214 1215 1216 1217 1218 1219 1220 1221 1222 1223
1169			
1170			
1171			
		Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 17359–17372.	1224 1225 1226 1227

1340	reasoning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 15012–15032, Miami, Florida, USA. Association for Computational Linguistics.	1398
1341		1399
1342		1400
1343		1401
1344	Xiangyu Peng, Prafulla Kumar Choubey, Caiming Xiong, and Chien-Sheng Wu. 2025. Unanswerability Evaluation for Retrieval Augmented Generation . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8452–8472, Vienna, Austria. Association for Computational Linguistics.	1402
1345		1403
1346		1404
1347		1405
1348		1406
1349		1407
1350		1408
1351	Suraj Prasai, Mengnan Du, Ying Zhang, and Fan Yang. 2025. Knowthyself: An agentic assistant for llm interpretability . <i>Preprint</i> , arXiv:2511.03878.	1409
1352		1410
1353		1411
1354	Jirui Qi, Gabriele Sarti, Raquel Fernández, and Arianna Bisazza. 2024. Model Internals-based Answer Attribution for Trustworthy Retrieval-Augmented Generation . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 6037–6053, Miami, Florida, USA. Association for Computational Linguistics.	1412
1355		1413
1356		1414
1357		1415
1358		1416
1359		1417
1360		1418
1361	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	1419
1362		1420
1363		1421
1364		1422
1365		1423
1366		1424
1367		1425
1368	Marco Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “why should I trust you?”: Explaining the predictions of any classifier . In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations</i> , pages 97–101. Association for Computational Linguistics.	1426
1369		1427
1370		1428
1371		1429
1372		1430
1373		1431
1374		1432
1375	Joel Rorseth, Parke Godfrey, Lukasz Golab, Divesh Srivastava, and Jaroslaw Szlichta. 2024. RAGE Against the Machine: Retrieval-Augmented LLM Explanations . In <i>2024 IEEE 40th International Conference on Data Engineering (ICDE)</i> , pages 5469–5472. ISSN: 2375-026X.	1433
1376		1434
1377		1435
1378		1436
1379		1437
1380		1438
1381	Fatai Sado, Chu Kiong Loo, Wei Shiung Liew, Matthias Kerzel, and Stefan Wermter. 2023. Explainable goal-driven agents and robots - A comprehensive review .	1439
1382		1440
1383		1441
1384	Hithesh Sankararaman, Mohammed Nasheed Yasin, Tanner Sorensen, Alessandro Di Bari, and Andreas Stolcke. 2024. Provenance: A Light-weight Fact-checker for Retrieval Augmented LLM Generation Output . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track</i> , pages 1305–1313, Miami, Florida, US. Association for Computational Linguistics.	1442
1385		1443
1386		1444
1387		1445
1388		1446
1389		1447
1390		1448
1391		1449
1392	Soumya Sanyal and Xiang Ren. 2021. Discretized integrated gradients for explaining language models . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10285–10299, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1450
1393		1451
1394		1452
1395		1453
1396		1454
1397		1455
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
	Tamar Rott Shaham, Sarah Schwettmann, Franklin Wang, Achyuta Rajaram, Evan Hernandez, Jacob Andreas, and Antonio Torralba. 2024. A multimodal automated interpretability agent . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	1406
		1407
		1408
		1409
		1410
		1411
		1412
	Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugging-gpt: Solving AI tasks with chatgpt and its friends in hugging face . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
	Danli Shi, Xiaolan Chen, Bingjie Yan, Weiyi Zhang, Pusheng Xu, Jianchen Yang, Ruoyu Chen, Siyu Huang, Bowen Liu, Xinyuan Wu, Meng Xie, Ziyu Gao, Yue Wu, Senlin Lin, Kai Jin, Xia Gong, Yih-Chung Tham, Xiujuan Zhang, Li Dong, and 8 others. 2025. A multimodal ai agent for clinical decision support in ophthalmology .	1421
		1422
		1423
		1424
		1425
		1426
		1427
	Soorya Ram Shingekar, Shayan Vassef, Abhay Goyal, Navin Kumar, and Koustuv Saha. 2025. Agentic ai framework for end-to-end medical data inference . <i>Preprint</i> , arXiv:2507.18115.	1428
		1429
		1430
		1431
	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflection: language agents with verbal reinforcement learning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1432
		1433
		1434
		1435
		1436
		1437
		1438
	Benjamin H Sigelman, Luiz André Barroso, Mike Burrows, Pat Stephenson, Manoj Plakal, Donald Beaver, Saul Jaspán, and Chandan Shanbhag. 2010. Dapper, a large-scale distributed systems tracing infrastructure .	1439
		1440
		1441
		1442
		1443
	Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. 2024. Rethinking interpretability in the era of large language models . <i>CoRR</i> , abs/2402.01761.	1444
		1445
		1446
		1447
	Maojia Song, Shang Hong Sim, Rishabh Bhardwaj, Hai Leong Chieu, Navonil Majumder, and Soujanya Poria. 2025. Measuring and enhancing trustworthiness of llms in RAG through grounded attributions and learning to refuse . In <i>The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025</i> . OpenReview.net.	1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455

1456	Renan Souza, Amal Gueroudji, Stephen DeWitt,	in chain-of-thought prompting . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1513
1457	Daniel Rosendo, Tirthankar Ghosal, Robert B. Ross,		1514
1458	Prasanna Balaprakash, and Rafael Ferreira da Silva.		1515
1459	2025. PROV-AGENT: unified provenance for tracking AI agent interactions in agentic workflows . In <i>IEEE International Conference on eScience, eScience 2025, Chicago, IL, USA, September 15-18, 2025</i> , pages 467–473. IEEE.		1516
1460			1517
1461	Viju Sudhi, Sinchana Ramakanth Bhat, Max Rudat, and Roman Teucher. 2024. RAG-Ex: A Generic Framework for Explaining Retrieval Augmented Generation . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24</i> , pages 2776–2780, New York, NY, USA. Association for Computing Machinery.		1518
1462			1519
1463			1520
1464			1521
1465			1522
1466			1523
1467			1524
1468			1525
1469			1526
1470			1527
1471			1528
1472	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks . In <i>Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017</i> , volume 70 of <i>Proceedings of Machine Learning Research</i> , pages 3319–3328. PMLR.		1529
1473			1530
1474			1531
1475			1532
1476			1533
1477			1534
1478			1535
1479	Yoo Yeon Sung, Hannah Kim, and Dan Zhang. 2025. Verila: A human-centered evaluation framework for interpretable verification of LLM agent failures . <i>CoRR</i> , abs/2503.12651.		1536
1480			1537
1481			1538
1482	Manan Suri, Puneet Mathur, Nedim Lipka, Franck Dernoncourt, Ryan A. Rossi, Vivek Gupta, and Dinesh Manocha. 2025. Follow the flow: Fine-grained flowchart attribution with neurosymbolic agents . <i>CoRR</i> , abs/2506.01344.		1539
1483			1540
1484			1541
1485			1542
1486			1543
1487	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, and 38 others. 2022. Lamda: Language models for dialog applications . <i>CoRR</i> , abs/2201.08239.		1544
1488			1545
1489			1546
1490			1547
1491			1548
1492			1549
1493			1550
1494			1551
1495	Zineddine Tighidet, Jiali Mei, Benjamin Piwowarski, and Patrick Gallinari. 2024. Probing language models on their knowledge source . In <i>Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP</i> , pages 604–614, Miami, Florida, US. Association for Computational Linguistics.		1552
1496			1553
1497			1554
1498			1555
1499			1556
1500			1557
1501			1558
1502	Jan Tolsdorf, Alan F. Luo, Monica Kodwani, Junho Eum, Mahmood Sharif, Michelle L. Mazurek, and Adam J. Aviv. 2025. Safety perceptions of generative AI conversational agents: Uncovering perceptual differences in trust, risk, and fairness . In <i>Twenty-First Symposium on Usable Privacy and Security, SOUPS 2025, Seattle, WA, USA, August 10-12, 2025</i> , pages 93–112. USENIX Association.		1559
1503			1560
1504			1561
1505			1562
1506			1563
1507			1564
1508			1565
1509			1566
1510	Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. 2023. Language models don't always say what they think: Unfaithful explanations		1567
1511			1568
1512			1569
			1570
			1571
			1572
			1573
			1574
			1575
			1576
			1577
			1578
			1579
			1580
			1581
			1582
			1583
			1584
			1585
			1586
			1587
			1588
			1589
			1590
			1591
			1592
			1593
			1594
			1595
			1596
			1597
			1598
			1599
			1600
			1601
			1602
			1603
			1604
			1605
			1606
			1607
			1608
			1609
			1610
			1611
			1612
			1613
			1614
			1615
			1616
			1617
			1618
			1619
			1620
			1621
			1622
			1623
			1624
			1625
			1626
			1627
			1628
			1629
			1630
			1631
			1632
			1633
			1634
			1635
			1636
			1637
			1638
			1639
			1640
			1641
			1642
			1643
			1644
			1645
			1646
			1647
			1648
			1649
			1650
			1651
			1652
			1653
			1654
			1655
			1656
			1657
			1658
			1659
			1660
			1661
			1662
			1663
			1664
			1665
			1666
			1667
			1668
			1669
			1670
			1671
			1672
			1673
			1674
			1675
			1676
			1677
			1678
			1679
			1680
			1681
			1682
			1683
			1684
			1685
			1686
			1687
			1688
			1689
			1690
			1691
			1692
			1693
			1694
			1695
			1696
			1697
			1698
			1699
			1700

1682
1683
1684
1685
1686
1687

Evgeny Kharlamov, and Steffen Staab. 2025b. [ArgRAG: Explainable Retrieval Augmented Generation using Quantitative Bipolar Argumentation](#). In *Proceedings of The 19th International Conference on Neurosymbolic Learning and Reasoning*, pages 697–718. PMLR. ISSN: 2640-3498.