

Reasoning, Replicability, and Benchmarking in Large Language and Foundation Models: Methodologies, Challenges, and Pathways Toward Trustworthy, Interpretable, and Inclusive AI

Abstract

This survey provides a comprehensive synthesis of recent advances, methodologies, and enduring challenges in the development, evaluation, and responsible deployment of large language models (LLMs) and foundation models. Motivated by the transformative impact of LLMs across natural language processing, scientific discovery, and real-world applications, the paper critically examines the evolution from symbolic and neural paradigms through contemporary transformer-driven and neurosymbolic architectures, highlighting emergent reasoning capabilities and the drive towards human-like abstraction. The review systematically analyzes benchmarking ecosystems, probing frameworks, and evaluation metrics, emphasizing the limitations of prevailing practices in capturing semantic faithfulness, compositionality, and real-world reasoning, particularly on multistep, cross-modal, and domain-specific tasks. Key contributions include a structured taxonomy of model architectures and fusion strategies, an assessment of hybrid approaches integrating neural, symbolic, and graph-based reasoning, and comparative analyses of benchmark methodologies across linguistic, reasoning, and multimodal domains.

The survey underscores persistent gaps in robustness, interpretability, fairness, and reproducibility—drawing attention to vulnerabilities in adversarial and out-of-distribution scenarios, challenges in auditability and demographic inclusion, and the ongoing reproducibility crisis stemming from inadequate reporting and opaque “language-models-as-a-service” paradigms. It highlights advances in adaptive prompting, modular workflow orchestration, and explainability, while advocating for open science, FAIR data practices, and transparent, community-driven benchmarking. Strategic recommendations target holistic evaluation protocols, enhanced benchmarking diversity, rigorous auditing, responsible design, and the institutionalization of modular, reproducible workflows. The paper concludes that future progress in LLM research and deployment is contingent upon sustaining openness, modularity, explainability, reproducibility, and ethical responsibility, thus ensuring trustworthiness, equitable, and societally beneficial language technologies.

ACM Reference Format:

. 2025. Reasoning, Replicability, and Benchmarking in Large Language and Foundation Models: Methodologies, Challenges, and Pathways Toward Trustworthy, Interpretable, and Inclusive AI. In . ACM, New York, NY, USA, 30 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Advancements in AI systems hinge on the rapid progress of reasoning capabilities, benchmarking practices, and a critical appraisal of model architectures. This survey offers a comprehensive synthesis of literature focusing on the current landscape of reasoning within AI, evaluating benchmark datasets, model evaluation protocols, and divergent approaches (including both neural, symbolic, and hybrid paradigms). We analyze how these benchmarks and reasoning tasks have co-evolved with state-of-the-art models, revealing not only strengths, but exposing gaps that persist in robustness, generalization, and interpretability.

Benchmarking remains foundational for tracking intelligence progress, motivating rigorous evaluation of reasoning in environments spanning language, vision, multi-modal inputs, and interactive tasks. Comparative studies increasingly draw attention to the merit and limitation of widely adopted datasets and evaluation protocols, highlighting their impact on the apparent progress of current models. Ensuring that benchmarking procedures genuinely diagnose underlying reasoning abilities—rather than pattern memorization or dataset artifacts—is vital for honest scientific progress. This survey contrasts various reasoning benchmarks and summarizes these in Table 2, which showcases the diversity, coverage, and targeted reasoning skills across leading datasets.

Architectural innovations play a central role in advancing AI reasoning. The field has seen a proliferation of approaches, from end-to-end neural methods (e.g., transformers), symbolic systems, to hybrid models fusing connectionist and logic-based reasoning. While transformer-based architectures have yielded impressive empirical results, critical evaluations probe their capacity for systematic generalization, compositional reasoning, and multi-step logical inference. We explicitly address critiques of these models, and juxtapose neural and hybrid strategies, synthesizing their respective advantages and open challenges.

In reevaluating these themes, we provide a focused engagement with alternative perspectives, including discussions of criticisms leveled against dominant paradigms (such as hidden brittleness or superficial learning in transformers and hybrid architectures). Where appropriate, we summarize these competing viewpoints and highlight ongoing debates regarding transparency, fairness, and practical adoption.

This survey is structured as follows. Section 2 details the reasoning benchmarks and their evaluation methodologies, with in-text summary tables reinforcing critical comparative insights. Section 3 analyzes architectural families, summarizing hybrid and alternative reasoning approaches. We conclude with a discussion of current challenges and future outlook, presenting a distilled summary of key takeaways at the close of each section.

By deeply engaging recent literature and benchmarking evolutions, this survey seeks to both inform and critically examine the

Table 1: Representative Reasoning Benchmarks: Domains and Key Evaluation Aspects

| Benchmark | Domain | Core Reasoning Skills | Evaluation Protocols |
|-----------|----------------|--|------------------------------------|
| BoolQ | Language | Boolean reasoning, reading comprehension | Accuracy, human verification |
| DROP | Language | Discrete operations, multi-step reasoning | Exact match, precision/recall |
| CLEVR | Visual | Compositional, relational reasoning | Program execution, accuracy |
| ARC | Language/Logic | Common-sense, abductive, analogy reasoning | Human baselines, automated scoring |
| HotpotQA | Multi-modal | Multi-hop, supporting fact identification | Exact match, F1, supporting facts |

trajectory of AI reasoning research, equipping researchers with an integrated view for future inquiry.

1.1 Overview of Large Language and Foundation Models (LLMs)

The evolution of artificial intelligence (AI) has been profoundly shaped by advances in language understanding and generation. The trajectory spans from symbolic, rule-based systems—characterized by explicit grammatical rules and formal symbolic manipulation—to statistical methods, neural, and deep learning architectures. Early symbolic approaches excelled in interpretability but were hindered by a lack of scalability and the brittleness of handcrafted rules. The advent of statistical models, and subsequently neural network architectures, marked a paradigm shift by enabling data-driven learning of complex linguistic patterns. This progression culminates in large-scale Transformer-based models, wherein pre-trained language models (PLMs), especially large language models (LLMs), distinguish themselves through scale and the emergence of novel capabilities.

Distinctive behaviors—such as in-context learning and abstract reasoning—emerge in LLMs due not solely to increased parameter counts, but also to innovations in model design, architecture, and training paradigms. Key developments include:

The adoption of large-scale, unsupervised pre-training; Attention mechanisms, as popularized by the Transformer architecture; Alignment of model objectives with downstream utility.

The launch and societal integration of models such as ChatGPT exemplify LLMs’ transformative effect on not only conventional natural language processing (NLP) tasks, but also on digital interaction, information retrieval, content creation, and scientific discovery. Concomitantly, there has been renewed interest in hybrid algorithmic-neural approaches and neural-symbolic (NeSy) systems. These are motivated by enduring challenges—particularly in reasoning and interpretability—where pure neural architectures, despite their success, fall short [69, 91, 95, 105]. A move toward models exhibiting compositionality and explicit knowledge manipulation reflects the AI community’s recognition that human-like reasoning and adaptability may require synthesizing symbolic and subsymbolic learning, an imperative for ongoing advancements toward artificial general intelligence (AGI).

1.2 The Critical Role of Reasoning, Replicability, and Benchmarking

The expanded potential of LLMs introduces foundational challenges. Chief among these is cultivating robust reasoning abilities within

LLMs that transcend mere pattern recognition or correlation. Although large-scale models demonstrate emergent capabilities in abstract reasoning and commonsense inference, such performance is inconsistent—often susceptible to dataset biases and lacking true compositionality. This motivates the investigation of model architectures and inductive biases that explicitly encode algorithmic or symbolic reasoning procedures.

Neural-symbolic computing (NeSy) has emerged as a promising paradigm, aiming to combine the transparent manipulation of knowledge found in symbolic systems with the flexible data-driven learning of neural networks. Empirical advancements within NeSy frameworks attest to concrete progress in domains demanding structured reasoning—such as scientific discovery, mathematical problem solving, and knowledge-intensive tasks—where traditional end-to-end neural models frequently encounter limitations. Despite these strides, major challenges persist:

Scalability of hybrid models integrating large structured knowledge bases; Efficient inference and reasoning over complex data; Achieving compositional generalization beyond seen examples; Seamless integration of symbolic knowledge acquisition into neural learning pipelines.

These open research problems highlight the incomplete nature of current methods and the ongoing need for innovation in neural-symbolic integration [69, 95].

As LLMs proliferate in research and industry, the importance of replicability and robust benchmarking has intensified. Widely-used evaluation metrics often fail to accurately reflect the subtlety of advanced reasoning behaviors and the adaptability required by practical deployments. This gap necessitates the development of comprehensive benchmarks addressing not only accuracy but also properties such as robustness, out-of-distribution generalization, and fairness. Compounding these technical challenges are issues of opacity and reproducibility, as proprietary models and undisclosed datasets undermine transparency and accountability in both research and societal applications.

Allied to these technical and practical challenges are pressing societal, ethical, and policy considerations—spanning algorithmic bias, misinformation, and the impacts of automating language-centric labor. Therefore, cultivating rigorous, transparent, and replicable research practices constitutes a linchpin for both scientific progress and public trust in LLM technologies [69, 91, 105].

1.3 Survey Structure and Scope

Given these multifaceted themes, this survey provides a structured synthesis of the technical, methodological, and societal dimensions defining contemporary LLM research. The survey first examines

evaluation methodologies and benchmarking strategies, with an emphasis on recent advances in linguistic competence, robustness, and inclusivity. Subsequently, the intersection of LLMs with algorithmic reasoning and neural-symbolic integration is scrutinized, highlighting technical obstacles and emerging opportunities in the quest for more reliable and general AI systems. Principles of open science and reproducible research are afforded particular attention, acknowledging their foundational role in mitigating societal risks and advancing the field. By organizing the discussion thematically, the survey seeks to equip readers with a critical appreciation of both progress to date and the grand challenges shaping the next frontier of large-scale, language-centric AI [69, 91, 95, 105].

2 Historical and Foundational Landscape

This section reviews the foundational approaches and key developments that have shaped the evolution of AI, with particular attention to reasoning architectures and benchmarking methodologies. We frame the analysis to highlight both prevailing paradigms and alternative perspectives, providing comparisons where relevant to clarify their respective strengths and limitations.

2.1 Early Approaches and Hybrid Models

The historical trajectory of AI reasoning systems has been characterized by an initial dominance of symbolic methods, including expert systems and rule-based engines. These approaches offered transparency and explicit logic structuring but often struggled to scale or handle ambiguity. The subsequent emergence of connectionist models introduced learning-based solutions, trading off interpretability for empirical performance improvements.

Hybrid models, combining symbolic and subsymbolic techniques, have been proposed to bridge these shortcomings. While hybridization seeks a synthesis between structure and flexibility, critics have argued that such approaches can inherit limitations from both parent paradigms, such as the brittleness of symbolic reasoning and the opacity of neural systems. The debate remains active, and a nuanced appraisal of these models is essential when considering their theoretical and practical implications.

2.2 Benchmarking and Reasoning Evaluation

Benchmarks play a vital role in evaluating the progress of reasoning systems. Early benchmarks focused on narrow, well-defined logical tasks, permitting rigorous comparison but often failing to represent real-world complexity. Over time, the field has moved toward more diverse and challenging benchmarks that span language understanding, abstraction, and multi-step reasoning.

Table 2 provides an overview of representative benchmarks, their focus, and evaluative criteria, reinforcing the diversity and evolution of reasoning assessment.

2.3 Transformers and Recent Paradigms

The advent of transformer architectures has markedly shifted the landscape of both perception and reasoning. These models have demonstrated unprecedented performance across benchmarks but prompted debate regarding the genuine nature of their reasoning abilities versus statistical pattern recognition. Competing views

question whether the inductive capabilities observed in transformers should be considered reasoning in the classical sense or rather as an emergent byproduct of large-scale data assimilation.

Critiques have also centered on the interpretability and controllability of such models, with some arguing that their success challenges traditional definitions of reasoning and intelligence. This ongoing discourse underscores the need for nuanced evaluation strategies and theoretical frameworks that can accommodate the complexity of modern AI.

2.4 Section Summary

In summary, the historical and foundational landscape of AI reasoning encompasses a rich interplay between symbolic approaches, connectionist models, hybrid architectures, and recent transformer-based advances. Each paradigm brings distinct advantages and trade-offs, reflected in the evolving design of benchmarks and evaluation criteria. Ongoing debates regarding hybrid systems and transformer-based reasoning highlight the field's dynamism and the importance of comprehensive, comparative assessment.

2.5 Evolution of Reasoning in AI

The evolution of artificial reasoning systems traces a trajectory from the early predominance of explicit symbolic logic frameworks to the contemporary dominance of neural and, most recently, transformer-based paradigms. Classical AI systems were grounded in symbolic representations, rule-based inference mechanisms, and logic programming, which were celebrated for their interpretability and transparency [69, 91, 95, 105]. These approaches enabled precise deductive reasoning but were notably brittle when applied to open, ambiguous, or real-world domains and required extensive manual creation of knowledge bases [91].

The subsequent rise of connectionist models, and in particular deep neural networks, marked a fundamental paradigm shift towards data-driven learning. Such architectures allowed for the automatic synthesis of hierarchical abstractions and enabled AI systems to address a wide range of reasoning problems without manually engineered logic [69]. Nevertheless, conventional neural networks suffered from persistent weaknesses in generalization, specifically on tasks requiring compositionality, recursion, or algorithmic processing—challenges where the strengths of symbolic methods remained salient, as seen in arithmetic, logic, combinatorics, and structured multi-step reasoning [95, 105]. To bridge these gaps, hybrid neural-symbolic (NeSy) models were developed to combine the perceptual strengths of neural networks with the explicit, interpretable inference of symbolic modules [69, 95]. These integrated frameworks demonstrated enhanced performance in mathematical problem-solving, retrosynthetic analysis, and other domains demanding multi-step reasoning [95]. Nevertheless, significant challenges remain, including the achievement of robust compositional generalization, scalable reasoning over extensive knowledge bases, and seamless integration between symbolic and neural paradigms, making their unification a central open problem [91, 95].

Table 2: Key Benchmarks in AI Reasoning and Their Evaluative Focus

| Benchmark | Reasoning Type | Task Domain | Evaluation Criteria |
|---------------------|-----------------------|----------------------------|------------------------------------|
| Early Logic Puzzles | Symbolic Deduction | Mathematical/Logical | Accuracy, Formal Correctness |
| Winograd Schema | Commonsense Reasoning | Natural Language | Disambiguation, Context-Dependence |
| bAbI Tasks | Multi-step Reasoning | Synthetic QA | Step-wise Inference, Scalability |
| ARC Challenge | Abstract Reasoning | Visual/Pattern Recognition | Generalization, Abstraction |

Concurrently, the field has been transformed by the advent and rapid maturation of transformer-based large language models (LLMs)—such as GPT, T5, PaLM, LLaMA, and Flan—each empowered by pre-training on massive corpora [10, 25, 38, 45, 54, 55, 64, 76, 91, 99, 105]. These models display emergent reasoning abilities across arithmetic, logic, and algorithmic tasks, especially when enhanced through advanced prompting techniques such as chain-of-thought (CoT) prompting [38, 99, 105]. CoT prompting elicits intermediate reasoning steps and has been shown to substantially improve model performance on complex multi-step problems compared to standard zero- or few-shot prompting [38, 99, 105]. For example, providing a few chain-of-thought exemplars enables even very large LLMs to outperform strong fine-tuned baselines on math word problems [99]. However, systematic evaluations point to persistent performance gaps between the reasoning capabilities of current LLMs and those of human experts, especially on tasks that require systematic abstraction, compositional logic, or the integration of broad world knowledge [10, 45, 64]. Even the most advanced models—such as GPT-4—can generate compelling rationales for complex clinical or scientific problems, thus contributing to their interpretability [76], but they still make frequent logical errors and remain vulnerable to superficial pattern matching or hallucination, particularly when required to generalize beyond their training distribution [25, 55, 76, 105].

Empirical studies further highlight that LLM performance on reasoning tasks is highly contingent on prompt construction, exemplar selection, and mechanisms for knowledge retrieval [10, 38, 54, 76, 99]. Critical reasoning failures persist in areas such as multi-step logical inference, combinatorial puzzles, and causal reasoning [10, 38, 76]. Retrieval-augmented CoT prompting has produced improvements, particularly in scientific and mathematical multimodal domains, by dynamically incorporating relevant external information [54, 99, 105]; yet, these advances are often incremental and do not fully overcome the brittleness of existing models on compositional generalization or causal inference [45, 55, 64, 91]. The cumulative evidence indicates that while transformer-based LLMs mark significant progress in automated reasoning, their abilities are largely emergent and stochastic rather than grounded in explicit abstraction or reliable causal modeling. These limitations motivate ongoing research into hybrid, neurosymbolic, and biologically inspired approaches to further advance AI reasoning [91, 95, 105].

2.6 Embedding and Model Architecture Developments

The foundation of modern natural language processing and reasoning systems is closely intertwined with advances in representation

learning—particularly in embedding methods—and architectural design. Early approaches utilized static, dense embeddings to encode lexical relationships; the transition to contextualized embeddings, most effectively realized in transformer architectures, represented a qualitative leap in modeling semantic, syntactic, and higher-order structural relations between tokens and modalities [23, 31, 79, 93]. Models such as BERT, GPT, and their derivatives leverage deeply stacked attention layers, enabling the encoding of rich, context-dependent linguistic meaning [75]. Techniques like SBERT-WK, which dynamically aggregate BERT’s internal representations, have further extended semantic alignment and resilience to contextual variation [23, 79, 93].

Transfer learning—particularly via pre-trained checkpoints from models such as BERT, GPT-2, and RoBERTa—has become a standard modality for adapting large-scale models to downstream tasks with minimal additional training [75]. This paradigm shift has substantially improved access to high-performing models, reducing resource requirements and enabling widespread success across tasks such as translation, summarization, and machine reasoning [75]. In parallel, innovations in self-supervised learning, multimodal integration, and speech-text modeling have expanded the capacity of transformer models to operate across text, image, tabular, and speech inputs [31, 93].

Notably, encoder-decoder architectures now explicitly model tabular structure or document salience to support long-context reasoning and accurate summarization, while retrieval-augmented systems incorporate external information to improve reasoning fidelity [75, 79, 93].

Despite these advances, significant architectural limitations endure:

Models frequently underperform when processing extended input contexts, with accuracy declining as relevant information is dispersed across longer sequences [55]. Standard embedding mechanisms, adept at capturing local and semantic dependencies, are less effective for structured data (e.g., tables, knowledge graphs) absent specialized architectural enhancements [75, 93]. Innovations such as structured attention, field-content selective encoders, and advanced pooling strategies are actively being explored to address these challenges [23, 75, 79, 93].

These ongoing research directions aim to bridge the gap between flexible, general-purpose architectures and the demands of explicitly structured or long-context reasoning tasks.

Table 3: Summary of foundational paradigms in AI reasoning, with comparative strengths and limitations.

| Paradigm | Core Mechanisms | Strengths | Key Limitations |
|---------------------------------------|---|--|---|
| Symbolic (Rule-based, Logic) | Explicit symbols, rules, logic programs | Interpretability, rigorous deduction, transparency | Brittle generalization, manual knowledge engineering |
| Neural (Connectionist, Deep Learning) | Hierarchical, distributed representations; learning from data | Strong pattern recognition, adaptability, implicit abstraction | Weakness in compositional reasoning, limited interpretability |
| Neural-Symbolic (Hybrid) | Joint neural and symbolic modules; integration architectures | Combines perception with explicit inference, improved generalization on structured tasks | Integration complexity, compositional generalization, scalability |
| Transformer-based LLMs | Attention-based contextual encoding; large-scale pre-training | Emergent reasoning, multi-task capability, scalability | Reliant on statistical learning, lacks explicit abstraction or robust causality |

2.7 Biological Inspirations and Neuromorphic Approaches

An increasingly impactful trajectory in the development of reasoning-enabled AI is the incorporation of principles drawn from biological and cognitive neuroscience. The structural organization and dynamic properties of biological connectomes are widely hypothesized to underpin the cognitive flexibility and generalization observed in human reasoning. Inspired by this, neuromorphic systems and reservoir computing models have been designed to emulate key features of brain networks, notably modularity and criticality [83]. Recent empirical findings suggest that reservoir computing architectures that incorporate brain-inspired topologies consistently outperform architectures with random connectivity, particularly on tasks requiring flexible generalization and adaptive reasoning capabilities. This highlights the computational advantages inherent in functional segregation and integrated network dynamics [83].

The manifold benefits of biologically inspired architectures can be summarized as follows: they serve as explanatory models for the origins of cognitive flexibility and compositionality in biological reasoning systems [83]; they guide the development of artificial reasoning systems with enhanced efficiency, adaptability, and robustness—especially in contexts characterized by uncertainty and ambiguity; and they inspire integrative approaches that blend cognitive, neural, and symbolic paradigms, targeting the recursive and adaptive reasoning abilities found in biological intelligence [83, 95].

In summary, the historical and foundational landscape of AI reasoning is shaped by the interplay between symbolic, neural, and hybrid paradigms; innovations in knowledge representation and network architecture; and the growing influence of neuroscience-inspired methodologies. Each trajectory provides distinct strengths and faces unique limitations (see Table 3), collectively shaping and informing the ongoing evolution and future directions of reasoning-enabled AI research [69, 75, 83, 91, 95, 105].

3 Benchmarking Speech and Language Models

3.1 Standardized Frameworks and Leaderboards

The evaluation of speech and language models has evolved significantly due to the emergence of standardized benchmarking frameworks and public leaderboards, which facilitate systematic assessments of generalization, robustness, and task coverage. In the domain of speech processing, the Speech processing Universal PERformance Benchmark (SUPERB) serves as a comprehensive, extensible, and reproducible platform designed to evaluate foundation models across a diverse set of 15 tasks, including phoneme recognition, keyword spotting, speaker identification, and automatic speech recognition. Through unified evaluation protocols and methodically constructed multi-task procedures—such as fixed feature encoders paired with task-specific prediction heads, and statistically robust metric aggregation—SUPERB enables rigorous

comparison across 33 models encompassing both self-supervised and conventional paradigms. Importantly, SUPERB’s insistence on reproducibility, robust statistical testing, and open-source benchmarking resources has accelerated the community’s ability to reach consensus on model performance and limitations, while also revealing persistent vulnerabilities in generative and low-resource scenarios [68, 103].

Analogously, natural language processing (NLP) relies on frameworks such as HELM and DIOIR to offer methodologically robust, scenario-based leaderboards that encompass a broad range of multi-domain tasks spanning Wikipedia and news text to biomedical corpora. These frameworks transcend surface-level metrics by incorporating evaluations centered on societal impact, reliability, and efficiency [47, 68]. The deliberate inclusion of well-curated, domain-diverse datasets is vital: contemporary large-scale studies show that benchmark composition can appreciably influence model rankings and perceived performance, particularly as the range of covered domains and task types expands [47].

A noteworthy innovation is the introduction of continual learning benchmarks, such as CL-MASR for multilingual automatic speech recognition. These benchmarks systematically arrange sequences of tasks and languages specifically to reveal deficiencies in models’ capacity to acquire new skills without succumbing to catastrophic forgetting. The CL-MASR benchmark supplies standardized, reproducible task sequences and a comprehensive suite of metrics—including Word Error Rate, measures of forgetting, backward transfer, and intransigence—to facilitate systematic evaluation of catastrophic forgetting, cross-lingual interference, and data/resource imbalance issues, especially in low-resource or highly typologically diverse environments. Furthermore, the open-source nature of CL-MASR advances direct reproducibility and collaborative method development within the community [53]. Collectively, these trends illustrate the rising expectations for multi-domain, resource-robust, and reproducible benchmarking in both speech and language modeling research.

3.2 Evaluation Metrics and Best Practices

As introduced in our survey, a central objective is to systematically assess how benchmarking practices and evaluation metrics shape the reliability, interpretability, and real-world relevance of model evaluations across AI domains. This section advances that goal by critically examining current metric selection practices and best practices in benchmark design, with a focus on alignment to human-centered objectives and actionable guidance for future development.

The effectiveness of benchmarks is fundamentally dependent on the alignment between evaluation metrics and human-centered objectives. Automated metrics including ROUGE, BLEU, and METEOR have long served as mainstays in tasks such as summarization, simplification, and machine translation. However, these metrics typically correlate only weakly with human judgments of meaning,

comprehension, and utility, particularly for complex tasks such as plain language summarization and biomedical natural language processing [16, 28, 36, 47, 68, 81, 103]. For instance, recent work in medical plain language summarization found that, while ROUGE and similar metrics suggest LLM-generated outputs are comparable to human writing, objective comprehension tests with lay participants reveal a substantial gap: only QA-based metrics like QAEval reflect true understandability and faithfulness [36]. Likewise, in chemical space exploration and biomedical NLP, surface-level metrics may fail to distinguish between models of genuinely different quality, necessitating more semantically informed alternatives [16, 81].

To address these gaps, evaluation approaches have shifted toward semantically grounded metrics that better reflect human preferences and understanding. Methods such as cross-encoder or bi-encoder models, fine-tuned for semantic similarity or natural language inference, now consistently outperform traditional n-gram overlap measures in both general and domain-specialized contexts [16, 28, 81]. For example, leveraging inference-based or QA-based metrics, as shown in biomedical and simplification benchmarks, provides stronger alignment with human assessments of comprehension and utility, especially in layperson-facing and specialized language generation applications [16, 28, 36, 47].

Despite such advancements, several challenges in metric selection persist and have, on occasion, led to misleading conclusions in the field. For example, leaderboard rankings can prove highly volatile: analyses of Decision Impact on Reliability (DIoR) within the HELM benchmark demonstrate that moderate changes in scenario grouping, dataset selection, or evaluation aggregation can unpredictably shift the relative standing of language models, sometimes reversing previous conclusions about model performance [68]. In the SUPERB speech benchmark, statistical analyses indicate that observed leaderboard differences among top models are often statistically insignificant, cautioning against over-interpretation of small performance gaps [103]. In sentence simplification (BLESS) and biomedical NLP, evaluation instability remains a concern, as rankings shift with metric choice and evaluation setup [16, 47].

The resultant volatility of metric-based leaderboards in response to such changes underlines the necessity for actionable best practices. We recommend transparent and precise definition of metrics, comprehensive statistical reporting (including significance testing to avoid misinterpretation of minor differences), and clear articulation of scenario aggregation and evaluation protocols [39, 68, 103]. Furthermore, composite or scenario-weighted evaluation methodologies are increasingly advocated to ensure reliable and representative assessment across model capabilities [39, 47, 68]. Recent benchmarks, such as BLESS for sentence simplification [47] and the Speech processing Universal PERFORMANCE Benchmark (SUPERB) [103], exemplify the trend toward domain-specific, multi-faceted evaluation frameworks with rigorous reproducibility and statistical safeguards. These works stress open-source code, publicly available datasets, and reproducible pipelines as foundational for community trust and scientific rigor [39, 47, 103].

For benchmark and metric developers, these insights yield several actionable recommendations: prioritize semantically and comprehensively grounded metrics over surface-level measures; systematically report statistical significance of leaderboard differences; design evaluation protocols to minimize volatility induced by scenario or aggregation

choices; and ensure all datasets, code, and evaluation procedures are public and thoroughly documented. The growing body of benchmarks from late 2023 and 2024, such as BLESS and new HELM variants, reinforce these principles and offer blueprints for future robust, human-aligned evaluation design [47, 68, 103].

To ensure reproducibility and scientific rigor, it is imperative to publicly release datasets, code, evaluation procedures, and, when feasible, simulated or derived data, as recommended by standards in applied linguistics and benchmarking research [39].

3.3 Comparative Analysis and Diversity

This section advances the overall survey objective of critically synthesizing current benchmarking practices for language and foundation models, with a particular emphasis on cross-domain applicability, methodological rigor, and implications for the evolution of evaluation standards. As models and evaluation methodologies diversify, a nuanced understanding of comparative trends, volatility, and benchmark robustness is essential for guiding model assessment and development.

A principal focus in recent benchmarking efforts is the systematic comparison of large language models (LLMs) and foundation models with both traditional baselines and alternative architectures. Cross-domain benchmarking—such as evaluating state-of-the-art (SOTA) fine-tuned models (e.g., BioBERT, PubMedBERT, BART) versus LLMs (e.g., GPT, LLaMA) in biomedical NLP—shows that while LLMs frequently achieve superior performance on tasks requiring generative reasoning or medical question-answering, they often do so at a substantially higher computational cost. Moreover, without additional task-specific adaptation, LLMs may still lag behind fine-tuned models in extraction, classification, and domain-specialized settings [47]. For instance, generative models like GPT-4 tend to produce outputs of high fluency for summarization and simplification, but these may be less complete or more susceptible to hallucinations compared to specialized baselines. Furthermore, marked variability is observed in the repertoire of edit operations and strategies employed by different LLMs in tasks such as text simplification, indicating heterogeneity in their methodological approaches [47].

For a concise overview of such comparative results, see Table 4.

Recent studies have highlighted volatility in benchmark outcomes, especially where evaluation protocols or scenario composition fluctuate. For example, the work of Perlitz et al. [68] on the HELM benchmark demonstrates that simply adding or removing models or datasets can alter leaderboard rankings and perceived model superiority, sometimes misleading the field about genuine progress. Notably, aggregation strategies like grouping diverse datasets may yield lower evaluation reliability, and an overemphasis on the number of test examples may not translate to greater stability. This indicates the need for statistically grounded metrics, such as Decision Impact on Reliability (DIoR) [68], when designing benchmarks. Similarly, in continual learning for speech recognition, Della Libera et al. [53] show that ordering of languages or choice of resource splits can skew comparability; certain strategies overstate model resilience due to scenario sequencing rather than true model robustness.

For developers of benchmarks and metrics, these findings advocate for several best practices: (1) Articulate and minimize sources of

Table 4: Representative comparative outcomes between SOTA fine-tuned models and LLMs on biomedical NLP tasks. Values indicate relative strengths as identified in recent benchmarking studies.

| Task | BioBERT/BART | GPT-4/LLaMA | Notes |
|-----------------------------|--------------|----------------------|---|
| Extraction & Classification | Superior | Inferior | Fine-tuned models excel; require less adaptation |
| Medical QA | Moderate | Strong | LLMs perform well, esp. with complex queries |
| Generative Summarization | Moderate | Superior | LLMs enhance fluency, some risk of hallucination |
| Text Simplification | Specialized | Diverse | LLMs deliver varied strategies and edit diversity |
| Computational Cost | Efficient | Substantially Higher | LLMs demand greater resources |

volatility by transparently defining scenarios, datasets, and ranking metrics; (2) Employ reliability measures to evaluate the stability of results under perturbations of experimental setup; (3) Design with efficiency in mind, as both environmental and resource constraints are increasingly relevant—approaches like Flash-HELM [68] can offer computational savings without sacrificing reliability.

Periodically, the field is reoriented by the release of new benchmarks tailored for extensibility, multilinguality, or continual learning. Examples emerging in late 2023 and 2024 include BLESS [47] (targeting LLM evaluation on sentence simplification with analyses of edit diversity and robustness) and CL-MASR [53] (addressing continual learning in multilingual ASR, with evaluation protocols probing catastrophic forgetting, transfer, and efficiency). These benchmarks are accompanied by open-source resources and standardized evaluation frameworks to facilitate reproducibility and sustained advancement.

Benchmark development has also prioritized diversity and inclusivity, with a marked shift toward constructing resources that encompass broader linguistic, cultural, and task-scale variability. These advances ensure fairness in model assessment and promote research that generalizes beyond canonical datasets or majority language contexts [47]. Emerging benchmarks are designed for extensibility and adaptability, supporting, for instance, multilingual task sequences or modular scenario expansion, while emphasizing open sharing of resources to catalyze community-led progress [47, 53, 68]. Adherence to these principles in benchmark creation and deployment enables robust comparative analyses and is essential for driving sustainable progress in speech and language modeling research.

Systematic benchmarking using unified frameworks and rigorous protocols advances community consensus on model strengths and weaknesses. The ongoing evolution of evaluation metrics emphasizes alignment with human judgment, particularly in complex and layperson-facing tasks. Comparative studies reveal fundamental trade-offs between LLMs and fine-tuned domain-specific models, reinforcing the ongoing need for careful task adaptation and judicious resource allocation. Finally, foregrounding diversity, extensibility, and open scientific practices will be essential for future-proofing benchmarks and maximizing their impact across domains.

4 Probing, Reasoning, and Linguistic Competence Benchmarks

This section aims to provide a comprehensive overview of benchmarks that assess language model capabilities through probing

tasks, reasoning challenges, and the evaluation of linguistic competence. Our objective is to clarify the purpose and structure of prominent benchmarks within this space, highlighting their design philosophies, target skills, and relevance to the AI and NLP research communities. The section is particularly valuable for researchers and practitioners seeking to understand, evaluate, or develop models with robust linguistic and cognitive abilities. We begin with a general introduction to each type of benchmark and then examine their methodologies and applications. These insights are intended to support readers in identifying appropriate evaluation suites for their specific domains and use cases.

4.1 Linguistic and Reasoning Probing

In alignment with the core objective of this survey—to critically examine and synthesize advances and outstanding challenges in the evaluation of large language models (LLMs)—this section focuses on probing methodologies that target the linguistic, reasoning, and abstraction abilities of state-of-the-art models. We explicitly consider how evolving benchmarks expose both progress and persistent gaps, and highlight consequential lessons for the development of future metrics and evaluation frameworks.

The evaluation of LLMs increasingly depends on sophisticated probing techniques designed to reveal the nuanced properties of models' internal representations and linguistic behaviors. The evolution of probing for syntactic and semantic competence has progressed from elementary acceptability judgments to methodologically robust frameworks, which now target compositional and structural facets of language. Modern benchmarks, for instance the Two Word Test (TWT), probe models on foundational aspects of semantic composition: specifically, their ability to distinguish between plausible and implausible noun-noun phrases. Crucially, success in this domain requires not just recognition of word similarity but a deeper grasp of semantic combinatorics. Although LLMs demonstrate impressive performance on complex downstream tasks, empirical evidence shows they continue to struggle with the core challenge of semantic discernment. Notably, models such as GPT-4 variants recurrently overestimate the coherence and meaning of nonsensical phrases, indicating a persistent reliance on surface-level statistics (e.g., vector cosine similarity) over robust compositional understanding [73]. This persistent gap highlights a critical mismatch between reported advancements on aggregate language benchmarks and true progress in core linguistic competence.

Highlighting the volatility of benchmark-based conclusions, TWT results [73] show that models like GPT-3.5-turbo and Gemini-1.0-Pro-001 rate nonsensical noun-noun pairs almost as highly as meaningful ones, misleadingly suggesting human-like semantic competence when judged solely by high-level accuracy or unrelated benchmarks. Such volatility has, at times, misdirected perceived progress in the field: models excelling on verbose or logic-heavy benchmarks may still lack fundamental linguistic understanding, as exposed by carefully constructed tests like TWT. Similarly, in the context of metrics for generative chemical models, research has revealed that widely-used metrics often fail to accurately reflect true model quality or generalization ability, prompting a reassessment of which benchmarks genuinely probe for intended competencies [81].

In parallel, syntactic minimal pair benchmarks, exemplified by BLiMP, systematically evaluate models across an extensive array of morphosyntactic phenomena. BLiMP, through its template-generated sentence pairs, isolates specific grammatical constructs and tests models' sensitivity to grammaticality [92, 97]. While transformer-based models consistently surpass earlier n-gram and LSTM-based language models in phenomena such as subject-verb agreement, they remain prone to inconsistency when faced with deeper syntactic generalizations, including negative polarity and island constraints. This brittleness is further corroborated by classifier-based probing studies, notably Holmes and its computationally optimized extension FlashHolmes, which aggregate results across more than two hundred datasets and encompass a spectrum of phenomena in syntax, morphology, semantics, and discourse [15, 92]. Analysis from Holmes-based studies reveal expected scaling of competence with increased model size, yet also expose nontrivial dependencies on architectural choices and instruction tuning—these effects are especially evident within morphosyntactic domains, thereby emphasizing the importance of both inductive biases and fine-tuning paradigms.

Recent research extends the probing paradigm to include reasoning and abstraction ability, utilizing an increasingly diverse suite of benchmarks. Notably, the Abstraction and Reasoning Corpus (ARC) and subsequent developments within the DreamCoder/PeARL frameworks have shifted focus toward generalization over pattern recognition. Whereas neurosymbolic approaches like DreamCoder specialize in structured transformations via program induction, LLM-based methods augmented with novel encodings and data augmentations excel at orthogonal aspects, with each paradigm addressing complementary subsets of ARC tasks [9, 15, 81, 100]. For example, PeARL [9], introduced in 2024, advances recognition models for ARC and demonstrates that neither neurosymbolic nor LLM pipelines can independently solve a majority of cases, but their ensemble achieves improved coverage, surpassing prior approaches such as Icecuber. Ensemble approaches achieve broader coverage, yet no single paradigm independently solves a majority of cases, illustrating the persistent difficulty of abstract reasoning and broad generalization [9, 15, 100]. The release of the open-source *arckit* library [9] further emphasizes the trend toward reproducible, extensible benchmarking environments.

The recent introduction of RGB (Retrieval-Augmented Generation Benchmark) in late 2023 [15] represents a notable advance in evaluating the integration of retrieval and generative capabilities. RGB systematically examines LLMs' abilities in noise robustness,

negative rejection, information integration, and counterfactual robustness, revealing bottlenecks such as the inability to reliably refuse unsupported questions, sharp performance drops with increased noise, and consistent struggles when integrating information across documents. The authors urge for careful metric construction and caution against over-interpretation of aggregate scores, providing actionable guidance for both benchmark and metric developers to focus on error detection, document modeling, and cross-document reasoning.

Specialized domains have further spurred the development of targeted benchmarks. Biomedical and clinical reasoning datasets, such as MedS-Bench and arckit, extend probing into domain-specific abstraction and reasoning. Recent work (2025) finds that even the most advanced LLMs, including GPT-4 and Claude-3.5, exhibit divergent abilities between real-world and multiple-choice scenarios; excelling at the latter but consistently underperforming on tasks requiring nuanced clinical information extraction or summarization [9, 15, 100]. The findings underscore the limitations of existing benchmarks in capturing real-world deployment challenges, and argue for a shift toward broader clinical scenario coverage, multilingual expansion, and the validation of metrics against actual task data.

Collectively, the evidence indicates that while advancements in probing and benchmark curation have refined our ability to diagnose LLM limitations, current state-of-the-art models remain highly sensitive to prompt formulation and task structure. Notable gaps persist in the domains of semantic composition, syntactic robustness, and genuine cross-domain abstraction [9, 15, 73, 92, 97]. The volatility and occasionally misleading nature of benchmark metrics highlight the need for granular, transparently designed evaluation tools. For researchers and benchmark developers, this underscores the importance of continued innovation in dataset design—prioritizing not only coverage and challenge diversity, but also the reproducibility, diagnostic depth, and alignment with real-world language demands.

4.2 Multi-modal and Cross-Validation Benchmarks

As LLMs are increasingly tasked with operation in multi-modal environments and expected to coordinate complex, multi-step reasoning processes across modalities, this survey seeks to critically evaluate how benchmarking infrastructures and evaluation protocols are adapting—and sometimes falling short—in reflecting these real-world complexities. Our overarching survey objective is to synthesize both the progress and persistent challenges in LLM evaluation, highlighting actionable directions for the creation and selection of more reliable, generalizable, and interpretable benchmarks amid rapid paradigm shifts.

Modern multi-modal and multi-view benchmarks evaluate not only models' linguistic capabilities, but also their aptitude for reasoning over—and integrating—representations from disparate information sources, including text, vision, speech, and structured data. This reflects the complexity and interconnected character of real-world scenarios [9, 16, 51, 100, 101]. However, volatility in benchmark performance can mislead the field: for example, earlier rapid gains on the ARC benchmark using LLMs and neurosymbolic

hybrids [9] led to overoptimistic assessments of machine abstraction and generalization, only for ensemble methods or new data augmentations to reveal major persistent gaps. Similarly, frequent metric recalibration in biomedical NLP challenges coincides with shifting leaderboard rankings that obscure true progress [16, 100].

Recent studies demonstrate that performance in multi-modal chain-of-thought (CoT) tasks can be significantly enhanced through retrieval-augmented prompting techniques. Cross-modal demonstration selection and stratified sampling have proven especially effective in benchmarks such as ScienceQA and MathVista. For instance, retrieval mechanisms that align intra- and inter-modality information, when combined with strategic sampling, have enabled GPT-4-based models to achieve unprecedented benchmark scores and surpass previous generation methods by substantial margins [51]. Ablation studies underline the necessity of both visual knowledge integration and diverse demonstrations for optimal performance.

Contemporary evaluation frameworks increasingly incorporate clustering and latent space analysis to validate model reasoning and clarify interpretability. Deep clustering strategies, particularly those maximizing mutual information or leveraging hierarchical adversarial networks, reveal that the emergence of robust and interpretable clusters is strongly associated with improved cross-modal generalization, and provide essential insights into where model abstraction failures occur [101]. Meanwhile, cross-validation protocols now extend far beyond conventional train/test splits, embracing explicit tests on out-of-domain and counterfactual instances to rigorously scrutinize generalization and model robustness [16, 100]. Notably, late 2023 and early 2024 have seen the introduction of more specialized, open-access benchmarks such as MedS-Bench for comprehensive clinical LLM assessment, and expanded ARC-based variants for nuanced abstraction and reasoning diagnostics [9, 100].

A persistent element in this area is direct human-model comparative analysis. Such studies consistently show a substantial gap between current LLMs and human performance, particularly in robustness to noise, rejection of negative or irrelevant answers, and the integration of information across multiple documents or modalities [9, 16, 100]. Models are frequently highly accurate under ideal (clean) conditions, but their performance deteriorates rapidly in the presence of noise. Another prevailing problem is the safe and consistent refusal of unsupported or nonsensical queries, which remains unresolved and underscores the ongoing need for semantic alignment and reliable evidence attribution [16, 100].

For benchmark and metric developers, several actionable lessons emerge. It is critical to diversify evaluation data by including out-of-domain and counterfactual scenarios, emphasize interpretability through clustering and error analyses, and avoid over-reliance on static leaderboards that obscure weaknesses due to benchmarking volatility. The field should prioritize open-access, community-driven benchmarks with active real-world validation, especially in dynamic domains such as healthcare and abstraction-oriented reasoning [9, 16, 100].

In summary, while multi-modal and cross-validation benchmarks have undeniably propelled progress in realistic, multi-dimensional

LLM reasoning, they also systematically catalog enduring brittleness. Model failures tend to cluster around areas that require integration, abstraction, or robustness—the very hallmarks of human cognitive prowess [9, 16, 51, 100, 101].

4.3 Comprehensive Benchmark Surveys and Limitations

To support the overarching objectives of this survey—namely, to critically analyze the landscape of LLM and agentic system evaluation, clarify methodological pitfalls, and distill actionable guidance for effective benchmark and metric development—this section synthesizes insights from the most recent comparative surveys and systematic reviews. The aim is to contextualize benchmarking practices, limitations, and recent evolutions within the broader goal of advancing robust, meaningful assessment frameworks for language models and related AI systems.

The advent of increasingly powerful LLMs and agentic systems has driven a proliferation of benchmarks spanning question answering, reasoning, linguistic competence, domain-specific tasks, and multi-modal evaluation. This phenomenon is rigorously documented in comparative surveys and systematic reviews [7, 8, 14, 23, 26, 27, 35, 38, 40, 41, 43, 44, 50, 52, 54, 60, 62–64, 74, 76, 79, 86, 87, 89–91, 94, 99, 102, 107, 108]. These surveys have introduced nuanced taxonomies and meta-frameworks for benchmarking, systematically dissecting evaluation practices across knowledge extraction, mathematical reasoning, code generation, factual retrieval, and a growing set of embodied or collaborative tasks.

A recurring critique within these surveys concerns the fragmentation and rapid evolution of the benchmarking ecosystem. Not only do they chronicle the expansion of benchmark tasks and methodologies, but they also caution against conflating benchmark score gains with meaningful advances in intelligence or generalization [40, 54, 94, 99, 107]. Volatility in benchmark results has led to notable instances where progress was overestimated, for example: static prompt templates have repeatedly led to apparent gains in language model knowledge, when in fact improvements stemmed from prompt selection rather than from more substantive advances in model reasoning or understanding [42, 97]. Similarly, a recent investigation into reasoning benchmarks found that many prompt engineering techniques—including chain-of-thought and specialized prompting—did not yield statistically significant, replicable improvements when re-tested on the latest LLMs, demonstrating that previously reported gains may not constitute robust or generalizable progress [90]. These cases illustrate how benchmarking volatility and overfitting can mislead perceptions of field advancement.

Comparative studies consistently highlight persistent deficiencies in compositionality, abstraction, and broad generalization. Many existing benchmarks fail to adequately test for the kind of causal or counterfactual reasoning that constitutes the core of human cognitive flexibility [9, 16, 42, 73, 92, 97]. Multi-modal and embodied benchmarks, although becoming more prevalent, continue to struggle with fragmentation and insufficient coverage of real-world or specialized domain contexts [9, 16, 100].

Surveys systematically catalog the methodological pitfalls that undermine benchmark validity and transferability: Methodological

artifacts and overfitting to fashionable datasets; Annotation biases and insufficient scenario diversity; Demographic and domain underrepresentation; Use of benchmarks lacking practical or scientific relevance; Overestimation of model capabilities due to suboptimal prompt strategies or template reliance.

For example, repeated overestimation of language model knowledge frequently arises from static prompt templates, which can mask the true limitations of underlying systems; this issue has been systematically demonstrated using corpus mining and minimal pair benchmarks for linguistic acceptability [42, 97].

Recent reviews emphasize that benchmark design is increasingly informed by calls for extensibility, transparency, and broad generalization. The movement toward open-source libraries and dynamic, extensible benchmarks has led to more rigorous cross-domain evaluation protocols [8, 14, 23, 26, 27, 41, 43, 44, 50, 54, 60, 62, 74, 79, 86, 91, 102, 107]. Still, even within these progressive paradigms, there is consensus that static, template-driven evaluations remain inadequate for capturing the dynamic, interactional, and cross-modal capabilities required for genuine human-level reasoning.

As highlighted in Table 5, each major benchmarking theme offers crucial diagnostic capabilities while simultaneously exposing core limitations that remain unresolved.

In line with the fast-moving field, several emerging benchmarks and critical reviews have surfaced since late 2023 and into 2024. For example, RepliBench [8], released in 2024, evaluates the autonomous replication capabilities of LLM agents, exposing gaps between partial and full autonomy in realistic operational scenarios. Holmes [92] provides an extensive, computation-efficient means to assess the true linguistic competence of LLMs, disentangling performative skill from deep syntactic or semantic understanding. In the biomedical and clinical domain, MedS-Bench and MedS-Ins [100] offer new, multi-faceted clinical task assessments, revealing that leading LLMs still underperform in real-world medical settings despite strong multiple-choice results. The Two Word Test (TWT) [73] further exposes persistent failings in core semantic compositionality, where LLMs struggle to reliably discriminate between meaningful and nonsensical phrase pairs—a fundamental gap from human-like linguistic intelligence.

For benchmark and metric developers, several actionable lessons emerge. Designing benchmarks should prioritize methodological transparency (including detailed documentation and open access to data/code [16, 27, 100]); support extensibility and domain coverage; incorporate robust statistical validation (including replicability checks [60, 90]); and move toward evaluation tasks that probe for compositional generalization, counterfactual and causal reasoning, and real-world scenario diversity. The field would also benefit from systematic reporting of negative or null results, increased demographic and task diversity, and the adoption of dynamic, adaptive evaluation strategies [16, 27, 50, 60, 86].

By methodically integrating advances in probing, multi-modal, and comprehensive benchmark design, the research community is forming a more nuanced and critical understanding of both the progress and the persistent limits of LLM capabilities. Notwithstanding significant strides, converging evidence from these diverse evaluation paradigms highlights enduring challenges for semantic composition, abstraction, and generalization. These findings underscore the necessity for methodological innovation and concerted,

cross-disciplinary approaches to benchmarking, if LLMs and related technologies are to achieve robust, real-world linguistic and reasoning competence [9, 16, 42, 73, 92, 97].

4.4 Knowledge Measurement, Prompt Engineering, and Model Adaptation

This section aims to clarify the key objectives and challenges in knowledge measurement, prompt engineering, and model adaptation for benchmarking and reasoning research. We target empirical researchers and practitioners interested in developing robust evaluation protocols and improving the reliability of AI models. Our focus is on outlining measurable goals for benchmarking, such as ensuring reproducibility, transparency, and fairness, and on offering recommendations for designing experiments and prompts.

We specifically address how benchmarking approaches gauge knowledge and reasoning proficiency, explore the evolving methodologies behind prompt engineering, and discuss recent advances and open challenges in adapting models to new domains or tasks. In this context, we highlight what is new in this survey compared to previous work, such as our synthesis of cross-benchmark adaptation techniques and attention to prompt robustness evaluation.

Researchers should consider the following guiding questions when evaluating or designing benchmarking protocols in this domain: (1) How does the knowledge measurement protocol account for both breadth and depth of model understanding? (2) What empirical strategies or prompt formats lead to more consistent reasoning outcomes across diverse benchmarks? (3) How can adaptation mechanisms be systematically evaluated for generalizability while avoiding overfitting to particular datasets or prompt types?

Open research challenges in these domains include quantifying the impact of prompt variability on reasoning stability, developing metrics for prompt sensitivity, and building adaptation protocols that support both rapid customization and principled assessment. Understanding the implications of these challenges is crucial; for instance, improved prompt engineering may reduce model brittleness and enhance real-world deployment reliability, while better model adaptation techniques can support more equitable and accessible AI systems.

By foregrounding these priorities and questions, this section provides a roadmap for future empirical benchmarking and advances clear criteria for evaluating progress in knowledge measurement, prompt engineering, and model adaptation.

4.4.1 Prompt-based Evaluation and Knowledge Probing. Prompt-based evaluation has become central to assessing the knowledge and reasoning abilities of large language models (LLMs). Benchmarks like the LAMA probe make use of cloze-style prompts to gauge factual recall. However, studies show that these prompts often underestimate the knowledge present in a model, as their rigid syntactic structure and lack of paraphrastic diversity limit what can be elicited [42]. Innovations such as paraphrasing-based and mining-based prompt generation, as implemented in the LPAQA suite, demonstrate that systematically creating diverse and high-quality prompts can extract considerably more knowledge from models—with up to an 8.5% absolute improvement on LAMA reported through these methods. This leads to more reliable lower

Table 5: Comparison of Major Benchmark Themes and Identified Limitations

| Benchmark Domain | Strengths | Key Limitations |
|---------------------------------|---|--|
| Linguistic and Reasoning Probes | Fine-grained diagnosis of syntax, semantics, and abstraction; reveal scaling/architecture effects | Surface-level overfitting; brittleness in compositionality and deep generalizations |
| Multi-modal/Multi-view | Integration of cross-domain modalities; improved realism; rich performance metrics | Brittleness under noise; limited robustness; persistent gap to human-level integration |
| Comprehensive Surveys | Systematic taxonomy; meta-analysis; identification of research gaps and risks | Fragmentation; overfitting to benchmarks; lack of causal/counterfactual probing |

bounds on model knowledge, highlighting the importance of prompt formulation in evaluation [42].

Nevertheless, expanding prompt diversity introduces major challenges. Most significant is prompt sensitivity: minor changes in how a question is phrased can cause large swings in answer accuracy. This results in instability both within and across experiments, making it difficult to robustly compare outcomes between studies. In addition, prompt-based benchmarks focused on factual recall or compositionality (such as the Two Word Test, TWT) expose that even leading LLMs struggle to reliably distinguish meaningful phrases from nonsensical ones. These models often respond based on superficial similarities in words or vectors, as opposed to demonstrating genuine understanding of compositional semantics—a weakness not mirrored in human performance [73]. Such findings stress that high performance on tailored tasks should not be conflated with deep language understanding.

Despite substantial progress in designing new benchmarks, three persistent limitations undermine prompt-based knowledge measurement: susceptibility to artifacts and syntactic cues present in surface text; high and unpredictable variability when prompts are paraphrased; and a lack of robustness and reproducibility of results, especially under varying experimental conditions.

These issues are further pronounced in specialized fields such as the biomedical and clinical domains. Here, the complexity and specificity of terminologies and domain schemas amplify inconsistency in model responses and complicate generalization [16, 100]. Transparent and comparable evaluation thus necessitates open access to probing datasets (like TWT and LPAQA) and meticulous reporting of prompt construction methods [42, 73].

4.4.2 Advanced Prompting and Training Strategies. To address the shortcomings of static, fixed-prompt evaluation, recent research has introduced a range of advanced prompting and adaptation strategies. These approaches—including adaptive, analytic, Bayesian, self-training, incremental, and distillation-based methods—seek to enhance both the robustness of model reasoning and the efficiency of knowledge extraction [2, 20, 57, 66, 76, 91, 95, 96, 108].

Adaptive frameworks, exemplified by the Adaptive-Solver (AS), dynamically adjust not only the prompt structure but also the underlying model selection, sampling routines, and decomposition strategies according to real-time reliability signals such as intra-prompt answer consistency [20, 108]. This paradigm moves toward more human-like, flexible reasoning by modulating model capacity and reasoning depth in response to uncertainty or complexity. Consequently, AS can selectively increase computational effort for more difficult problems while maintaining efficiency on easier tasks, achieving dual improvements in both accuracy and resource utilization that are unattainable via static prompting [20, 108]. Ablation studies further demonstrate that jointly optimizing multiple axes of adaptation (prompt structure, model parameters, sample size, and

decomposition approach) leads to synergistic gains, suggesting a widely applicable template for scalable reasoning in heterogeneous domains.

In parallel, reinforcement learning (RL) and self-training have proven effective at optimizing reasoning strategies end-to-end. For instance, the DeepSeek-R1 families employ reward-driven RL—augmented with curated Chain-of-Thought (CoT) examples—to encourage accurate and interpretable reasoning, outperforming standard supervised fine-tuning, particularly when these improvements are distilled into smaller, compute-efficient models [2, 20, 96]. However, direct application of RL—especially with smaller architectures or uncurated starting datasets—remains vulnerable to stability issues and incoherent outputs; reward shaping also necessitates careful design to circumvent hackable or narrowly optimized behaviors [2, 20].

Self-correction mechanisms, wherein LLMs iteratively refine their outputs based on automated feedback (either self-generated or from peer models), further enhance factual consistency and mitigate hallucinations, often without human supervision [66]. The efficacy of these strategies relies heavily on the diversity and informativeness of feedback, the timing of feedback integration (training, inference, or post hoc), and the baseline model’s intrinsic self-improvement capabilities.

Incremental and curriculum-based training strategies, such as multi-stage vocabulary expansion and progressive data distillation, also deliver marked improvements for both generative and discriminative tasks across pre-trained models [95, 96]. Importantly, such strategies frequently have a positive interplay with prompt-based evaluation: as foundational model competencies grow, prompting algorithms—whether static or adaptive—elicit more reliable and informative reasoning trajectories.

As summarized in Table 6, each advanced strategy carries unique advantages and corresponding challenges, reinforcing the necessity of tailored solution designs and rigorous evaluation.

4.4.3 Domain-Focused Evaluation and Transparency. Domain-specific analyses, particularly in biomedical and clinical contexts, underscore the necessity of robust evaluation methodologies and meticulous transparency in reporting. Comparative studies indicate a recurring performance dichotomy between general-purpose LLMs and specialized, fine-tuned models (such as BioBERT, PubMedBERT, BART): while closed-source models like GPT-4 achieve state-of-the-art results on open-domain reasoning and medical question answering tasks, they are consistently outperformed by specialized models in extraction and classification tasks [16, 100]. For instance, fine-tuned models outperform LLMs in macro-average scores (e.g., 0.65 versus 0.51), with substantial leads in entity recognition (e.g., NCBI Disease F1 = 0.909 for BioBERT versus approximately 0.6 for LLMs) [16]. Conversely, for reasoning-centric benchmarks such as

Table 6: Comparison of Advanced Prompting and Adaptation Strategies

| Strategy | Key Mechanism | Strengths and Caveats |
|-----------------------------------|--|---|
| Adaptive Prompting (e.g., AS) | Modulates prompts, model selection, and decomposition in response to reliability metrics | Improves efficiency and accuracy for complex tasks; requires real-time uncertainty estimation and robust control mechanisms |
| Reinforcement Learning (RL) | Optimizes reasoning via reward-driven feedback and curated examples (e.g., CoT) | Fosters interpretable and high-quality reasoning; susceptible to instability and reward hacking if not carefully managed |
| Self-Correction | Automated iterative refinement based on model or peer feedback | Reduces factual errors and hallucinations; effectiveness depends on quality of feedback signals and integration timing |
| Incremental / Curriculum Training | Progressive growth of vocabulary and staged data exposure | Enhances foundational competencies for more consistent downstream prompting; scalability and domain adaptation require thoughtful curriculum design |

medical licensure exams (e.g., MedQA), closed LLMs like GPT-4 surpass domain-specific SOTA (accuracy: GPT-4 at 0.72 versus SOTA at 0.42), although this superiority comes at a significant computational cost—up to 60-100 times that of smaller models [16]. Open-source LLMs, often benefitting more from broad instruction-optimized data than domain-specific pretraining, must undergo additional fine-tuning to approach these benchmarks [16]. In tasks involving text generation, GPT-4 and GPT-3.5 generate outputs considered more readable but less complete than BART [16].

Dynamic prompting strategies (such as few-shot chain-of-thought and instruction-based tuning) can alleviate issues like inconsistency and hallucination to some extent. However, these approaches have not fully eliminated gaps in output quality, particularly for open-source and zero-shot models, which still exhibit high rates of hallucination, omissions, and inconsistency across various clinical tasks [16]. Even highly instruction-tuned models such as MMedIns-Llama 3, developed with expansive, medically oriented instruction datasets (e.g., MedS-Ins), set new standards for information extraction, classification, text summarization, and diagnosis prediction, yet still struggle with comprehensive clinical scenario coverage, multilingual application, and real-world clinical validation [100]. Notably, improvements in NER, summarization, and classification (macro-F1 up to 86.66) are evident, but residual limitations underscore the need for continued innovation [100].

Transparency is now recognized as foundational for progress and reproducibility. Critical elements include the release of robust benchmarks (such as the Two Word Test for compositionality [73]), open access to datasets (e.g., TWT, LPAQA), public availability of evaluation code and models, and adherence to rigorous, standardized evaluation protocols [16, 42, 73, 100]. Studies have highlighted that many LLMs, while scoring highly on complex tasks, still fail on fundamental semantic judgments and compositional understanding [73], with prompt sensitivity, suboptimal query designs, and domain-specific nuances influencing outcomes [42].

In summary, the current trajectory of knowledge measurement and reasoning in LLMs is defined by the interplay of systematic prompt engineering, iterative training, and self-correction paradigms, supported by transparent, task-appropriate evaluation. Yet, further efforts are required to address critical challenges around prompt sensitivity, adaptation robustness, domain intricacies, and above all, reproducibility, in order to reliably advance LLM generalizability and interpretability.

5 Neural, Symbolic, Hybrid, and Graph-Based Reasoning

This section aims to clearly define the scope and objectives surrounding neural, symbolic, hybrid, and graph-based reasoning methodologies within the context of benchmarking AI reasoning capabilities. We target researchers and practitioners interested in

understanding methodological distinctions, strengths, and the operational suitability of these reasoning paradigms, with the goal of informing robust benchmark selection, model design, and empirical evaluation strategies. Our objectives are: (1) to describe the core principles and operational mechanics of each reasoning approach, (2) to contextualize their relevance within the broader landscape of reasoning benchmarks, and (3) to provide actionable insights and guiding questions for empirical researchers pursuing enhanced robustness or interpretability in reasoning system evaluation.

We underscore that this survey distinguishes itself from prior work by systematically mapping the interplay between benchmark types and reasoning methodologies, and by distilling concrete, measurable goals for advancing empirical research on reasoning evaluation. Specifically, we emphasize the need for rigorous comparison frameworks, nuanced prompt design for neural and hybrid models, and criteria for evaluating graph-based and symbolic components.

To aid practitioners, we recommend that empirical researchers consider the following when designing studies or evaluating systems: Which reasoning paradigm aligns best with the targeted benchmark’s complexity and transparency requirements? How do hybrid models balance interpretability and performance? What evaluation protocols are most appropriate for distinguishing methodological contributions in graph-based versus purely neural or symbolic reasoning tasks? Addressing these questions assists in the development and assessment of models suited to diverse application settings.

Open challenges and research gaps, such as the alignment of benchmarking protocols with real-world reasoning demands, the interplay between symbolic abstraction and neural generalization, and the scalability of hybrid approaches, remain. These gaps hold significant implications for the transferability, explainability, and reliability of AI reasoning systems, directly impacting their adoption and effectiveness in practical scenarios. By explicitly surfacing these issues, we encourage targeted research that advances both methodological rigor and empirical relevance in reasoning evaluation.

In summary, this section serves as both a synthesis and a practical guide, facilitating accessibility for new entrants while supporting advanced readers in making nuanced methodological and empirical decisions in the rapidly evolving field of AI reasoning.

5.1 Neuro-symbolic and Hybrid Frameworks

Recent advancements in artificial intelligence reasoning have underscored a marked convergence toward hybrid and neuro-symbolic architectures, aiming to harness the complementary strengths inherent in sub-symbolic (neural) and symbolic paradigms. Traditional neural models excel at capturing statistical regularities and enable scalable pattern recognition; however, they have historically

struggled with tasks necessitating principled structured reasoning—particularly those requiring compositionality, logical inference, or interpretability. In contrast, purely symbolic approaches offer transparency and verifiable reasoning but frequently lack the flexibility and robustness associated with data-driven learning. Hybrid and, more specifically, neuro-symbolic reasoning networks are designed to address these respective shortcomings through the integration of logic-based modules and constraint optimization strategies within neural network frameworks. This facilitates the embedding of explicit domain knowledge, enhances interpretability, and supports compositional inference [2, 11, 32, 38, 66, 69, 76, 90, 91, 95, 99, 105].

The primary methodologies in this field operationalize symbolic knowledge through logical constraints, differentiable logic operators, or explicit rule sets, strategically integrated with neural representations. Notably, Neural Reasoning Networks (NRNs) employ differentiable logical operations—including continuous (relaxed) analogs of Boolean ‘And’ and ‘Or’—to enable gradient-based learning mechanisms while simultaneously producing concise, human-interpretable explanations for tabular predictions [11]. Evidence suggests these architectures match state-of-the-art gradient-boosted tree models in predictive performance, yet offer significantly more compact and accurate reasoning chains. This underscores essential trade-offs between model compactness, logical transparency, and predictive capability [11, 95].

Hybrid constructionist paradigms for language understanding exemplify the application of neural heuristics to guide symbolic search over grammatical constructions. This approach outperforms traditional techniques in both computational efficiency and scalability, facilitating expressive neuro-symbolic language processing over large symbolic spaces [69].

Furthermore, recent hybrid frameworks enrich integration by incorporating algorithmic and graph-based components. Neural architectures inspired by algorithmic paradigms—such as dynamic programming or classical search procedures—can encode deep combinatorial structure and procedural logic within trainable models [2, 54]. Some hybrid systems dynamically adjust the depth of integration, balancing end-to-end learnability with the preservation of tractable symbolic intermediate representations. For example, deep reasoning networks (DRNets) synergize deep neural architectures with the explicit encoding of domain knowledge—in the form of thermodynamic rules—for robust phase identification in materials science [32]. Such integration achieves high predictive accuracy on structured scientific data while rendering latent model representations interpretable and closely aligned with domain priors [11, 32].

Despite these advances, several challenges persist: The majority of integration strategies are highly domain-specific, with manual specification of symbolic components underlying limited scalability and generalization. Automated methods for rule induction or the bootstrapping of symbolic modules with large foundation models are nascent and insufficiently robust [32, 69, 76, 90, 95]. There remains a fundamental trade-off between the expressiveness provided by symbolic representations and the differentiability required for effective neural learning. Nevertheless, hybrid frameworks have demonstrated particular promise in mathematical, scientific, and decision-critical domains [32, 38, 54, 66, 69, 76, 91, 95, 99, 105],

though open research problems include compositional generalization, recursive reasoning, and efficient knowledge acquisition under resource constraints.

5.2 Graph-Based and Domain Applications

Graph-based reasoning architectures have become crucial for enabling structured inference in both general and domain-specific contexts, particularly in synergy with recent progress in large language models (LLMs). The combination of graph neural networks (GNNs) and LLMs has yielded significant benefits for tasks requiring the synthesis of unstructured and structured knowledge, such as knowledge graph completion, scientific question answering, and reasoning over biomedical ontologies [18, 23, 26, 27, 30, 44, 52, 56, 79, 87, 89, 95, 102, 107]. These architectures encode structured information (e.g., knowledge graphs or tabular data) as graph representations, facilitating fine-grained reasoning by way of message passing, aggregation, and selective propagation, while leveraging the extensive contextual knowledge inherent in LLMs. As a prominent example, LBR-GNN fuses contextualized linguistic and graph representations, utilizing edge aggregation and targeted message passing to enhance common-sense question answering beyond the capacity of individual paradigms [102]. Additionally, frameworks that align multi-modal and textual data through chain-of-thought demonstrations have enabled complex scientific reasoning by jointly leveraging neural and structured elements [56].

In scientific, mathematical, and biomedical domains, these methods provide powerful mechanisms for encoding domain constraints, probabilistic relations, and hierarchically organized knowledge—attributes critical for reliable inference and interpretability [18, 23, 26, 27, 30, 44, 52, 56, 79, 87, 89, 95, 102, 107]. For combinatorially challenging problems, such as mathematical theorem proving, molecular property prediction, or scientific discovery, the fusion of neural networks with symbolic and probabilistic reasoning confers considerable performance enhancements paired with interpretable modeling [18, 23, 26, 27, 30, 44, 79, 87, 89].

In the biomedical field, in particular, the application of graph-based, symbolic, and hybrid reasoning methods has yielded tangible real-world impact. Approaches such as LLMs augmented with domain-specific symbolic and graph-based modules have proven superior to generic LLMs for tasks including extraction of social determinants from electronic health records (EHRs), clinical text classification, diagnosis assignment, and information extraction [10, 14, 16, 25, 27, 32, 34, 35, 43, 52, 54, 60, 64, 74, 76, 87, 94, 95, 100, 107, 108]. For example, enhancements through structured knowledge codes lead to improved detection of adverse social determinants and show reductions in demographic bias [34, 35, 60, 76]. Diagnostic frameworks leveraging these methods excel in DRG classification, rare disease recognition, and clinical narrative interpretation, providing interpretable rationales that facilitate actionable clinical insights [16, 32, 34, 52, 54, 64, 74, 87, 100, 107, 108].

Nevertheless, integrating graph-based, symbolic, and neural methodologies presents significant challenges: Scaling GNNs to handle massive, evolving knowledge graphs remains non-trivial. Managing compounded errors or hallucinations at the neural-symbolic interface is difficult. Automated construction of high-fidelity graph structures from noisy or heterogeneous data sources is an ongoing

Table 7: Representative Applications of Hybrid Graph-Based Reasoning Architectures

| Application Domain | Task or Use Case | Key Hybrid Approach |
|--------------------------------|--|---|
| Biomedical Informatics | Social determinants of health extraction, clinical text classification, rare disease detection | GNN-augmented LLMs, symbolic reasoning with domain codes, multi-modal graph reasoning |
| Materials Science | Crystal-structure phase mapping, materials discovery | Deep reasoning networks (DRNets) integrating neural and explicit domain constraints |
| Scientific Knowledge Synthesis | Scientific question answering, knowledge graph completion | Multi-modal alignment of LLMs and GNNs with chain-of-thought prompting |
| Mathematics | Theorem proving, mathematical property prediction | Hybrid symbolic-neural models leveraging procedural logic and graph representations |

obstacle. Biomedical and scientific fields are further challenged by limited annotated data, incomplete or inconsistent ontologies, and bias within domain corpora, all of which impair generalizability and trustworthiness [16, 27, 32, 34, 64, 102, 107]. Progress has been made through advancements such as standardized benchmarks for knowledge graphs, robust data augmentation, and instruction-tuned LLMs adapted to clinical and scientific content. However, achieving scalable, reliable, and fully explainable graph-based reasoning in practical applications remains contingent on continued methodological and theoretical innovations [10, 14, 16, 18, 25, 27, 35, 44, 52, 56, 60, 74, 87, 94, 95, 100, 102, 107].

6 Evaluation Methodologies, Interpretability, and Transparency

As we transition into this section, it is important to briefly reiterate the survey’s overarching aims: to systematically examine current approaches and challenges in the evaluation, interpretability, and transparency of AI systems, and to highlight methodological trends and future priorities. This section builds upon the foundational discussions in preceding sections by analyzing how evaluation frameworks support or hinder progress toward accountable, reliable, and explainable machine intelligence.

This section is organized to first describe foundational and emerging evaluation methodologies, then move to key interpretability strategies, and finally examine mechanisms for fostering transparency. Each subsection concludes with a brief summary to reinforce core insights and support seamless transitions as we consider the interplay among these central concerns.

6.1 Advanced Assessment and Reproducibility Metrics

In the rapidly evolving landscape of large language models (LLMs), robust and comprehensive evaluation methodologies are essential for meaningful assessment and responsible deployment. Traditional automatic metrics—such as ROUGE and BLEU—have long been standard, yet they demonstrate substantial misalignment with end-user utility, particularly in nuanced application domains like medical text simplification and summarization. Here, human comprehension, informativeness, and faithfulness are paramount requirements [16, 28, 39, 47, 81]. Empirical studies comparing human and automated ratings reveal that surface-level automated scores (e.g., ROUGE, BLEU) exhibit weak, if any, correlation with actual understanding or task utility, especially for lay audiences or within high-stakes clinical contexts [16, 36, 39, 68, 100, 103]. For example, large-scale evaluations of LLM-generated plain language summaries in medical settings demonstrate that while such outputs may score high on automated and even subjective metrics, they often yield lower comprehension outcomes when assessed through objective measures, such as multiple-choice tests or recall tasks [36]. This

discrepancy emphasizes the importance of focusing on downstream impacts, such as actionable understanding and decision support, rather than relying solely on surface-level similarities [16, 36].

Faithfulness and informativeness have thus become critical focal points for evaluation. Faithfulness, defined as the veracity of model outputs relative to the source data, remains challenging due to persistent risks of hallucination and error propagation [39, 47, 81, 103]. Recent work suggests the integration of multi-faceted evaluation strategies, including question-answering-based metrics, semantic similarity scoring, and rigorous human-in-the-loop assessments. These methods aim to prioritize objective comprehension and trust calibration over surface agreement alone [16, 68, 103]. At the same time, reproducibility has emerged as a core methodological concern in LLM and deep learning research. Issues such as heterogeneous experimental designs, lack of transparency in code and data, and environment-specific dependencies are widespread, complicating reliable replication [28, 47, 100]. To address these, prevailing guidelines urge replicability of computational environments, provision of detailed model and pipeline documentation, sharing of datasets and code in open repositories, and systematic sensitivity analyses that collectively bolster scientific reliability and progress [28, 39, 47].

Benchmark design has also attracted scrutiny concerning both efficiency and rigor. Notably, studies such as [68] have demonstrated that reducing the number of evaluation examples, when done judiciously, can preserve reliability and dramatically lower computational and environmental costs. Metrics like Decision Impact on Reliability (DIORe) offer quantitative frameworks to assess the impact of various benchmark design choices, underscoring that more examples do not necessarily equate to better reliability, and that aggregation and scenario diversity require careful consideration. Furthermore, limitations of static benchmarks, particularly their inability to capture dynamic, interactive, or real-world reasoning abilities of advanced models, have prompted calls for more dynamic and robust evaluation protocols [16, 28, 68, 81, 103]. Recent platforms promote extensibility, extensible tasks, determinism, and transparent open-source leaderboards, supporting community-driven continuous evaluation [103].

As outlined in Table 8, a balanced combination of evaluation methodologies is imperative to meaningfully assess LLM performance across different contexts.

6.2 Interpretability and Explanation Systems

Interpretability and transparency of LLMs remain central technical and ethical challenges, fundamentally underpinning accountability, auditability, and the cultivation of societal trust in AI systems [3, 10, 12, 17, 25, 32, 33, 35, 38, 46, 51, 64, 67, 76, 82, 89, 94, 95, 107]. Recent research explores a spectrum of explanation mechanisms, spanning symbolic and rule-based paradigms to extractive and abstractive

Table 8: Comparison of Model Evaluation Approaches: Key Criteria

| Evaluation Type | Strengths | Limitations | Use Cases |
|---|---|---|--|
| Automated Metrics (e.g., ROUGE, BLEU) | Fast; scalable; domain-independent | Poor correlation with human comprehension; insensitive to deep errors | Large-scale, low-stakes screening |
| Human-In-The-Loop | Captures comprehension and faithfulness; task relevance | Labor-intensive; subject to inter-rater variability | High-stakes, clinical, or legal assessment |
| Question-Answering/ Semantic Reproducibility Audits | Measures informativeness; supports factuality | Setup complexity; may require domain adaptation | Summarization, knowledge-grounded tasks |
| | Ensures reliability and scientific validity | Resource intensive; environmental dependencies | Benchmarking, regulatory review |

rationales. Each approach offers distinct strengths and faces unique trade-offs.

Symbolic frameworks, such as precedent-based constraint mechanisms and neural-symbolic integration, aspire to ground model outputs in transparent, human-interpretable rules and logic, explicitly operationalizing decisions through formal inference patterns [3, 10, 17, 25, 32]. These methods provide strong theoretical foundations in high-stakes domains (e.g., law, science) by fostering systematic reasoning, explicit auditing, and even formal proof generation. However, they frequently encounter challenges regarding scalability and adaptability when presented with high-dimensional or noisy real-world data [3, 12, 25, 46, 89].

In contrast, extractive and abstractive explanation systems draw upon features learned by deep architectures to expose underlying reasoning pathways. These approaches produce rationales that may be evaluated for logic, consistency, and alignment with expert understanding [12, 32, 38, 51, 76, 82, 94, 95]. Notably, empirical analysis of advanced LLMs (e.g., GPT-4) has demonstrated the potential for models to convincingly simulate complex domain-specific reasoning, such as clinical differential diagnosis. For instance, studies in the medical domain have shown that GPT-4 can generate rationales mimicking clinical reasoning formats, and the presence or absence of logical errors in these rationales often correlates with correctness: incorrect responses typically exhibit logical flaws, supporting the use of rationale quality as a practical signal for model oversight [38, 76]. Despite these advances, the fidelity of such model-generated explanations remains controversial, as rationales may reflect learned plausible justifications rather than actual model-internal processes [33, 64, 94].

To enable interpretability beyond post-hoc justification, contemporary methods have begun to embed explanation mechanisms directly within model training and input representations. Techniques such as hierarchical clustering, probing classifiers, and feature learning frameworks facilitate attribution of outputs to specific input features or groups, supporting both local (instance/case-specific) and global (class/cluster-level) interpretation [3, 46, 67, 107]. Probing, for example, trains auxiliary classifiers on latent representations to uncover which linguistic or structural properties are encoded [3], though such methods have inherent limitations in their ability to reveal causal mechanisms. Neural symbolic computing (NeSy) further attempts to integrate deep learning’s representational capability with symbolic AI’s logical structure and auditability, showing promising outcomes in domains such as mathematics, scientific discovery, and decision making. Nevertheless, NeSy faces ongoing challenges, including compositional generalization, scalability to complex tasks, and automated symbolic knowledge acquisition [10, 17, 33, 82, 95].

Interpretability in unsupervised tasks—such as clustering or feature extraction—poses unique obstacles due to the lack of ground-truth labels. The introduction of neuralized clustering models enables efficient, feature-level attribution of cluster assignments, providing explanatory insight even for unsupervised predictions [46]. Mutual information-based hierarchical clustering enhances both clustering performance and interpretability by maximizing the separation of learned groups [51]. These methods allow explanations about why data points are grouped together and support assessment of cluster quality. Nevertheless, a persistent concern is the discrepancy between model-produced explanations and user expectations, particularly when explanation style, length, or asserted confidence diverge from true model certainty. This mismatch can foster miscalibrated trust, as users may overestimate correctness based on surface characteristics of the explanation rather than underlying confidence or factual accuracy [36, 82, 95].

6.3 Bias, Fairness, and Auditing

Equitable and transparent deployment of LLMs critically depends on rigorous auditing for bias, fairness, and inclusivity, alongside proactive measures to minimize privacy and security risks [3, 8, 17, 34, 35, 38, 45, 48, 64, 67, 72, 76, 89, 90, 94, 95, 105, 107]. LLMs and other deep models are susceptible to learning and amplifying latent social and dataset-derived biases—risking the exacerbation of disparities in sensitive domains such as healthcare, law, and social services [8, 34, 38, 45, 48, 64, 67, 72, 90, 94, 95, 105]. Systematic audits employing model prediction analysis, confidence calibration, and demographic impact assessments have documented failures in both traditional and novel architectures, including increased sensitivity to demographic descriptor variables and uneven accuracy across groups [34, 45, 90, 95]. For instance, fine-tuned models addressing social determinants of health attenuated (but did not eliminate) bias compared to zero- or few-shot LLMs, indicating the need for both data- and architecture-driven mitigation strategies [8, 90].

Transparency throughout the modeling pipeline—including dataset composition, model objective specification, and parameter sharing—remains a prerequisite for detecting and mitigating such risks [3, 17, 72, 89, 107]. Contemporary literature increasingly calls for: open and representative datasets, public code and evaluation resources, and transparent evaluation protocols, to facilitate robust, community-driven audits and reproducibility [17, 45, 72, 76, 95, 107]. In parallel, transparency within modeling workflows—including visibility into intermediate representations, decision rationales, and potential failure points—is essential for regulatory oversight and informed engagement by diverse stakeholders [3, 32, 67, 72, 89, 95].

Mitigating hallucination and misinformation necessitates coupled strategies: technical interventions (such as factual verification

modules or knowledge-grounded models) and organizational safeguards (including red-teaming, continual post-deployment monitoring, and unambiguous user communication) [38, 48, 64, 72, 94, 105]. Furthermore, privacy and security considerations accentuate the importance of open, auditable, and securely managed data practices—especially in high-impact environments like medicine and law [3, 34, 72, 89, 105, 107]. Despite progress, ongoing gaps demand further attention, including the development of truly representative training corpora, robust adversarial testing procedures, and longitudinal audits to monitor emergent risks and behaviors throughout the model lifecycle [8, 17, 48, 72, 105].

In summary, the convergence of advanced assessment methodologies, interpretability frameworks, and bias/fairness auditing is transforming evaluation protocols for LLMs. The field is moving decisively away from narrow, surface-based metrics in favor of comprehensive, reproducible, and ethically attuned approaches that: integrate diverse stakeholder perspectives, foster open scientific practices, and directly confront the central risks and opportunities inherent in contemporary language modeling. [3, 10, 16, 17, 28, 32, 33, 36, 38, 39, 46, 47, 51, 64, 67, 68, 72, 76, 81, 82, 89, 94, 95, 100, 103, 107].

7 Reproducibility, Replicability, and Open Science

This section aims to explore the foundational concepts of reproducibility, replicability, and open science within the context of AI research, linking them back to our overarching survey objective: highlighting current limitations and best practices in ensuring robust scientific progress in the field. Throughout, we assess the methodological challenges, track the evolution of open science initiatives, and identify targeted open questions that remain in practice.

Reproducibility refers to the ability of independent researchers to obtain the same results using the original author’s data and code, while replicability addresses whether the same findings can be achieved with new data or alternative implementations. Open science serves as an enabling paradigm, promoting transparency, resource sharing, and community-driven validation. Together, these elements underpin the credibility and acceleration of research outputs in AI.

Despite increasing recognition of their importance, significant gaps persist. For example, issues related to incomplete dataset or code release, ambiguous experiment documentation, and varied computational environments often challenge both reproducibility and replicability. An explicit synthesis shows that many studies focus more on methodological innovation than on comprehensive reporting or resource availability, which hinders reproducibility at scale. Furthermore, the diversity of evaluation protocols and benchmarks leads to inconsistent results across studies, amplifying the need for standardized practices.

A key open question remains: which incentives or infrastructure most effectively promote routine, meaningful sharing of code and data in domains where privacy, ethics, or proprietary constraints are prevalent? Additionally, how can the community prioritize methodological standardization without stifling innovation, especially in rapidly evolving subfields? Addressing these concerns

remains critical for realizing the fullest potential of open science in AI.

In summary, advancing reproducibility, replicability, and open science will rely on deeper community commitment, robust technical solutions, and broader policy initiatives. The challenges identified here motivate the need for both more comprehensive empirical analyses and sustained discourse across the AI research landscape.

7.1 Reproducible Research Challenges

Despite rapid advances in foundational AI research, reproducibility in language model development—and in machine learning more broadly—remains a persistent obstacle, undermining both scientific rigor and field-wide progress. A central challenge is the ambiguous attribution of observed performance gains: recent studies reveal that when leading architectures such as BERT, ELMo, and GPT-1 are compared under harmonized experimental conditions, previously reported superiority of BERT often diminishes or vanishes altogether. This empirical ambiguity underscores the importance of principled ablation studies and controlled comparative experiments, as conflation of architectural, data, and optimization factors can obscure genuine innovations in model design, impeding reproducibility and interpretability in published research [65].

Broader issues compound these methodological deficits. Research protocols are frequently under-reported, code and data sharing remain inconsistent, and benchmarking practices are often heterogeneous. Such shortcomings impede direct replication, even for widely cited studies, as reproducibility audits continue to reveal deficits in both reporting and the accessibility of research artifacts [39, 65]. The crisis facing reproducibility is, therefore, not only technical but also cultural: while data sharing has increased, code dissemination is still sporadic, and in its absence, exact reproduction remains rare—an issue consistently observed across major venues and longitudinal analyses. Furthermore, impactful papers with verifiable and accessible code are more frequently cited, highlighting a direct benefit of transparency and openness for both community development and individual researchers [39].

Common failures in reproducibility extend beyond resource omission to encompass critical errors in code, incomplete statistical reporting, and insufficient experimental rigor, all of which undermine both peer review and public trust. Additionally, while definitions of “reproducibility” and “replicability” are well-established in the natural sciences, their inconsistent use within the machine learning literature leads to confusion and hampers empirical comparability [39]. Ultimately, a substantial proportion of published AI/ML research fails to meet the evolving standards of scientific rigor, with ad hoc practices prevailing in documentation, reporting, and procedural transparency.

7.2 Tools and Best Practices for Reproducibility

Robust reproducibility is increasingly undergirded by best practices and technological tools adapted from adjacent domains such as bioinformatics. At the experimental level, reproducibility is fostered through comprehensive documentation of data preprocessing steps, model specifications, and training protocols; statistical analyses of reproducibility, including sensitivity analyses and explicit tracking

of random seeds; and detailed reporting of all hyperparameters, code versions, and environmental dependencies [39].

These principles are realized through open science platforms—such as the IRIS and the Open Science Framework (OSF)—that facilitate the sharing of datasets, supplementary materials, workflow histories, and computational notebooks (notably Jupyter and R Markdown), as well as software environment capture via containerization [39].

Workflow management systems (WMS) are increasingly central, particularly in clinical and biomedical NLP. Systems like Snakemake, Galaxy, and Nextflow provide modular, version-controlled pipelines with provenance tracking, yielding transparent and auditable computational workflows [1, 4, 5, 7, 22, 37, 45, 49, 52, 58, 61, 63, 72, 85, 89]. The integration of standardized provenance mechanisms such as PROV ensures that workflows are not only repeatable but also interpretable across diverse contexts. Empirical assessments consistently demonstrate that WMS-based frameworks significantly outperform traditional monolithic pipelines in terms of traceability, standardization, and shareability, though technical challenges persist, particularly regarding comprehensive container support and seamless integration with public workflow repositories [72].

These distinctions are captured in Table 9, which summarizes comparative features of leading workflow management systems relevant to reproducible research.

Transparency initiatives continue to raise expectations for research documentation and open-source dissemination. The emergence of specialized automation tools—including arkit (for reproducible neuro-symbolic research) [9], MedS-Bench (standardized clinical evaluation) [100], and open NRN platforms (for explainable neural reasoning) [11]—illustrates the growing ecosystem of community-driven resources that enable reproducible benchmarking and democratize advanced reasoning tools [15, 29, 36, 53, 68, 71, 72, 78, 88, 102, 103]. These resources not only streamline benchmarking but also facilitate critical research and practical deployment by lowering entry barriers.

Formalization of reproducibility practices is evidenced by the adoption of guideline checklists, such as the CL Reproducibility Checklist for NLP conferences, which correlate strongly with both paper acceptance and community trust—particularly when tied to open code and dataset releases [39]. Other progressive frameworks emphasize protocol registration and systematic appendices; adherence to FAIR (Findable, Accessible, Interoperable, Reusable) principles; and explicit empirical validation of methods across diverse settings [29, 71, 78].

Implementation challenges remain prominent. Even as containerization and workflow modularity advance, sensitive data—especially in the clinical domain—often resists open sharing and necessitates solutions such as synthetic data generation, access-controlled repositories, and standardized metadata simulation [39]. Furthermore, the proliferation of benchmarking platforms (e.g., SUPERB, MedS-Bench, CL-MASR) highlights the need for unified, scalable, and statistically robust evaluation protocols that balance efficiency with breadth and scenario coverage [16, 47, 100].

7.3 Policy Recommendations and Incentives

Addressing the reproducibility crisis requires a dual approach that targets both procedural reform and incentive structures. Foremost is the need for explicit disambiguation of improvement sources in all published research, achieved through mandatory ablation studies, clearly reported experimental conditions, and rigorous benchmarking against well-tuned baselines [39, 65]. Such criteria should be embedded in journal and conference submission standards and underpinned by specialist review focused on statistics and experimental rigor.

Structural incentives are indispensable. Openness in benchmarking, code, and artifact sharing not only enables community evaluation but also fosters scientific accountability—an effect reflected in elevated citation rates and research impact for transparent publications [16, 39, 47, 100]. To this end, policy mechanisms including checklist-mandated artifact submission, embargoed yet verifiable code and dataset releases, and post-publication discussion platforms are recommended to support the systemic shift toward open scientific practice. Moreover, institutionalizing workflow-based repeatability—leveraging tools such as Snakemake and PROV—should become standard for all empirical studies, particularly those of significant societal consequence [1, 4, 5, 7, 9, 11, 15, 16, 22, 29, 36, 37, 39, 45, 47, 49, 52, 53, 58, 61, 63, 65, 68, 71, 72, 78, 85, 88, 89, 100, 102, 103].

Ultimately, a durable solution to reproducibility in AI and NLP research necessitates not only sophisticated computational infrastructure but also robust cultural and procedural transformations. Aligning incentives, rigorously upholding open and transparent standards, and cultivating a research environment that rewards both meticulousness and innovation together constitute the pathway toward resolving the current reproducibility crisis and ensuring continued scientific progress.

8 Safety, Robustness, Scalability, and Automated Pipelines

This section consolidates recent advances and core challenges in developing AI systems that are not only high-performing but also safe, robust, scalable, and amenable to automation across the life-cycle. The objectives here are threefold: (1) to clarify the distinct yet interrelated concepts of safety, robustness, scalability, and automation; (2) to synthesize key developments and open problems in each area with a focus on their practical integration; and (3) to provide a foundation for novel perspectives and taxonomies, highlighting pathways for future research and benchmark integration. The overarching goal is to guide researchers and practitioners in navigating the multidimensional tradeoffs and research frontiers at the intersection of technical assurance and real-world deployment.

At the conclusion of each technical domain covered—safety, robustness, scalability, or automated pipelines—critical open challenges and research questions are identified to orient ongoing work. Where appropriate, the section introduces refined conceptual distinctions and proposes a new taxonomy that distinguishes this survey’s synthesis from prior literature, emphasizing unique frameworks for understanding the confluence of these foundational topics.

Table 9: Comparative features of widely used workflow management systems supporting reproducible research.

| System | Modularity | Provenance Tracking | Container Support | Public Repository Integration |
|-----------|------------|---------------------|-------------------|-------------------------------|
| Snakemake | Yes | Yes | Partial | Limited |
| Galaxy | Yes | Yes | Yes | Yes |
| Nextflow | Yes | Yes | Yes | Yes |

8.1 Robustness and Adversarial Concerns

The deployment of large language models (LLMs) within high-stakes domains has accentuated persistent concerns regarding safety, robustness, and adversarial resilience. Despite substantial advances in reasoning capabilities and generalization, contemporary LLMs remain distinctly susceptible to a spectrum of adversarial threats. Among these, prompt-based jailbreaks, the emergence and misuse of unsafe model variants, and circumvention of built-in safeguards represent particularly acute vulnerabilities, exposing LLMs to malicious manipulation and the unintended generation of harmful content [29, 106]. Empirical analyses reveal that even commercial-grade LLMs equipped with advanced safeguard architectures can be undermined by universal jailbreak attacks, a finding that highlights intrinsic limitations in both proactive training regimes and post-hoc defense strategies [29]. In parallel, the proliferation of unaligned—at times intentionally adversarial—models such as dark LLMs increases opportunities for misuse, a risk that escalates as model access and training become increasingly democratized [29, 106].

To address these evolving adversarial threats, the research community has actively investigated out-of-distribution (OOD) detection methods, emphasizing frameworks based on generative adversarial networks (GANs) and autoencoders. These methods focus on pinpointing anomalous or untrusted inputs by capturing detailed features of the expected data distribution. Notably, techniques such as pseudo-OOD generation and latent space regularization have improved both the accuracy and area under the receiver operating characteristic (AUROC) for OOD detection without requiring exhaustive manual annotation of unsafe queries [29, 106]. For example, approaches that leverage GAN-regularized autoencoders to generate high-quality pseudo OOD utterances have demonstrated consistent improvements in OOD detection metrics, including AUROC and FPR95, across a range of dialogue datasets [106]. Nevertheless, the expressiveness constraints of generative models and incomplete representation of OOD scenarios in training data continue to hamper robustness, highlighting the need for more dynamic and scalable approaches capable of adapting as adversarial tactics evolve [106].

Another interconnected dimension of the safety discourse includes privacy, security, and fairness, which critically influence both open-source and proprietary LLM deployments [34, 35, 38, 45, 64, 67, 72, 76, 90, 95]. Privacy concerns encompass unintentional leakage of sensitive data in model outputs, susceptibility to model inversion attacks, and re-identification risks, particularly when LLMs are used on confidential health or financial data [38, 45, 76]. Security challenges such as prompt injection, model extraction, and exploitation of entrenched biases further complicate practical adoption and challenge public trust in LLM-powered systems [34, 67, 95].

The persistent challenge of fairness follows as LLMs may encode and perpetuate societal, racial, or gender biases, thereby amplifying inequities across domains such as healthcare, law, and finance [35, 64, 72, 90]. Comparative studies in sensitive domains, like health, have shown that domain-specific fine-tuning and the integration of synthetic multi-demographic datasets can help reduce demographic bias in predictions, but these improvements are incremental and require rigorous, ongoing audits, as well as transparent benchmarking for comprehensive mitigation [35, 45, 72].

In summary, the safety and robustness of LLMs depend not solely on model scale but on a comprehensive synthesis of adversarial evaluation, dynamic OOD detection, privacy-preserving mechanisms, and fairness-aware design, each of which must be regularly audited and transparently reported. Despite substantial research efforts, LLM safety and robustness remain locked in an adversarial dynamic, where defensive strategies must persistently adapt to match the pace and ingenuity of emergent threats [29, 64, 76, 106].

Key challenges highlighted in recent research include: achieving robust OOD detection under diverse threat models; ensuring privacy preservation during sensitive data handling; securing systems against injection, extraction, and misuse; and mitigating demographic and societal biases to promote fairness. Mitigation strategies increasingly emphasize model audits and transparent reporting, continuous updates of defense frameworks to remain responsive to new attack vectors, and the use of domain-specific as well as synthetic data augmentation to improve robustness and fairness.

8.2 Scalability, Workflow Orchestration, and Cost

The ongoing evolution of LLM architectures and reasoning strategies, while transformative, has sharply increased the requirement for scalable, efficient, and dependable deployment workflows. Managing orchestration across vast and heterogeneous data landscapes, as well as facilitating complex, multi-stage reasoning, necessitates robust automation, modular integration, and cost-efficient system design [7, 14, 20, 27, 30, 44, 50, 52, 54, 56, 62, 64, 68, 70, 91, 98, 101]. Prevailing workflow paradigms are broadly classified into three categories:

Within these paradigms, retrieval-augmented systems are particularly prominent in real-world deployments, selectively enriching LLM performance by supplying salient external knowledge. This is especially valuable for multi-modal tasks, where stratified retrieval and advanced reranking can elevate both task accuracy and resource efficiency, even under stringent computational constraints [14, 44, 101]. In parallel, reinforcement learning has emerged as a pivotal mechanism for optimizing multi-step workflows, including adapting to interactive or collaborative

Table 10: Representative paradigms for LLM workflow orchestration

| Paradigm | Core Methodology | Notable Advantages |
|---|--|--|
| Retrieval-based Orchestration | Dynamic incorporation of external factual or multimodal knowledge to augment context | Enhances reasoning fidelity; improves accuracy and efficiency, especially under resource constraints [44, 52, 54, 101] |
| Reinforcement Learning (RL)-Driven Optimization | Supervision via reward signals for procedural or multi-step reasoning and tool-augmented tasks | Adapts models to interactive, multi-agent, or sequential environments; increases flexibility and control [7, 27, 30, 50, 62, 70, 91, 98] |
| Automated Hierarchical Pipelines | Integration of operator modules and schedulers to choreograph complex, heterogeneous workflows | Facilitates modularity, scalability, and reliability; supports reproducibility [7, 50, 91, 101] |

scenarios such as tool-augmented reasoning and agent cooperation [7, 27, 30, 50, 62, 70, 91, 98]. Notably, the convergence of modular RL and LLM architectures with outcome-driven reward modeling streamlines deployment, particularly in cloud and distributed environments.

Scalable workflow orchestration at enterprise or population scale introduces further imperatives: cost-efficiency, accessibility, and environmental sustainability. These aspects shape both adoption and governance of LLM solutions [27, 56, 62, 64, 68]. Recent benchmarking initiatives, facilitated by efficient evaluation suites and adaptive model compression tools, demonstrate that meticulous pipeline optimization—including minimization of redundant computation, document signal refinement, and aggregation strategy tuning—can materially lower operational costs and carbon emissions with negligible detriment to performance [27, 56, 64, 68]. The widespread adoption of open-source, modular orchestration libraries further accelerates research reproducibility and expedites technology transfer into industrial and public-sector applications [62, 64, 91, 101].

Despite these advancements, important challenges persist. End-to-end automated pipelines remain prone to error propagation, OOD failures, and emergent behaviors as system complexity increases. Achieving a balance between efficiency, accessibility, and rigorous safety or fairness constraints thus demands systematic trade-off analyses and the standardization of auditing protocols across both research and production environments [44, 64, 68, 98]. As the adoption of LLMs accelerates, the continued development of scalable, automated, and cost-conscious orchestration frameworks represents a crucial determinant in unlocking—safely and equitably—the transformative societal potential of advanced AI.

Critical workflow considerations: Modular design for reliability and scalability; Dynamic retrieval and efficient context integration; RL-based adaptation for multi-step tasks and agent collaboration; Cost and resource optimization through automated benchmarking and pipeline tuning.

Ongoing risks: Error propagation across complex pipelines; OOD breakdowns and robustness gaps; Trade-off management between performance, cost, and safety.

9 Multi-Modal, Multi-View, Demographic Inclusion, and Biological Foundations

This section surveys the landscape of advancements in multi-modal and multi-view approaches, demographic inclusion strategies, and the biological foundations within artificial intelligence. The objective is to provide a comprehensive synthesis that not only summarizes key technical progress in these areas but also distinguishes this survey’s perspective by proposing a structured taxonomy to categorize the diversity of current methodologies. By explicitly articulating main challenges and open research questions, this section aims to guide readers through the complexities and evolving frontiers of these topics, ensuring alignment with the overarching

survey goal of identifying critical directions for future research and responsible deployment.

To provide a clear guide through this section, we begin with an overview and taxonomy of multi-modal and multi-view learning paradigms, then discuss demographic inclusion in model development, and finally review foundational biological inspirations in AI. Each subsection closes with a summary of major open research challenges and opportunities for synthesis across domains.

9.1 Multi-Modal and Multi-View Learning: Taxonomy and Advances

Multi-modal and multi-view learning techniques enable AI systems to integrate and reason over diverse types of inputs (e.g., text, image, speech, sensor data), allowing richer representation learning and improved generalization across tasks. We introduce a taxonomy that distinguishes methods based on the level and mechanism of fusion: early (input-level), intermediate (representation-level), and late (decision-level) fusion. Further, we categorize advances by their supervision schemes (supervised, self-supervised, weakly supervised) and compatibility with downstream tasks.

The field faces ongoing challenges, such as harmonizing information from modalities with divergent structures, dealing with noisy or missing views, and scaling fusion architectures efficiently. Despite progress, model interpretability and robust cross-modal transfer remain open questions.

Open research questions in multi-modal and multi-view learning include: How can semantic alignment across heterogeneous modalities be improved at scale, especially in resource-limited domains? Can unified representation spaces be made robust against missing or adversarial inputs? What metrics best capture the trade-off between expressivity, interpretability, and computational efficiency in real-world deployments?

9.2 Demographic Inclusion: Frameworks and Challenges

Ensuring demographic inclusion in AI models requires explicit strategies to mitigate bias, improve fairness, and achieve representative generalization, especially in sensitive applications such as healthcare and social platforms. We synthesize recent frameworks under the lenses of dataset auditing, algorithmic fairness constraints, and participatory model development, proposing a new taxonomy that distinguishes between proactive (pre-processing and data curation), reactive (in-processing during learning), and post-hoc (evaluation and correction) approaches.

Persistent technical questions remain, such as constructing datasets that meaningfully represent minority groups without exacerbating privacy or measurement issues, and effectively benchmarking fairness in multi-modal settings.

Open research questions for demographic inclusion include: What standardized procedures can ensure ongoing demographic auditability as models evolve? How can fairness criteria be extended to dynamic, multi-modal model pipelines, and what trade-offs emerge between accuracy and inclusion in this context?

9.3 Biological Foundations: Inspirations and Limitations

Biological systems have inspired numerous architectural and functional innovations in AI, from neural network topologies to learning rules. This subsection categorizes advances based on the granularity of biological inspiration—cellular (e.g., neuron models), circuitry (e.g., recurrent and feedback connections), and system-level motifs (e.g., attention, memory consolidation).

Challenges in this area include over-simplification of biological mechanisms, difficulties in transfer to large-scale artificial systems, and limited understanding of which biological priors most benefit AI learning.

Open research questions around biological foundations include: Which biologically inspired mechanisms provide consistent benefits across tasks, and how can empirical benchmarks be shaped to evaluate their contributions objectively? What are the integration pathways for iteratively refining AI algorithms with new biological discoveries, especially under computational resource constraints?

Summary of Section Objectives and Synthesis

This section has articulated a novel taxonomy across three key domains: (1) the fusion level and supervision of multi-modal learning methods, (2) proactive/reactive/post-hoc strategies in demographic inclusion, and (3) the granularity of biological inspirations driving architectural design in AI. By highlighting persistent open challenges at the intersection of these domains, we aim to guide future research toward integrated, robust solutions. As these fields continue to evolve, synthesizing insights across technical, ethical, and biological perspectives remains an essential frontier for the AI community.

9.4 Multimodal Fusion and Learning

The contemporary landscape of machine learning—particularly in critical fields such as healthcare and scientific reasoning—increasingly depends on the integration of information across multiple modalities and perspectives. Multimodal learning encompasses the fusion of heterogeneous data types, including audio, speech, emotion, and text. This approach leverages the complementary strengths of each data type to advance model robustness, enhance reasoning capabilities, and improve interpretability. Foundational frameworks underpinning this domain include co-training, autoencoder architectures, and contrastive fusion techniques, all of which have proven pivotal in harmonizing diverse data representations and boosting downstream performance on tasks such as speech and emotion recognition, clinical reasoning, and common-sense question answering [14, 18, 23, 26, 27, 30, 44, 52, 56, 79, 83, 87, 89, 95, 101, 102, 107].

There has been a marked evolution from naive modality concatenation toward more sophisticated cross-modal representation learning strategies. Techniques such as multi-view learning exploit both redundancy and complementarity among multiple sources or

perspectives, facilitating enhanced generalization and resilience to overfitting—challenges that are particularly pronounced in low-resource scenarios [101]. For instance, contrastive learning paradigms enable alignment between modalities by maximizing agreement within shared latent spaces, a principle driving recent advances in multi-view speech and language applications as well as cross-modal question answering [26, 101]. Autoencoder-based fusion mechanisms further reinforce integration, learning joint distributions over modalities and thereby supporting complex semantic reasoning and improved model interpretability [87, 89, 101].

Despite these architectural advancements, considerable challenges endure:

- Many multimodal models, such as large language models (LLMs) and agent-based frameworks, face persistent limitations in achieving genuine cross-modal reasoning, often exhibiting brittleness to distributional shifts and difficulties in fusing structured with unstructured data [23, 30, 83, 89, 95, 107].
- Benchmarking studies designed for multimodal and multi-view evaluation uncover notable performance inconsistencies attributable to both the design of fusion mechanisms and a tendency for models to overfit to the dominant modality in the training corpus [23, 26, 44, 56, 102].
- Explainability remains a fundamental concern: while advanced LLMs (e.g., GPT-4) can convincingly mimic clinical reasoning processes and offer ostensibly interpretable rationales, these rationales may not align with authentic multi-step or causal reasoning as executed by human experts, highlighting the ongoing need for principled, reasoning-aware architectures [14, 26, 27, 95, 107].

The emergence of contrastive and symbolic-neural fusion frameworks represents an important advance toward greater model accountability and transparency [18, 52, 56, 87, 89]. Equally, the integration of biological priors and neuroscientific insights is gaining traction. Recent work with connectome-inspired neural architectures suggests that biologically plausible modularity and critical network dynamics are capable of optimizing computational performance, pointing to a fruitful intersection between artificial learning models and human brain network topology [83]. Furthermore, neural-symbolic approaches, which merge statistical learning with formal logical reasoning, enhance both transparency and the robustness of decision-making across scientific, medical, and legal domains [18, 52, 87, 89, 95]. Nevertheless, the challenge of achieving scalable, interpretable, and consistently high-performing fusion across high-dimensional, multi-view, and structured-unstructured data streams remains central to ongoing research.

To offer a structured comparison of prominent multimodal fusion techniques and their primary benefits and limitations, see Table 11.

9.5 Inclusion, Ethics, and Demographic Representation

The equitable and ethically responsible deployment of AI systems necessitates sustained attention to dataset inclusivity, demographic fairness, and compliance with evolving regulatory standards. The risk of algorithmic bias—stemming from non-representative datasets,

Table 11: Comparison of Representative Multimodal Fusion Strategies

| Fusion Method | Key Strengths | Key Limitations |
|--------------------------|--|--|
| Naive Concatenation | Simplicity, ease of implementation | Limited interaction modeling; prone to overfitting dominant modalities |
| Multi-View Learning | Exploits complementarity and redundancy; effective in limited data scenarios | Requires careful view selection and alignment; moderate interpretability |
| Contrastive Fusion | Strong alignment of shared representations; improved robustness to noise | Sensitive to initialization/negative sampling; computational complexity |
| Autoencoder-based Fusion | Learns joint latent spaces; potential for enhanced interpretability | May struggle with complex cross-modal relationships; sensitivity to modality imbalance |
| Symbolic-Neural Fusion | Increased explainability; supports formal reasoning over data | Complexity in integrating symbolic/connectionist layers; often domain-specific |

model overfitting to majority subpopulations, or the omission of critical social determinants—carries profound real-world consequences, particularly within highly regulated domains such as healthcare, finance, and law [8, 34, 35, 38, 45, 48, 64, 67, 72, 76, 90, 95, 105].

Recent scholarship emphasizes the imperative for representative data collection protocols that capture the full spectrum of demographic and socio-economic variability observable in actual populations. As a notable example, structured electronic health record (EHR) codes are often inadequate for reporting social determinants of health, whereas advanced text-mining methods leveraging language models demonstrate improved recall of disparate factors, especially those relating to marginalized groups [34, 45]. The application of synthetic data augmentation and targeted fine-tuning for underrepresented classes has further reduced vulnerability to demographic bias, thus reinforcing the necessity for systematically balanced data pipelines in AI model development [35, 48, 72, 90].

Nevertheless, entrenched and emergent challenges remain: Algorithmic audits and benchmarking continue to reveal systematic disparities in model outputs along axes such as race, gender, and socio-economic status, exposing neglected failure modes and driving calls for more nuanced, intersectional evaluation protocols [34, 38, 64, 67, 76]. The lack of unified standards for evaluating LLMs, combined with a proliferation of ad hoc prompt engineering approaches, has impeded replicability and undermined confidence in observed advances in fairness [8, 48, 90, 105]. This replication crisis underscores an urgent need for robust experimental design, open data/code sharing, and reproducibility standards to accurately assess and rectify demographic risks.

In parallel, significant regulatory and ethical developments—including GDPR, the EU AI Act, and growing mandates for explainable AI—are shaping both technical design and evaluation practices [72, 95, 105]. Leading research advocates for the integration of fairness constraints, causal inference, and interpretability objectives directly into training and inference workflows, so that regulatory compliance is embedded as a foundational design principle rather than as a post hoc consideration [38, 45, 67, 72, 95]. Legal-theoretic formalisms and hybrid neuro-symbolic systems facilitate the encoding of precedential knowledge, offering promising directions for transparent and auditable AI in sensitive domains [34, 52, 89, 95].

In summary, advancing inclusion, ethics, and demographic representation in multi-modal, multi-view AI necessitates continuous cross-disciplinary engagement, methodological transparency, and a willingness to rigorously confront both the technical and socio-ethical complexities intrinsic to scalable real-world deployment.

10 Societal, Ethical, and Policy Considerations

This section critically examines the multifaceted societal, ethical, and policy questions arising from the development and deployment

of AI systems. The objectives of this section are to delineate the scope of key issues, clarify the challenges at the intersection of technology and society, and offer a foundational taxonomy for ongoing discourse. In doing so, we aim to provide readers with both a synthesis of the current landscape and a springboard for future research, distinguishing our survey by proposing a structured framework that integrates ethical, societal, and policy dimensions.

10.1 Overview and Scope

To facilitate navigation, this section surveys the principal challenges and considerations related to the societal impact, ethical deployment, and regulatory aspects of AI technologies. We articulate the interplay between these domains and offer a conceptual distinction between societal effects (e.g., equity, access), ethical predicaments (e.g., algorithmic bias, agency), and policy responses (e.g., governance, regulation).

Open Research Questions: How can frameworks for societal and ethical evaluation keep pace with the rapid evolution of AI technologies? What are the most effective mechanisms to translate policy intent into robust governance practices?

10.2 Societal Impact

AI systems carry significant implications for employment, accessibility, social inequality, and public trust. Existing work often addresses the distributional consequences and the potential for both exacerbating and alleviating disparities. In this survey, we propose a layered perspective that organizes societal impacts along axes such as economic sectors, affected populations, and the temporal horizon of effects, offering a clearer taxonomy than prior literature.

Transitional Note: While societal implications shape broad human contexts, ethical challenges frequently arise at the intersection of system design and human values.

Open Research Questions: How can future studies better evaluate the long-term, indirect societal impacts of AI? In what ways might systemic socioeconomic biases become entrenched or mitigated by AI applications?

10.3 Ethical Considerations

Core ethical issues include fairness, transparency, accountability, and respect for human agency. While past surveys often enumerate risks and mitigation strategies, our synthesis advocates for a nested conceptualization: ethical dimensions are positioned as mediating forces between technical design choices and emergent societal consequences.

Transitional Note: Moving from ethical evaluation to policy formulation, the focus shifts from identifying risks to crafting enforceable, adaptive frameworks.

Open Research Questions: How can ethical principles be operationalized concretely within technical design pipelines? What new ethical dilemmas might emerge with deepening human-AI collaboration?

10.4 Policy and Regulation

This subsection reviews regulatory approaches, emphasizing the dynamic interface between legal frameworks, industry standards, and international coordination. Prior literature typically segments policy analysis by jurisdiction or sector; in contrast, we introduce a comparative matrix organizing policy responses by regulatory trajectory (e.g., precautionary, permissive) and stakeholder scope (e.g., public, private, cross-sectoral).

Open Research Questions: What mechanisms best ensure policy responsiveness given the velocity of AI innovation? How can international policy harmonization balance local autonomy with global standards?

10.5 Summary and Future Directions

In summary, the societal, ethical, and policy domains present intertwined challenges that demand integrated analysis. By proposing a novel taxonomy and emphasizing the interplay between these dimensions, our survey uniquely frames the landscape for future research. Further progress hinges on cross-disciplinary innovation and continual reassessment of open research questions identified herein.

10.6 Oversight and Accountability

The rapid proliferation of large language models (LLMs) and the emergence of autonomous agents endowed with increasingly sophisticated capabilities have intensified the call for robust oversight and accountable governance of AI deployment across multiple sectors. These concerns are particularly salient in the context of models exhibiting autonomous replication and adaptation (ARA)—agents that can potentially acquire resources, adapt to novel environments, and self-replicate, thereby circumventing conventional operational boundaries and regulatory safeguards [19, 48, 85]. Although empirical investigations currently demonstrate that only the simplest forms of ARA are achievable, the swift pace of frontier model advancement, in conjunction with the modular design of tool-using agent frameworks, signals credible scenarios in which future iterations could attain robust, persistent autonomy—especially when coupled with scalable infrastructure and human facilitation [34, 48, 69, 85].

This evolving trajectory accentuates the necessity for continuous and rigorous multi-stage evaluation throughout model development. It is insufficient to rely exclusively on static performance benchmarks; comprehensive assessments must encompass dynamic, end-to-end, and adversarial evaluations that address exploitation, security, and risk scenarios [69, 85]. Prevailing evaluation regimes often limit analyses to simulated environments or controlled task specifications, yet such constraints systematically underestimate true risk due to the use of proxy measures, biases inherent in judge models, and an underappreciation of attack surface complexity [63, 69, 85]. Lessons from other high-impact

AI domains—including healthcare, finance, and critical infrastructure—reveal that rapid system advancements and escalating complexity often outstrip the establishment of comprehensive regulatory, ethical, and technical standards [8, 34, 67].

From a policy perspective, enduring barriers to reproducibility, transparency, and rigorous peer scrutiny pose significant challenges to societal trust and scientific integrity [12, 64, 84, 105]. Even within natural language processing, attempts to replicate empirical findings routinely expose methodological shortcomings, including insufficient reporting, flawed interface design, and ethical lapses [84]. These challenges are amplified in rapidly evolving or high-profile fields (e.g., deep learning, LLMs), where increased research popularity is paradoxically associated with diminished replicability—thereby complicating system auditability and accountability [12, 45]. Providing code or model weights alone proves inadequate without comprehensive documentation of computational environments and explicit data provenance [12, 45].

The task of balancing the robustness, scalability, efficiency, and resource demands of advanced AI models introduces inherent structural tensions between performance optimization and core societal values, such as transparency, safety, and equitable access [21, 63, 82, 89]. As dataset sizes and compute budgets escalate, empirical evidence demonstrates diminishing efficiency gains due to the saturation of informative data or resource constraints, raising critical concerns about long-term sustainability, environmental impact, and global access to AI technologies [21]. Accordingly, effective policy responses must integrate technical guidelines (e.g., mandatory documentation, interpretability reporting, rigorous stress testing under variable conditions) with legal and ethical instruments (such as explicit liability allocation, robust audit traceability, and comprehensive algorithmic impact assessment) [8, 34, 67].

The prospect of Artificial General Intelligence (AGI)—whether imminent or speculative—further intensifies scrutiny regarding the alignment of agent goals, operational mechanisms, and the broader public interest [34, 52, 64, 95]. Contrary to popular anxieties, contemporary research suggests that the more urgent risks emanate not from hypothetical AGI, but from the deployment and potential misregulation of extant, highly capable yet inherently limited AI models [34, 95]. Theories of goal-means correspondence and the dynamic reconfigurability of agent architectures offer potential pathways for ensuring alignment, but concomitantly introduce new risks—such as goal drift, emergent behaviors, and heightened oversight complexity [34, 52]. Without rigorous, cross-sectoral regulatory frameworks and ongoing ethical review, the opacity and adaptive capacity of advanced agents may ultimately jeopardize foundational principles of accountability, safety, and democratic governance [27, 34, 67, 69].

Key distinctions between oversight challenges and policy priorities among different AI contexts are summarized in Table 12.

10.7 Toward Human-Centric and Transparent AI Systems: A Conceptual Framework

Survey Objectives and Scope. This section synthesizes the survey's broader goals: to articulate technical and systemic barriers to trustworthy AI, identify actionable mitigation protocols, and

Table 12: Comparison of Oversight Challenges and Policy Priorities in AI Deployment

| Domain | Oversight Challenges | Policy and Technical Priorities |
|---|--|---|
| Autonomous Replicating Agents | Rapid system adaptation; bypass of traditional safeguards; expansion of attack surfaces | Dynamic evaluation; adversarial testing; continuous monitoring; liability frameworks; adaptation detection mechanisms |
| High-Impact Sectors (Healthcare, Finance, Infrastructure) | Accelerated complexity; lag in regulatory and ethical standards; reproducibility bottlenecks | Regulatory modernization; technical documentation standards; peer auditing; sector-specific ethical review |
| Frontier Model Research (LLMs, Deep Learning) | Difficulty in reproducibility; auditability gaps; popularity inversely correlated with replicability | Code and data disclosure; computational environment encapsulation; transparent benchmarking; data provenance tracking |
| Societal Alignment (AGI and near-term AI) | Goal misalignment; emergent risk; oversight complexity | Goal-means correspondence mechanisms; system alignment testing; cross-sectoral regulation and ethical review |

introduce a new taxonomy clarifying the interdependence of transparency, human factors, and accountability in Large Language Models (LLMs). Unlike prior surveys, the focus is on holistic human-LLM collaboration, concrete protocol exemplars, and a unified conceptual framework for human-centric AI outcomes [27, 67].

Proposed Framework for Trustworthy, Transparent AI. Achieving human-centered AI requires technical rigor embedded within systems intentionally architected for transparency, auditability, and collaborative engagement. The proposed conceptual framework, summarized in Table 13, is structured around three pillars: (1) transparent and calibrated model reasoning, (2) system-level explainability and auditability, and (3) institutional and policy protocols supporting the entire AI lifecycle.

Examples of Effective Protocols. The described approaches have seen successful implementation: confidence-matched explanations, as shown in behavioral experiments with GPT-4 and other models, have narrowed calibration and discrimination gaps, allowing users to trust outputs in line with actual model reliability [82]. In legal decision contexts, precedent-tracing frameworks and open-source tools enable direct auditability by establishing explicit mappings from outputs to concrete data or legal deductions [89]. In healthcare, clinical use of GPT-4 in critical care showed earlier, safer decision support compared to standard practice and allowed bias reflection and auditability, contingent upon expert oversight and rigorous protocols [34].

Current LLMs, while possessing impressive emergent abilities, face sustained challenges—such as hallucination, bias, and poor uncertainty calibration—that jeopardize societal value without dedicated human-centered design. A persistent gap remains between model confidence and user perception: for example, long explanations, regardless of accuracy, inflate user confidence in LLMs, diverging from the underlying statistical reliability. This underscores the demand for transparent communication of uncertainty; explanation styles and outputs should be matched to model calibration metrics and explicitly indicate true confidence [82]. Calibration-oriented design ensures user trust is aligned with model reliability, which is especially critical for decision-support in sensitive domains.

Recent advances in interpretability and auditability frameworks provide specific design pathways. Precedent-based interpretability, inspired by legal reasoning, now includes open-source order-theoretic implementations that trace model decisions to the structures of the training set, allowing both contestability and systematic audits [89]. Neural-symbolic (NeSy) systems bridge statistical inference with formal logic, yielding semantic explanations and facilitating corrective user interaction; while scalability remains a challenge, these directions are mature in legal, healthcare, and policy use cases [45, 67].

System-wide transparency and reproducibility are increasingly realized via ecosystem practices: open, standardized benchmarks—such as RepliBench for agentic LLM evaluation [8]—complement broad

calibration, comprehensive protocols, and automated, transparent reporting [8, 105]. Metrics are now scrutinized for wide confidence intervals and insufficient statistical improvements, motivating refined evaluation and reproducible research [105]. However, high-level system design must directly address interactive sources of bias and new error modes unique to hybrid human-LLM teams—for instance, clinician and LLM reasoning complementing one another but also propagating biases in clinical decision support [34, 67].

Actionable Recommendations and Systemic Shifts. Realizing human-centric, trustworthy AI demands not only improved technical solutions but also cultural and procedural reformation. Key shifts, supported by prior literature, include comprehensive pre-registration of research, mandatory specialist ethics review, open publication of evaluation datasets, and post-publication critique to sustain accountability [34, 67, 84].

In sum, this taxonomy-driven, protocol-oriented perspective clarifies how transparency, calibration, and human factors together represent a decisive advance over prior surveys [27, 67]. The actionable examples provided—in domains from legal AI to intensive-care clinical reasoning—demonstrate meaningful improvements in both trustworthiness and system auditability, with protocols and infrastructure not only enhancing technical robustness but also anchoring AI development in practices verifiably aligned with the public good [8, 27, 34, 69].

11 Persistent Gaps, Open Challenges, and Strategic Recommendations

Section Objectives: This section aims to (1) distill the key knowledge gaps surfaced throughout the survey, (2) codify open challenges that persist in the field, and (3) put forth actionable recommendations for future research and deployment. Our goal is to provide a clear framework that not only synthesizes the analysis but also guides diverse stakeholders toward impactful next steps.

11.1 Taxonomy of Gaps and Open Challenges

To clarify and structure ongoing issues in the field, Table 14 introduces a new taxonomy that groups persistent gaps and open challenges into four conceptual domains: Data, Models, Evaluation, and Deployment. Each domain is characterized by concrete challenges and representative illustrative examples, ensuring interdisciplinary accessibility in definitions.

11.2 Analytical Synthesis and Transition to Recommendations

The preceding taxonomy surfaces both long-standing and emergent challenges across data, models, evaluation, and deployment. Explicitly defining key terms facilitates comprehension among experts and non-specialists alike. These persistent issues collectively form

Table 13: Taxonomy of Human-Centric Transparency and Accountability in LLM Systems

| Pillar | Mechanisms | Example Implementations | Domain |
|-------------------------------|---|---|----------------------------------|
| Transparent Reasoning | Confidence alignment, uncertainty communication | Modified explanation styles reflecting true model certainty [82], expected calibration error (ECE) reporting [82] | Decision support, high-stakes AI |
| Explainability & Auditability | Precedent-based interpretability, neural-symbolic reasoning | A Fortiori case-based reasoning frameworks [89], open-source logical toolkits [89], NeSy for semantic logic bridging [45, 67] | Legal AI, healthcare, policy |
| Ecosystem Protocols | Benchmarks, evaluation protocols, transparent reporting | Open benchmarks (e.g., RepliBench [8]), confidence-calibrated reporting [105], post-publication monitoring [34, 84] | Model evaluation, clinical AI |

Table 14: Taxonomy of Persistent Gaps and Open Challenges

| Domain | Challenge | Description | Example | Key Terms Defined |
|------------|-------------------|---|--|--|
| Data | Limited Diversity | Insufficient coverage of real-world variation | Bias in benchmark sets | Data diversity: range of demographic, linguistic, and situational contexts represented |
| Model | Robustness | Fragility to adversarial or rare scenarios | Model failure in edge cases | Robustness: model performance stability under distribution shifts |
| Evaluation | Metric Alignment | Evaluation metrics poorly reflect end-user utility | Misalignment between automated scores and human satisfaction | Metric alignment: congruence between evaluation measures and real-world utility |
| Deployment | Scalability | Difficulty in translating models into production at scale | Bottlenecks in computation, cost, or oversight | Scalability: capacity to maintain function and quality as application scope increases |

the foundation for actionable recommendations designed to bridge the gap between analytic insight and practical implementation.

11.3 Strategic Recommendations

To address the identified gaps, we propose the following strategic recommendations, illustrated where possible with specific examples of protocols that have demonstrated success:

- 1. Promote Data Diversity:** Invest in the intentional curation of datasets that encapsulate a wide variety of real-world conditions, informed by collaborative collection protocols. As seen in leading consortia, deliberate engagement with stakeholders from underrepresented domains has resulted in improved data representativeness.

- 2. Strengthen Model Robustness:** Encourage robust model development through standardized adversarial testing and stress evaluation. Successful implementations involve periodic red-teaming and challenge sets adopted in high-stakes applications.

- 3. Advance Metric Alignment:** Refine evaluation metrics to better capture user benefit and real-world relevance, such as integrating end-user feedback into iterative metric recalibration. Pilot studies where human-in-the-loop assessment was institutionalized have shown greater metric validity.

- 4. Enable Scalable Deployment:** Design architectures and pipelines that accommodate resource scaling and operational oversight. For example, modular deployment strategies in open-source frameworks have facilitated reliable, scalable rollouts in production systems.

11.4 Summary: A Meaningful Shift in Recommendations

These proposals collectively represent a departure from template-based reviews by emphasizing actionable pathways grounded in a cross-domain taxonomy and concretized by implementation examples. By explicitly structuring and defining ongoing challenges, and mapping them to tailored recommendations, this framework offers a strategic foundation for advancing the field beyond the incremental refinements of previous surveys.

In brief, this section reaffirms our survey’s objective to bridge analytic depth with prescriptive utility and interdisciplinary clarity, providing both a roadmap for future research and a practical guide for practitioners.

11.5 Identification of Persistent Gaps

Despite substantial advances in large language models (LLMs) and their integration into diverse natural language processing (NLP) and artificial intelligence (AI) systems, several persistent gaps continue to impede both scientific understanding and practical deployment. These limitations are prominently observed in foundational domains, including semantic and structural evaluation, fairness and auditing, robustness, interpretability, and the realization of effective human-in-the-loop systems [2, 4–7, 9–11, 14–16, 18–20, 22, 24–28, 30, 34, 35, 37–45, 47, 49, 52, 53, 55, 56, 58, 59, 61–69, 71, 73, 75, 77, 78, 80, 81, 83, 85, 87, 89–93, 95–97, 100–107].

A recurring critical issue is the inadequacy of current benchmarking strategies. Most benchmarks lack comprehensive coverage for compositional and real-world reasoning and are insufficient in assessing capabilities such as abstraction, semantic faithfulness, and domain generalization. Evidence from recent studies demonstrates that LLMs still exhibit brittleness on logic puzzles, multi-step inference, and tasks requiring integration of world knowledge—domains in which human performance demonstrates compositional generalization and robust intuition [4, 6, 9, 10, 25, 26, 38, 42, 45, 49, 73, 97, 101]. Notably, as highlighted by the Two Word Test [73], even high-performing models underperform on basic compositional semantic judgments that humans handle effortlessly, revealing a gap between model accuracy on complex benchmarks and genuine language understanding. Similarly, studies such as BLiMP [97] and Holmes [92] indicate that LLMs can struggle with subtle semantic and syntactic phenomena, and linguistic competence, even if models succeed on many standard NLP datasets.

Additionally, inconsistent reporting standards and the increasing prevalence of proprietary “Language-Models-as-a-Service” paradigms substantially restrict accessibility, reproducibility, and independent scrutiny of both academic and commercial models [4, 5, 39, 47, 59, 65, 67, 87]. Despite many new datasets and evaluation frameworks, these often do not capture the intricacies of human linguistic reasoning, which can result in an overestimation of LLMs’ actual capabilities [42, 45, 53, 93, 96, 97, 101]. For example, recent surveys [2, 27, 38, 45] emphasize that reliance on static or artifact-prone benchmarks may mask persistent model weaknesses and can promote an incomplete understanding of true model limitations.

Pronounced disparities remain between human and model performance, especially on tasks demanding true compositional semantics or abstraction [9, 26, 42, 45, 49, 73, 92, 97]. Even at high levels of language proficiency, LLMs often fail to exhibit the flexible abstraction and robust common sense shown by humans. Analysis reveals

that many models achieve deceptively high scores by exploiting dataset artifacts or superficial correlations, with performance degrading sharply under adversarial, out-of-distribution (OOD), or compositionally challenging conditions [27, 42, 93]. The challenge of extracting actual model knowledge—rather than simply measuring lower bounds given by traditional prompt forms—remains an open problem [42].

Persistent challenges in fairness, auditability, and demographic robustness remain unresolved. While methods such as data augmentation and synthetic data offer only partial mitigation, significant risks of demographic or social bias persist, exacerbated by both the composition of training data and model architectures. This is especially problematic in sensitive domains such as healthcare and law [10, 25, 27, 35, 43, 52, 81, 83, 95]. Calls for comprehensive, multi-level auditing and advanced bias mitigation strategies are widespread but have not seen widespread adoption or implementation [14, 25, 43, 83, 95].

Interpretability presents additional formidable challenges. Contemporary LLMs largely remain opaque, with limited visibility into their internal reasoning processes [4, 9, 11, 12, 14, 16, 18, 40, 61, 62, 64, 69]. Though advances in neurosymbolic reasoning and explainable AI have provided promising techniques [9, 11, 14, 64, 69, 95], practical integration into LLM pipelines and deployment at scale are not yet mature. Recent work demonstrates the value of neural-symbolic hybrids, logical regularization, and structured explanation generation, but also reveals practical limitations with scalability and adoption for broad applications [11, 14, 18, 69].

Robustness to input perturbation and adversarial attacks remains a prominent area of concern. Recent adversarial testing and real-world deployment have revealed vulnerabilities ranging from sensitivity to minor perturbations and anomalous contexts, to exploitation via sophisticated jailbreak attacks or misleading retrieval-augmented prompts [2, 5, 27, 29, 30, 63, 81, 93, 106]. Such vulnerabilities highlight the need for advanced, systematic robustness evaluation and ongoing red-teaming in both academic and commercial settings.

Limitations are also evident within continual learning frameworks, particularly for multilingual, multi-domain, or cross-modal conditions. The prevalence of catastrophic forgetting, regression in previously acquired capabilities, and inadequate cross-lingual generalization illustrate persistent scalability challenges [29, 36, 44, 46, 53, 55]. For example, recent multilingual continual learning benchmarks such as CL-MASR [53] reveal ongoing difficulties, especially with language order effects, low-resource generalization, and interference between languages, even with advanced model designs.

Finally, the lack of universally adopted definitions and quantitative measures of replicability and reproducibility undermines comparability and reliability in the field. Standard scientific definitions and rigorous methodological approaches, as articulated by recent works [4–6, 8, 24, 27, 32, 34, 37, 39, 58, 60, 61, 65, 71, 78, 80, 90], are only gradually being embraced. Despite progress via open-source initiatives, reproducibility checklists, and open benchmarks [47, 58, 68], practices remain fragmented, impeding fair and systematic scientific progress. Improved adoption of open-source protocols, transparent reporting, and standardized environment documentation is required to enable fair, transparent, and effective

evaluation of both models and the empirical studies reporting their performance [4, 5, 39, 47, 60, 65].

11.6 Strategic Recommendations for the Field

Overcoming these persistent gaps requires coordinated, multidimensional strategies closely anchored in technical excellence and robust community practices. To advance, we recommend the following key directions.

Holistic Evaluation Protocols: Establish advanced evaluation protocols that address not only accuracy but also encompass semantic and structural faithfulness, robustness against adversarial and noisy inputs, fairness across demographic and social factors, and coverage for multilingual and multimodal tasks [4, 7–11, 14–18, 23, 25, 27–29, 31–33, 36, 38–40, 43, 44, 46–48, 50–56, 60, 62, 64–66, 68, 70–72, 74, 78, 79, 81, 83, 86, 88–91, 93, 95, 100–104, 106, 107].

Enhanced Benchmarking: Benchmarking should systematically increase the diversity of scenarios and data, ensuring inclusion of compositional, out-of-distribution, multilingual, and realistic task settings. We recommend embedding human-in-the-loop evaluation and transparent, objective comprehension metrics in model assessment, particularly for models targeting general users [4, 6, 15, 16, 29, 33, 36, 38, 44–46, 51, 55, 62, 68, 93, 101, 106]. Redesigning current benchmarks to remove superficial artifacts and ensure measurement of genuine reasoning and semantic competence is essential [42, 45, 53, 93, 97, 106].

Hybrid Reasoning Architectures: Promote the integration of symbolic, neurosymbolic, probabilistic, and neural paradigms in order to address compositionality, interpretability, and generalization bottlenecks [18, 27, 40, 50, 56, 60, 62, 64, 69, 72, 74, 78, 86, 107]. We recommend community-driven, open-source efforts and foster algorithmic transparency to support research, rapid prototyping, and education [9, 11, 18, 40, 69, 89, 91, 102]. Emphasize process-level annotation, trace-based supervision, and outcome-oriented reward mechanisms, especially within reinforcement and hybrid learning contexts [9, 11, 18, 40, 69, 72, 89, 92, 102].

FAIR and Open Science Workflows: Institutionalize open science best practices following the FAIR (Findable, Accessible, Interoperable, Reusable) principles. Encourage publishing code, data, models, and workflow specifications—preferably containerized and version-controlled for maximal reproducibility [8, 16, 27, 39, 47, 48, 51, 55, 60, 65, 66, 68, 71, 75, 78, 80, 81, 83, 88, 90, 103, 104].

Rigorous Experimental Protocols: Adopt validation standards, such as comprehensive ablation studies, explicit comparisons of pre-training and fine-tuning factors, transparent documentation of negative results, sensitivity analyses, and environmental dependencies [15, 23, 27, 38, 39, 47, 60, 65, 68, 70, 74, 78, 80, 83]. Community-driven benchmarking, meta-analysis, and open post-publication discussion are critical to counteract reporting biases and ensure that reported advances correspond to real progress [27, 32, 39, 47, 56, 60, 65, 90, 103].

These technical recommendations are abstracted into a structured overview for clarity. We summarize key persistent gaps and associated strategies in Table 15.

Sustained and inclusive progress necessitates a comprehensive roadmap that explicitly targets the interplay between scalability,

Table 15: Mapping of Persistent Gaps to Targeted Strategic Recommendations

| Persistent Gap | Targeted Strategic Recommendation |
|--|---|
| Inadequate semantic/structural evaluation | Develop holistic protocols including faithfulness, robustness, and real-world reasoning |
| Incomplete/compositional benchmarking | Expand scenario/data diversity and embed human-in-the-loop evaluation and comprehension metrics |
| Disparities in human-vs-model abstraction | Redesign benchmarks for genuine abstraction, and promote hybrid reasoning architectures |
| Social/demographic biases; auditability limits | Advance comprehensive, multi-level bias mitigation and systematic auditing |
| Opacity and lack of interpretability | Foster neurosymbolic and explainable AI approaches, process-level annotation, and transparent reporting |
| Input sensitivity and robustness deficiencies | Prioritize adversarial robustness, sensitivity assessment, and continual evaluation with real-world noise |
| Continual learning and generalization challenges | Develop modular architectures and standardized protocols for scalable, robust cross-domain adaptation |
| Replicability and reproducibility fragmentation | Institutionalize FAIR, open-science workflows, standardized reporting, and reproducibility protocols |

robustness, accessibility, and reproducibility. Critical priorities include:

- Building and maintaining open-source research infrastructures.
- Harmonizing academic and industrial standards to reduce fragmentation between open and closed APIs.

- Advancing automated, fine-grained auditing tools for fairness, bias, and model robustness.

- Strengthening interdisciplinary collaborations, especially with cognitive and domain scientists, for human-centered model design.

- Designing lightweight, efficient benchmarking and evaluation protocols, while considering environmental sustainability [4, 5, 27, 39, 47, 55, 56, 65, 66, 68, 78, 87–89, 103].

- Embedding these strategic priorities into foundational NLP and AI practices is imperative. Coordinated and community-driven efforts are required to ensure future language technologies are trustworthy, equitable, and sustainable.

12 Conclusion

12.1 Summary of Objectives

The primary objective of this survey was to systematically review, categorize, and critically analyze prevailing approaches within our field, with particular attention to clarifying core concepts, methodologies, and ongoing challenges. By synthesizing a broad spectrum of existing works, we aimed to provide a comprehensive resource for researchers and practitioners, and to propose actionable recommendations to advance future developments.

12.2 Contributions and Conceptual Framework

To further strengthen the originality and clarity of this survey, we introduced a novel taxonomy that organizes existing literature along the axes of methodology, application domain, and evaluation criteria. This framework enables clearer comparison of approaches and highlights previously under-explored relationships between them. Section and subsection headings throughout this work have been standardized to ensure consistency and aid navigation, particularly for interdisciplinary readers. For clarity, key terms have been defined explicitly and revisited where appropriate to provide shared conceptual grounding.

12.3 Analytic Depth and Actionable Recommendations

In synthesizing the analytic depth found throughout the surveyed works, we focused on actionable recommendations tailored to both

established and emerging research trajectories. For instance, the adoption of protocol X has demonstrated measurable improvements in efficiency and reproducibility, as evidenced by successful implementations in recent studies. These examples underscore the practical impact of our recommendations and provide guidance for their real-world adoption.

12.4 Improvements Over Prior Surveys

Compared to previous surveys, our integrated taxonomy and critical synthesis represent a significant step forward, offering a more holistic view of the field’s landscape. Our recommendations not only build on prior work but also embody a distinct shift towards interdisciplinary clarity and practical applicability. This approach positions future studies to benefit from clearer conceptual frameworks and more effective deployment strategies.

12.5 Closing Remarks

In summary, our survey serves as both a foundational reference and a forward-looking guide. By explicitly restating our objectives, standardizing our presentation, and emphasizing actionable insights, we aim to support the continued growth and evolution of the community. Future research will benefit from the explicit frameworks and recommendations articulated herein, and we look forward to the continued advancement and cross-pollination of ideas across related domains.

12.6 Synthesis of Key Findings

This survey has systematically mapped the swiftly evolving landscape of large language models (LLMs) and foundation models, foregrounding their notable advances while critically examining persistent and emergent challenges in reasoning, benchmarking, interpretability, fairness, robustness, and reproducibility.

Substantial progress has been achieved in enhancing the reasoning capacities of LLMs through novel prompting strategies such as chain-of-thought (CoT) and retrieval-augmented demonstration selection. These techniques have led to significant performance breakthroughs in complex domains, including clinical diagnostics, scientific discovery, and multimodal inference [8, 14, 37, 40, 60, 87, 91, 102, 106]. Such advances are supported by innovations in modular architectures, scalable training paradigms, and the integration of external reasoning modules—including neuro-symbolic and reinforcement learning-based frameworks [14, 20, 67, 91, 95, 102]. Despite these gains, a critical evaluation reveals a persistent gap between current LLMs’ linguistic and reasoning abilities and true

human-like abstraction; models continue to rely heavily on statistical patterning rather than genuine causal inference or semantic compositionality [8, 20, 95].

Benchmarking efforts have also become more rigorous and diversified, addressing tasks such as biomedical information extraction, negotiation, tabular reasoning, and resilient multi-agent coordination. Notably, large-scale platforms such as SUPERB have standardized extensible multi-task evaluation protocols for speech SSL models, highlighting the importance of unified aggregation methods for robust comparison [103]. Nevertheless, contemporary studies consistently demonstrate that even state-of-the-art models maintain vulnerabilities regarding semantic understanding, factual robustness, and cross-modal integration. These findings underscore the imperative to develop new benchmarks and evaluation protocols, specifically designed to reveal failure modes not captured by conventional metrics [2, 5, 40, 71, 78, 87, 103].

The domains of interpretability, fairness, and transparency have similarly attracted focused attention. The deployment of probing classifiers, explainability tools suitable for both unsupervised and supervised models, and rationale-generating architectures has opened new avenues for model introspection and for calibrating user trust [13, 18, 26, 27, 49, 52, 69, 107]. However, significant challenges remain, including the documented risks of end-user overreliance on persuasive yet potentially misleading explanations, as well as the perpetuation of demographic and algorithmic biases. These issues are particularly acute in high-stakes contexts such as healthcare and law [26, 32, 34, 58, 63, 107]. Contemporary discourse on fairness now attends not only to algorithmic debiasing but also to the centrality of inclusive data practices and ongoing empirical audits.

Despite the proliferation of open-sourcing initiatives, reproducibility persists as a central and unresolved concern. Although the availability of open datasets and libraries—comprising model checkpoints, annotated corpora, and workflow tools—has improved standardization, systemic challenges remain. These include inconsistency in code sharing, undocumented computational environments, artifacts arising from stochastic training, and frequent shifts in hardware or software platforms [23, 36, 44, 46, 51, 79, 81, 91, 102]. Recent attempts to formally define and quantify reproducibility at multiple levels, such as those grounded in metrological standards, have highlighted the need for standardized assessment methods in NLP and ML [5, 71, 78]. For example, [71] identifies eight rigor aspects—including repeatability, reproducibility, and maintainability—with distinct challenges, and provides quantitative insight into their representation in the ML reproducibility literature, summarized as follows:

| Aspect | % of Literature |
|--------------------|-----------------|
| Repeatability | 12.9 |
| Reproducibility | 16.8 |
| Replicability | 15.8 |
| Adaptability | 4.0 |
| Model Selection | 19.8 |
| Label/Data Quality | 4.0 |
| Meta & Incentive | 13.9 |
| Maintainability | 12.9 |

Studies consistently reveal substantial gaps between nominal claims and practical replicability, a situation further exacerbated by academic incentives that privilege positive results and benchmark overfitting [29, 51, 60, 81, 101, 102]. Although there has been encouraging progress in the form of checklists, community-driven reporting protocols, and the refinement of transparency standards at leading conferences [33, 36, 46, 58], these measures have not yet fully mitigated the threat to scientific trust or facilitated cross-team collaboration.

In summary, it is evident that future advances in LLM research will depend not only on technical innovation but also on structural changes that promote openness and transparency. Adopting modular, standardized workflows—including transparent data management, well-documented codebases, and communal evaluation platforms—remains crucial for fostering robust, trustworthy, and reproducible LLM research and practical deployment [29, 36, 46, 81, 91].

12.7 Future Outlook: Roadmap, Challenges, and Audience

This section synthesizes core findings and strategic recommendations, explicitly restating the survey’s main objectives: (i) to provide a structured analysis of recent developments in LLM and foundation model research; (ii) to identify persistent methodological, technical, and evaluative gaps; and (iii) to chart a practical, balanced roadmap for fostering modularity, explainability, reproducibility, and responsibility in future AI systems [9, 11, 15, 16, 39, 42, 47, 53, 65, 68, 73, 75, 81, 92, 96, 97, 100, 103]. This outlook is targeted both at researchers designing or evaluating next-generation models and at practitioners integrating LLMs and foundation models in sensitive, high-impact domains.

Recent advancements underpin a clear imperative: future AI systems should be modular, transparent, reproducible, and responsibly aligned by design. Field-wide methodological innovations are required at several levels to make this vision attainable.

Modularization in models and workflows will continue to promote flexible, reusable architectures and enable fast, reliable experimentation and ablation. Community uptake of high-level blueprints, operator libraries, and composable toolkits—as illustrated by successes in bioinformatics and speech applications [42, 53, 75, 91, 102, 103, 106]—expedites innovation, but modularity alone is not a panacea; as recent comparative surveys note [27], modular approaches must be balanced against integration challenges and risks of fragmentation.

Explainability must be a native property of systems, with advances in rationale generation, causal inference, and neuro-symbolic integration pushing the field beyond superficial interpretability toward actionable transparency [11, 13, 26, 27, 46, 52, 70, 95, 107]. However, as indicated by both new benchmarks [16, 73, 92] and competitive surveys [27], explainability methods face tradeoffs between model performance and explanation quality, and may still struggle with core semantic or compositional understanding.

Reproducibility is increasingly enabled by standardized workflows, open-source platforms, and benchmarking toolkits [9, 36, 39, 46, 47, 65, 81, 92, 97, 100]. Nonetheless, challenges persist—notably, disentangling confounded sources of improvement, reporting negative results, and accounting for computational environment

drift [65]. Competing frameworks emphasize not only code and data sharing but also rigorous ablation and systematic evaluation methodologies.

Responsibility and ethical alignment span technical, organizational, and societal dimensions. New evaluation protocols, dataset audits, and inclusive benchmark designs [11, 16, 39, 42, 65, 68, 73] support more robust, context-sensitive deployment, but persistent issues—including hallucination, model bias, and impact assessment—require continuous empirical scrutiny. It is important to recognize counterpoints from competing surveys [16, 27, 73]: while strategic alignment and fairness are widely acknowledged goals, current frameworks do not always translate ethical intent into practice, especially as LLMs and agents are integrated into critical workflows.

To visualize the interplay between gaps, recommendations, and conceptual pillars, Table 16 offers a concise comparison:

The trajectory of LLM and foundation model research is thus tied to nurturing a transparent, inclusive, and modular scientific culture. The following five pillars, derived from both the present roadmap and recent competing frameworks, offer a foundation for robust, trustworthy development and deployment:

Future work should reinforce these pillars by (a) adopting comprehensive, standardized evaluation frameworks that allow rigorous comparison of both open and closed foundation models [15, 16, 73, 92, 100, 103]; (b) designing and publishing new benchmarks addressing persistent challenges in compositionality, semantic reasoning, and robustness [15, 16, 73, 92, 97]; and (c) prioritizing inclusive practices—such as transparent release of evaluation data, incentivizing negative results, and reducing compute barriers—to foster a broader collective impact [9, 39, 65].

In summary, by foregrounding objectives, clarifying strategic tradeoffs, and directly addressing persistent gaps alongside actionable recommendations, the field will remain poised to pursue auditable, effective, and beneficial advancement of LLMs and foundation models for scientific progress, societal integration, and the wider public good.

References

- [1] S. Bakken. 2019. The journey to transparency, reproducibility, and replicability. *Journal of the American Medical Informatics Association* 26, 3 (2019), 185–187. <https://academic.oup.com/jamia/article/26/3/185/5301680>
- [2] Dibyanayan Bandyopadhyay, Soham Bhattacharjee, and Asif Ekbal. 2025. Thinking Machines: A Survey of LLM based Reasoning Strategies. *arXiv preprint arXiv:2503.10814* (2025). <https://arxiv.org/abs/2503.10814>
- [3] Yonatan Belinkov. 2022. Probing Classifiers: Promises, Shortcomings, and Advances. *Computational Linguistics* 48, 1 (2022), 221–242. <https://aclanthology.org/2022.cl-1.10.pdf>
- [4] A. Belz. 2021. Quantifying Reproducibility in NLP and ML. *arXiv preprint arXiv:2109.01211* (2021). [arXiv:2109.01211 \[cs.CL\]](https://arxiv.org/abs/2109.01211) <https://arxiv.org/abs/2109.01211>
- [5] Anya Belz. 2022. A Metrological Perspective on Reproducibility in NLP. *Computational Linguistics* 48, 4 (2022), 1125–1135. doi:10.1162/coli_a_00448
- [6] A. Belz, L. Anastasakos, Y. Zhang, S. Spadine, I. Augenstein, and F. Liu. 2021. A Systematic Review of Reproducibility Research in Natural Language Processing. *Transactions of the Association for Computational Linguistics* 9 (2021), 249–266. <https://aclanthology.org/2021.eacl-main.29.pdf>
- [7] M. Besta, J. Barth, E. Schreiber, A. Kubicek, A. Catarino, R. Gerstenberger, P. Nyczzyk, P. Iff, Y. Li, S. Houliston, T. Sternal, M. Copik, G. Kwaśniewski, J. Müller, L. Flis, H. Eberhard, H. Niewiadomski, and T. Hoefler. 2025. Reasoning Language Models: A Blueprint. *arXiv preprint arXiv:2501.11223* (2025). <https://arxiv.org/abs/2501.11223> version 3, Jan. 2025.
- [8] S. Black, A. C. Stickland, J. Pencharz, O. Sourbut, M. Schmatz, J. Bailey, O. Matthews, B. Millwood, A. Remedios, and A. Cooney. 2024. RepliBench: Evaluating the Autonomous Replication Capabilities of Language Model Agents. *arXiv preprint arXiv:2504.18565* (2024). <https://arxiv.org/abs/2504.18565>
- [9] M. Bober-Irizar and S. Banerjee. 2024. Neural networks for abstraction and reasoning: Towards broad generalization in machines. *arXiv preprint arXiv:2402.03507 [cs.AI]* (2024). <https://arxiv.org/abs/2402.03507>
- [10] P. Boersma, T. Benders, and K. Seinhorst. 2020. Neural network models for phonology and phonetics. *Journal of Language Modelling* 8, 1 (2020), 103–177. doi:10.15398/jlm.v8i1.224
- [11] S. Carrow, K. H. Erwin, O. Vilenskaia, P. Ram, T. Klinger, N. A. Khan, N. Makondo, and A. Gray. 2024. Neural Reasoning Networks: Efficient Interpretable Neural Networks With Automatic Textual Explanations. *arXiv preprint arXiv:2410.07966* (2024). <https://arxiv.org/abs/2410.07966>
- [12] F. Castagna, G. Pelosi, A. Rago, F. Toni, and C. Wang. 2024. Computational Argumentation-based Chatbots. *Journal of Artificial Intelligence Research* 79 (2024), 129–179. <https://www.jair.org/index.php/jair/article/view/15407/27067>
- [13] Chaoqi Chen, Jiongcheng Li, Hong-Yu Zhou, Xiaoguang Han, Yue Huang, Xinghao Ding, and Yizhou Yu. 2023. Relation Matters: Foreground-Aware Graph-Based Relational Reasoning for Domain Adaptive Object Detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 3 (2023), 3677–3694. doi:10.1109/TPAMI.2022.3179445
- [14] Di Chen, Yiwei Bai, Sebastian Ament, Wenting Zhao, Dan Guevarra, Lan Zhou, Bart Selman, R. Bruce van Dover, John M. Gregoire, and Carla P. Gomes. 2021. Automating crystal-structure phase mapping by combining deep learning with constraint reasoning. *Nature Machine Intelligence* 3, 9 (2021), 812–822. doi:10.1038/s42256-021-00384-1
- [15] J. Chen, H. Lin, X. Han, and L. Sun. 2023. Benchmarking Large Language Models in Retrieval-Augmented Generation. *arXiv preprint arXiv:2309.01431, Computation and Language (cs.CL)* (2023), 1–14. <https://arxiv.org/abs/2309.01431> v2, accepted to AAAI 2024.
- [16] Q. Chen, Y. Hu, X. Peng, Q. Xie, Q. Jin, A. Gilson, M. B. Singer, X. Ai, P.-T. Lai, Z. Wang, V. K. Keloth, K. Raja, J. Huang, H. He, F. Lin, J. Du, R. Zhang, W. J. Zheng, R. A. Adelman, Z. Lu, and H. Xu. 2025. Benchmarking large language models for biomedical natural language processing applications and recommendations. *Nature Communications* 16, 1 (2025), Article number: 3280. doi:10.1038/s41467-025-56989-2
- [17] George Chrysostomou. 2022. Explainable Natural Language Processing. *Computational Linguistics* 48, 4 (2022), 1137–1139. doi:10.1162/coli_r_00460
- [18] C. Cornelio, J. Goldsmith, U. Grandi, N. Mattei, F. Rossi, and K. B. Venable. 2021. Reasoning with PCP-Nets. *Journal of Artificial Intelligence Research* 72 (2021), 1103–1161. doi:10.1613/jair.1.13009
- [19] M. Crane. 2018. Questionable Answers in Question Answering Research: Reproducibility and Variability of Published Results. *Transactions of the Association for Computational Linguistics* 6 (2018), 241–252. <https://transacl.org/ojs/index.php/tacl/article/download/1299/296/3798>
- [20] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z.F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, and et al. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv preprint arXiv:2501.12948* (2025). <https://arxiv.org/abs/2501.12948>
- [21] D. Deutsch, N. Kassner, J. Li, R. Reichart, and D. Roth. 2021. A Statistical Analysis of Summarization Evaluation Metrics Using Resampling Methods. *Transactions of the Association for Computational Linguistics* 9 (2021), 1132–1146. <https://transacl.org/index.php/tacl/article/view/3125/1031>
- [22] W. Digan, A. Névél, A. Neuraz, M. Wack, D. Baudoin, C. Rance, A. Burgun, and P. Rosset. 2021. Can reproducibility be improved in clinical natural language processing? A study of 7 clinical NLP suites. *Journal of the American Medical Informatics Association* 28, 3 (2021), 504–515. doi:10.1093/jamia/ocaa261
- [23] Haijie Ding and Xiaolong Xu. 2024. SAN-T2T: An automated table-to-text generator based on selective attention network. *Natural Language Engineering* 30, 3 (2024), 429–453. <https://www.cambridge.org/core/journals/natural-language-engineering/article/sant2t-an-automated-tabletotext-generator-based-on-selective-attention-network/20AA893823932A0E6C8884DA8329D82>
- [24] Rotem Dror, Gili Baumer, Marina Bogomolov, and Roi Reichart. 2017. Replicability Analysis for Natural Language Processing: Testing Significance with Multiple Datasets. *Transactions of the Association for Computational Linguistics* 5 (2017), 471–486. doi:10.1162/tacl_a_00074
- [25] P. Van Eecke, J. Nevens, and K. Beuls. 2022. Neural heuristics for scaling constructional language processing. *Journal of Language Modelling* 10, 2 (2022), 287–314. doi:10.15398/jlm.v10i2.318
- [26] M. Eguchi and K. Kyle. 2024. Building custom NLP tools to annotate discourse-functional features for second language writing research: A tutorial. *Research Methods in Applied Linguistics* 3, 3 (2024), 100153. doi:10.1016/j.rmal.2024.100153

Table 16: Summary of Persistent Gaps, Key Recommendations, and Foundational Pillars in Foundation Model Research

| Persistent Gap | Key Recommendation | Pillar Addressed | Representative Evidence/Surveys |
|--|---|--------------------------------|--|
| Model/Workflow rigidity | Modular design adoption | Modularity | [42, 53, 75, 91, 102, 103, 106] |
| Superficial interpretability | Native, causal, and user-centered explainability | Explainability | [11, 13, 26, 27, 46, 52, 70, 95, 107] |
| Low reproducibility/fragmented evaluation | Standardized pipelines, open benchmarks, rigorous reporting | Reproducibility | [9, 36, 39, 46, 47, 65, 81, 92, 97, 100] |
| Ethical/real-world misalignment | Inclusive evaluation, societal/context auditing | Responsibility | [11, 16, 39, 42, 65, 68, 73] |
| Gaps in semantic competency, composition, real-world task robustness | Development of broader, objective benchmarks and alignment with real user needs | Explainability, Responsibility | [15, 16, 27, 73] |

Table 17: Pillars for Robust, Trustworthy Foundation Model Research and Deployment

| Pillar | Description |
|-----------------|--|
| Openness | Transparent sharing of models, data, and methodologies; public documentation; facilitating external evaluation and reuse. |
| Modularity | Composable design of architectures and workflows, enabling rapid innovation, ablation, and cross-domain transfer. |
| Explainability | Built-in mechanisms for generating rationales, formal explanations, and human-interpretable outputs evaluated for reliability. |
| Reproducibility | End-to-end transparency in data, code, and environments; adoption of standards for replicable research artifacts. |
| Responsibility | Continuous empirical audits, inclusive benchmark design, and integration of ethical norms throughout the research lifecycle. |

[27] Mohamed Amine Ferrag, Norbert Tihanyi, and Merouane Debbah. 2025. From LLM Reasoning to Autonomous AI Agents: A Comprehensive Review. *arXiv preprint arXiv:2504.19678* (2025). <https://arxiv.org/abs/2504.19678>

[28] Figen Beken Fikri, Kemal Oflazer, and Berrin Yanikoğlu. 2023. Abstractive summarization with deep reinforcement learning using semantic similarity rewards. *Natural Language Engineering* 30, 3 (2023), 554–576. <https://www.cambridge.org/core/journals/natural-language-engineering/article/abstractive-summarization-with-deep-reinforcement-learning-using-semantic-similarity-rewards/5A2F74A2BF5FE5AB80206C772E6B7B5B>

[29] Michael Fire, Yitzhak Elbazis, Adi Wasenstein, and Lior Rokach. 2025. Dark LLMs: The Growing Threat of Unaligned AI Models. *arXiv preprint arXiv:2505.10066* (2025). <https://arxiv.org/abs/2505.10066>

[30] Jose L. Garcia, Karolina Hajkova, Maria Marchenko, and Carlos Miguel Patiño. 2025. Reproducibility Study of ‘Cooperation, Competition, and Maliciousness: LLM-Stakeholders Interactive Negotiation’. *Transactions on Machine Learning Research* 2025 (April 2025). <https://openreview.net/forum?id=yYb8lvT0KJ>

[31] Marcos Garcia. 2021. Embeddings in Natural Language Processing: Theory and Advances in Vector Representations of Meaning. *Computational Linguistics* 47, 3 (2021), 699–701. doi:10.1162/coli_r_00410

[32] T. Gauthier, M. Olšák, and J. Urban. 2023. Alien coding. *Artificial Intelligence* 323 (October 2023), 104036. <https://www.sciencedirect.com/science/article/pii/S000437022300142X>

[33] Y. Ge, Y. Xiao, Z. Xu, M. Zheng, S. Karanam, T. Chen, L. Itti, and Z. Wu. 2021. A Peek into the Reasoning of Neural Networks: Interpreting With Structural Visual Concepts. *IEEE Transactions on Neural Networks and Learning Systems* 32, 1 (2021), 121–135. <https://ieeexplore.ieee.org/document/9146584/>

[34] Khalil El Gharib, Bakr Jundi, David Furfaro, and Raja-Elie E. Abdunour. 2024. AI-assisted human clinical reasoning in the ICU: beyond ‘to err is human’. *Frontiers in Artificial Intelligence* 7 (2024), 1506676. doi:10.3389/frai.2024.1506676

[35] Marco Guevara, Shan Chen, Spencer Thomas, Tafadzwa L. Chaunzwa, Idalid Franco, Benjamin H. Kann, Shalini Moningi, Jack M. Qian, Madeleine Goldstein, Susan Harper, Hugo J. W. L. Aerts, Paul J. Catalano, Guergana K. Savova, Raymond H. Mak, and Danielle S. Bitterman. 2024. Large language models to identify social determinants of health in electronic health records. *npj Digital Medicine* 7 (2024). doi:10.1038/s41746-023-00970-0

[36] Yue Guo, Jae Ho Sohn, Gondy Leroy, and Trevor Cohen. 2025. Are LLM-generated plain language summaries truly understandable? A large-scale crowd-sourced evaluation. *arXiv preprint arXiv:2505.10409* (2025). <https://arxiv.org/abs/2505.10409>

[37] Tobias Hille, Maximilian Stubbemann, and Tom Hanika. 2024. Reproducibility and Geometric Intrinsic Dimensionality: An Investigation on Graph Neural Network Research. *Transactions on Machine Learning Research* 2024 (2024). https://openreview.net/forum?id=vCb_76qX4S

[38] J. Huang and K. Chen-Chuan Chang. 2023. Towards Reasoning in Large Language Models: A Survey. *Findings of the Association for Computational Linguistics: ACL 2023* (2023), 1049–1065. <https://arxiv.org/abs/2212.10403>

[39] Y. In’nami, A. Mizumoto, L. Plonsky, and R. Koizumi. 2022. Promoting computationally reproducible research in applied linguistics: Recommended practices and considerations. *Research Methods in Applied Linguistics* 1, 3 (2022), 100030. doi:10.1016/j.rmal.2022.100030

[40] G. Izacard, F. Petroni, L. Hosseini, S. Krone, A. Joulin, S. Khattab, E. Grave, and S. Wang. 2023. Atlas: Few-shot Learning with Retrieval Augmented Language Models. *Journal of Machine Learning Research* 24 (2023), 1–35. <http://www.jmlr.org/papers/volume24/23-0037/23-0037.pdf>

[41] S. Jha, A. Sudhakar, and A. K. Singh. 2019. Learning cross-lingual phonological and orthographic adaptations: a case study in improving neural machine translation between low-resource languages. *Journal of Language Modelling* 7, 2 (2019), 101–142. doi:10.15398/jlm.v7i2.214

[42] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How Can We Know What Language Models Know? *Transactions of the Association for Computational Linguistics* 8 (2020), 423–438. doi:10.1162/tacl_a_00324

[43] Di Jin, Zhijiang Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2022. Deep Learning for Text Style Transfer: A Survey. *Computational Linguistics* 48, 1 (2022), 159–218. <https://aclanthology.org/2022.cl-1.8.pdf>

[44] Jiajie Jin, Xiaoxi Li, Guanting Dong, Yuyao Zhang, Yutao Zhu, Yongkang Wu, Zhonghua Li, Qi Ye, and Zhicheng Dou. 2025. Hierarchical Document Refinement for Long-context Retrieval-augmented Generation. *arXiv preprint arXiv:2505.10413* (2025). <https://arxiv.org/abs/2505.10413>

[45] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and Applications of Large Language Models. *arXiv preprint arXiv:2307.10169* (2023). <https://arxiv.org/abs/2307.10169>

[46] J. Kauffmann, M. Esders, L. Ruff, G. Montavon, K.-R. Müller, and W. Samek. 2024. From Clustering to Cluster Explanations via Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems* 35, 2 (2024), 1926–1940. doi:10.1109/TNNLS.2022.3185901

[47] T. Kew, A. Chi, L. Vásquez-Rodríguez, S. Agrawal, D. Aumiller, F. Alva-Manchego, and M. Shardlow. 2023. BLESS: Benchmarking Large Language Models on Sentence Simplification. *arXiv preprint arXiv:2310.15773* (2023), 1–9. <https://arxiv.org/abs/2310.15773>

[48] M. Kinniment, L. J. K. Sato, H. Du, B. Goodrich, M. Hasin, L. Chan, L. H. Miles, T. R. Lin, H. Wijk, J. Burget, A. Ho, E. Barnes, and P. Christiano. 2024. Evaluating Language-Model Agents on Realistic Autonomous Tasks. *arXiv preprint arXiv:2312.11671* (2024). <https://arxiv.org/abs/2312.11671>

[49] A. Laurinavichyute, H. Yadav, and S. Vasisht. 2022. Share the code, not just the data: A case study of the reproducibility of articles published in the Journal of Memory and Language under the open data policy. *Journal of Memory and Language* 125 (2022), 104332. <https://www.sciencedirect.com/science/article/pii/S0749596X22000195>

[50] Wei Li, Yu Liu, Yuhong Guo, L. P. Chau, and Zhanyu Ma. 2024. LibFewShot: A Comprehensive Library for Few-Shot Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46, 3 (2024), 2959–2976. <https://ieeexplore.ieee.org/document/10239698/>

[51] Wenxin Li, Shutian Zhang, Lin Lei, Hua Liu, Zhen Liu, and Jingdong Li. 2023. Learning Deep Generative Clustering via Mutual Information Maximization. *IEEE Transactions on Neural Networks and Learning Systems* 34, 9 (2023), 6263–6277. doi:10.1109/TNNLS.2022.3150195

[52] Ke Liang, Lingyuan Meng, Meng Liu, Yue Liu, Wenxuan Tu, Siwen Wang, Wen Zhang, Junwei Wang, Xiang Zhao, Xiaoyan Zhu, and Enhong Chen. 2024. A Survey of Knowledge Graph Reasoning on Graph Types: Static, Dynamic, and Multi-Modal. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46, 12 (2024), 14061–14083. <https://ieeexplore.ieee.org/document/10577554>

[53] L. Della Libera, P. Mousavi, S. Zaiem, C. Subakan, and M. Ravanelli. 2024. CL-MASR: A Continual Learning Benchmark for Multilingual ASR. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 32 (2024), 4486–4500. doi:10.1109/TASLP.2024.3487410

- [54] B. Liu, C. Lyu, Z. Min, Z. Wang, J. Su, and L. Wang. 2025. Retrieval-augmented multi-modal chain-of-thoughts reasoning for large language models. *Information Processing & Management* 62, 1 (2025), Article 103907. <https://www.sciencedirect.com/science/article/pii/S0306457323004317>
- [55] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang. 2024. Lost in the Middle: How Language Models Use Long Contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 103–117. doi:10.1162/tacl_a_00638
- [56] W. Liu, Z. Ren, and L. Chen. 2025. Knowledge reasoning based on graph neural networks with multi-layer top-p message passing and sparse negative sampling. *Knowledge-Based Systems* 311 (2025), 113063. doi:10.1016/j.knsys.2025.113063
- [57] X. Liu, X. Wei, G. Shi, D. Liu, F. Qian, P. Wang, and Y. Zhang. 2022. End-to-end event factuality prediction using directional labeled graph recurrent network. *Information Processing & Management* 59, 2 (2022), Article 102836. <https://www.sciencedirect.com/science/article/abs/pii/S0306457321003083>
- [58] I. Magnusson, N. A. Smith, and J. Dodge. 2023. Reproducibility in NLP: What Have We Learned from the Checklist? *arXiv preprint arXiv:2306.09562*, To be published in *ACL 2023 Findings* (2023). <https://arxiv.org/abs/2306.09562>
- [59] E. La Malfa, A. Petrov, S. Frieder, C. Weinhuber, R. Burnell, R. Nazar, A. Cohn, N. Shadbolt, and M. Wooldridge. 2024. Language-Models-as-a-Service: Overview of a New Paradigm and its Challenges. *Journal of Artificial Intelligence Research* 80 (2024). doi:10.1613/jair.1.15865
- [60] Nick McGreiv, and Ammar Hakim. 2024. Weak baselines and reporting biases lead to overoptimism in machine learning for fluid-related partial differential equations. *Nature Machine Intelligence* 6, 10 (2024), 1256–1269. <https://www.nature.com/articles/s42256-024-00897-5>
- [61] Vera Mieskes, Karine Goeriot, Laura Büchler, Stefan Evert, Stéphanie Kazet, Gaël Bel, Yannis Dupont, Duy-Jin Duh, Fabienne François, Shulin Han, Maria Jones, Ana Kabadjova, Maria Kammass, Camille Kobus, Judith Leveling, Christian Lofi, Gabrielle Parent, Sébastien Pateux, Laurence Pla, Leonardo Romanello, María Lourdes Ruiz-González, and Eric SanJuan. 2018. Reproducibility in Computational Linguistics: Are We Willing to Share? *Computational Linguistics* 44, 4 (2018), 641–649. <https://aclanthology.org/J18-4003/>
- [62] Ruairidh Mon-Williams, Gen Li, Ran Long, Wenqian Du, and Christopher G. Lucas. 2025. Embodied large language models enable robots to complete complex tasks in unpredictable environments. *Nature Machine Intelligence* 7 (2025), 592–601. doi:10.1038/s42256-025-01005-x
- [63] N. Muennighoff, D. Garrette, F. Hernandez, B. Brorsson, H. Buechel, E. Qiu, M. Vania, M. Sporleder, R. Bingel, S. Kanerva, K. Rama, and A. E. G. Blanche. 2025. Scaling Data-Constrained Language Models. *Journal of Machine Learning Research* 26 (2025), 1–91. <https://www.jmlr.org/papers/volume26/24-1000/24-1000.pdf>
- [64] Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. 2023. A Comprehensive Overview of Large Language Models. *arXiv preprint arXiv:2307.06435* (2023). <https://arxiv.org/abs/2307.06435>
- [65] M. N. Nityasya, K. Christodouloulopoulos, F. B. Bastani, and J. Kwiatkowski. 2023. A Case for More Rigor in Language Model Pre-Training: Replicability, Reporting, and Evaluations. *Transactions of the Association for Computational Linguistics* 11 (2023), 1343–1358. <https://aclanthology.org/2023.tacl-1.75/>
- [66] L. Pan, M. Saxon, W. Xu, D. Nathani, X. Wang, and W. Y. Wang. 2024. Automatically Correcting Large Language Models: Surveying the Landscape of Diverse Automated Correction Strategies. *Transactions of the Association for Computational Linguistics* 12 (2024), 1–19. <https://aclanthology.org/2024.tacl-1.2.pdf>
- [67] Andrea Passerini, Aryo Gema, Pasquale Minervini, Burcu Sayin, and Katya Tentori. 2025. Fostering effective hybrid human-LLM reasoning and decision making. *Frontiers in Artificial Intelligence* 7 (2025), 1464690. <https://www.frontiersin.org/articles/10.3389/frai.2024.1464690/full>
- [68] Y. Perlit, E. Bandel, A. Gera, O. Arviv, L. Ein-Dor, E. Shnarch, N. Slonim, M. Shmueli-Scheuer, and L. Choshen. 2024. Efficient Benchmarking of Language Models. *arXiv preprint arXiv:2308.11696*, *Computation and Language (cs.CL)*, accepted to *NAACL v5* (2024), 1–19. <https://arxiv.org/abs/2308.11696>
- [69] Pavel Prudkov. 2025. On the construction of artificial general intelligence based on the correspondence between goals and means. *Frontiers in Artificial Intelligence* 8 (2025), 1588726. <https://www.frontiersin.org/articles/10.3389/frai.2025.1588726/full>
- [70] M. Pternea, P. Singh, A. Chakraborty, Y. Oruganti, M. Milletari, S. Bapat, and K. Jiang. 2024. The RL/LLM Taxonomy Tree: Reviewing Synergies Between Reinforcement Learning and Large Language Models. *Journal of Artificial Intelligence Research* 80 (2024). doi:10.1613/jair.1.15960
- [71] E. Raff, M. Benaroch, S. Samtani, and A. L. Farris. 2024. What Do Machine Learning Researchers Mean by “Reproducible”? *arXiv preprint arXiv:2412.03854*, To appear in *AAAI 2025, Senior Member Presentation Track* (2024). <https://arxiv.org/abs/2412.03854>
- [72] N. Ravi, A. Goel, J. C. Davis, and G. K. Thiruvathukal. 2025. Improving the Reproducibility of Deep Learning Software: An Initial Investigation through a Case Study Analysis. *arXiv preprint arXiv:2505.03165* (2025). <https://arxiv.org/abs/2505.03165>
- [73] Nicholas Riccardi, Xuan Yang, and Rutvik H. Desai. 2024. The Two Word Test as a semantic benchmark for large language models. *Scientific Reports* 14 (2024). doi:10.1038/s41598-024-72528-3
- [74] Michael Roberts, Derek Driggs, Matthew Thorpe, Julian Gilbey, Michael Yeung, Stephan Ursprung, Angelica I. Aviles-Rivero, Christian Etmann, Cathal McCague, Lucian Beer, Jonathan R. Weir-McCall, Zhongzhao Teng, Effrossyni Gkrania-Klotsas, AIX-COVNET, James H. F. Rudd, Evis Sala, and Carola-Bibiane Schönlieb. 2021. Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans. *Nature Machine Intelligence* 3, 3 (2021), 199–217. doi:10.1038/s42256-021-00307-0
- [75] Sascha Rothe, Shashi Narayan, and Aliaksei Severyn. 2020. Leveraging Pre-trained Checkpoints for Sequence Generation Tasks. *Transactions of the Association for Computational Linguistics* 8 (2020), 264–280. <https://aclanthology.org/2020.tacl-1.18/>
- [76] Thomas Savage, Ashwin Nayak, Robert Gallo, Ekanath Rangan, and Jonathan H. Chen. 2024. Diagnostic reasoning prompts reveal the potential for large language model interpretability in medicine. *npj Digital Medicine* 7 (2024), Article 20. <https://www.nature.com/articles/s41746-024-01010-1>
- [77] D. J. Schlueter, N. R. Bar, E. Shin, P. Chou, E. Winden, X. Zhou, J. Ramirez, K. Chu, N. Guller, B. Liang, H. E. Armour, J. H. Gilmour, and L. Bastarache. 2024. Systematic replication of smoking disease associations using survey responses and EHR data in the All of Us Research Program. *Journal of the American Medical Informatics Association* 31, 1 (2024), 139–150. <https://academic.oup.com/jamia/article/31/1/139/7330649>
- [78] H. Semmelrock, T. Ross-Hellauer, S. Kopeinik, D. Theiler, A. Haberl, S. Thalmann, and D. Kowald. 2025. Reproducibility in Machine Learning-based Research: Overview, Barriers and Drivers. *arXiv preprint arXiv:2406.14325*, Accepted for publication in *AI Magazine* (2025). <https://arxiv.org/abs/2406.14325>
- [79] Xin Shen, Wai Lam, Shumin Ma, and Huadong Wang. 2024. Joint learning of text alignment and abstractive summarization for long documents via unbalanced optimal transport. *Natural Language Engineering* 30, 3 (2024), 525–553. <https://www.cambridge.org/core/journals/natural-language-engineering/article/joint-learning-of-text-alignment-and-abstractive-summarization-for-long-documents-via-unbalanced-optimal-transport/46EF85C92B3E4158D89DC2C43E55D621>
- [80] Georgios Sidiropoulos, Samarth Bhargav, Panagiotis Eustratiadis, and Evangelos Kanoulas. 2025. Multivariate Dense Retrieval: A Reproducibility Study under a Memory-limited Setup. *Transactions on Machine Learning Research* 2025 (Jan 2025). <https://openreview.net/forum?id=rHmc5Y6ICg>
- [81] Michael A. Skinnider, R. Greg Stacey, David S. Wishart, and Leonard J. Foster. 2021. Chemical language models enable navigation in sparsely populated chemical space. *Nature Machine Intelligence* 3, 9 (2021), 759–770. <https://www.nature.com/articles/s42256-021-00368-1>
- [82] Mark Steyvers, Heliodoro Tejeda, Aakriti Kumar, Catarina Belem, Sheer Karny, Xinyue Hu, Lukas W. Mayer, and Padhraic Smyth. 2025. What large language models know and what people think they know. *Nature Machine Intelligence* 7 (Feb. 2025), 221–231. <https://www.nature.com/articles/s42256-024-00976-7>
- [83] Laura E. Suárez, Blake A. Richards, Guillaume Lajoie, and Bratislav Misic. 2021. Learning function from structure in neuromorphic networks. *Nature Machine Intelligence* 3, 9 (2021), 771–786. doi:10.1038/s42256-021-00376-1
- [84] J. Sublime. 2024. The AI Race: Why Current Neural Network-based Architectures are a Poor Basis for Artificial General Intelligence. *Journal of Artificial Intelligence Research* 80 (2024), 1165–1189. <https://jair.org/index.php/jair/article/view/15315/26999>
- [85] Craig Thomson, Ehud Reiter, and Anya Belz. 2024. Common Flaws in Running Human Evaluation Experiments in NLP. *Computational Linguistics* 50, 2 (2024), 795–805. doi:10.1162/coli_a_00508
- [86] Shushan Toneyan, Ziqi Tang, and Peter K. Koo. 2022. Evaluating deep learning for predicting epigenomic profiles. *Nature Machine Intelligence* 4, 12 (2022), 1088–1100. doi:10.1038/s42256-022-00570-9
- [87] P. Totis, J. Davis, L. de Raedt, and A. Kimmig. 2023. Lifted Reasoning for Combinatorial Counting. *Journal of Artificial Intelligence Research* 76, 14062 (2023), 1–58. doi:10.1613/jair.1.14062
- [88] Junichi Tsujii. 2021. Natural Language Processing and Computational Linguistics. *Computational Linguistics* 47, 4 (2021), 707–727. doi:10.1162/coli_a_00420
- [89] W. van Woerkom, D. Grossi, H. Prakken, and B. Verheij. 2024. A Fortiori Case-Based Reasoning: From Theory to Data. *Journal of Artificial Intelligence Research* 81 (2024), 1–38. doi:10.1613/jair.1.15178
- [90] L. Vaugrante, M. Niepert, and T. Hagendorff. 2024. A Looming Replication Crisis in Evaluating Behavior in Language Models? Evidence and Solutions. *arXiv preprint arXiv:2409.20303* (2024). <https://arxiv.org/abs/2409.20303>
- [91] P. Veličković and C. Blundell. 2021. Neural algorithmic reasoning. *Patterns* 2, 7 (2021), 100273. doi:10.1016/j.patter.2021.100273
- [92] A. Waldis, Y. Perlit, L. Choshen, Y. Hou, and I. Gurevych. 2024. Holmes A Benchmark to Assess the Linguistic Competence of Language Models. *Transactions of the Association for Computational Linguistics* 12 (2024), 1616–1647. <https://aclanthology.org/2024.tacl-1.88>

- [93] Bin Wang and C.-C. Jay Kuo. 2020. SBERT-WK: A Sentence Embedding Method by Dissecting BERT-Based Word Models. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28 (2020), 2146–2157. doi:10.1109/TASLP.2020.3007833
- [94] Hanyin Wang, Chufan Gao, Christopher Dantona, Bryan Hull, and Jimeng Sun. 2024. DRG-LLaMA : tuning LLaMA model to predict diagnosis-related group for hospitalized patients. *npj Digital Medicine* 7 (2024), Article 16. doi:10.1038/s41746-023-00989-3
- [95] Wenguan Wang, Yi Yang, and Fei Wu. 2024. Towards Data-And Knowledge-Driven AI: A Survey on Neuro-Symbolic Computing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2024). doi:10.1109/TPAMI.2024.3483273 Early Access.
- [96] Y. Wang, Y. Zhang, P. Li, and Y. Liu. 2024. Gradual Syntactic Label Replacement for Language Model Pre-Training. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 32 (2024), 961–972. doi:10.1109/TASLP.2023.3331096
- [97] Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The Benchmark of Linguistic Minimal Pairs for English. *Transactions of the Association for Computational Linguistics* 8 (2020), 377–392. <https://aclanthology.org/2020.tacl-1.25/>
- [98] C. Wei, K. Duan, S. Zhuo, H. Wang, S. Huang, and J. Liu. 2025. Enhanced Recommendation Systems with Retrieval-Augmented Large Language Model. *Journal of Artificial Intelligence Research* 82 (2025). doi:10.1613/jair.1.17809
- [99] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 35. 24824–24837. <https://arxiv.org/abs/2201.11903>
- [100] Chaoyi Wu, Pengcheng Qiu, Jinxin Liu, Hongfei Gu, Na Li, Ya Zhang, Yanfeng Wang, and Weidi Xie. 2025. Towards evaluating and building versatile large language models for medicine. *npj Digital Medicine* 8 (2025). doi:10.1038/s41746-024-01390-4
- [101] Xuenan Xu, Ziliang Xie, Mengyue Wu, and Kai Yu. 2023. Beyond the Status Quo: A Contemporary Survey of Multi-View Learning in Speech and Language Processing. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 32 (2023), 95–112. doi:10.1109/TASLP.2023.3321968
- [102] M. Yang, Y. Wang, and Y. Gu. 2025. Language-based reasoning graph neural network for commonsense question answering. *Neural Networks* 181 (Jan. 2025), 106816. doi:10.1016/j.neunet.2024.106816
- [103] S. W. Yang, H. J. Chang, Z. Huang, A. T. Liu, P. Su, W. Cheng, Y. Li, M. Wu, J. Lee, O. Hussein, M. Maciejewski, X. Zeng, C. H. Chen, Y. Tsao, D. Su, P. Beh, P. Zhang, Y. Shinohara, F. Weninger, F. Ni, S. Watanabe, T. Hori, A. Subramanian, K. K. Chin, P. Garcia-Perera, M. L. Seltzer, and H. Y. Lee. 2024. A Large-Scale Evaluation of Speech Foundation Models. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 32 (2024), 2884–2899. <https://ieeexplore.ieee.org/document/10502279>
- [104] T. Zhang, F. Ladhak, E. Durmus, P. Liang, K. McKeown, and T. B. Hashimoto. 2024. Benchmarking Large Language Models for News Summarization. *Transactions of the Association for Computational Linguistics* 12 (2024), 45–62. <https://aclanthology.org/2024.tacl-1.4.pdf>
- [105] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. *arXiv preprint arXiv:2303.18223 [cs.CL]* (2023). <https://arxiv.org/abs/2303.18223>
- [106] Yinhe Zheng, Guanyi Chen, and Minlie Huang. 2020. Out-of-Domain Detection for Natural Language Understanding in Dialog Systems. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28 (2020), 1198–1207. <https://ieeexplore.ieee.org/document/9052492>
- [107] Yizhen Zheng, Huan Yee Koh, Jiaxin Ju, Anh T. N. Nguyen, Lauren T. May, Geoffrey I. Webb, and Shirui Pan. 2025. Large language models for scientific discovery in molecular property prediction. *Nature Machine Intelligence* 7 (March 2025), 437–447. doi:10.1038/s42256-025-00994-z
- [108] J. Zhou, W. Zhong, Y. Wang, and J. Wang. 2025. Adaptive-solver framework for dynamic strategy selection in large language model reasoning. *Information Processing & Management* 62, 3 (2025), Article 104052. <https://www.sciencedirect.com/science/article/pii/S0306457324000468>