Original Research



Benchmarking Large Language Models on Diagnostic Inference Tasks in Medical Texts

Rendra Wirawan¹

¹Department of Computer Science, Universitas Teknologi Nusa Jaya, Jl. Cempaka Desa Karangjati, Kabupaten Temanggung, Jawa Tengah, Indonesia.

Abstract

The exponential growth of large language models has led to extensive research aimed at evaluating their capabilities for various specialized tasks, particularly in fields where interpretive clarity and diagnostic accuracy are of utmost importance. In medical contexts, the capacity to engage in diagnostic inference relies on multiple interconnected factors, including the ability to parse symptoms, correlate them with potential conditions, and address the nuances of domain-specific language. This paper explores the benchmarking of large language models on diagnostic inference tasks in medical texts, focusing on their performance when tasked with identifying complex disease processes and recommending appropriate clinical interventions. By systematically comparing several leading models, we aim to discern how their learned representations handle synonymy, polysemy, and context-dependent cues critical in medical discourse. Through a robust quantitative approach, our assessment encompasses both standard measures of precision and recall as well as more advanced evaluation metrics that capture the interpretive subtlety required by clinical practitioners. Furthermore, we present analytical perspectives centered on logical consistencies, semantic transparency, and cross-domain adaptability, evaluating the ability of these models to generalize to diverse clinical scenarios. Our results highlight key challenges and emergent strengths in the realm of automated medical reasoning, underscoring potential paths toward advancing large language models to robustly support real-world diagnostic workflows. The findings outlined herein may serve as a foundational basis for future research directed at integrating sophisticated inference mechanisms into medical text processing pipelines.

1. Introduction

The development of large language models has revolutionized natural language processing and propelled an abundance of applications across numerous fields, including medical informatics [1]. With large-scale pretraining on vast corpora, such models often demonstrate remarkable language understanding capabilities that extend beyond simple keyword matching. Nevertheless, the specific demands of diagnostic inference in medical texts call for more targeted analyses than those typically performed on generalpurpose language tasks [2]. Medical practitioners not only rely on precisely curated terminologies but also incorporate subtle variations in context and textual cues to differentiate among multiple overlapping conditions. As a result, benchmarks specifically designed to capture these complexities can offer valuable insights into models' limitations and guide improvements that align with real clinical needs.

Evaluation methodologies in this domain frequently span multiple paradigms, ranging from simple text classification to more sophisticated reasoning tasks [3]. For instance, tasks may require a model to read a patient report describing a constellation of symptoms and then deduce the most probable diagnosis. The central hypothesis is that large language models, due to their extensive training, may already capture a variety of linguistic and statistical patterns beneficial to diagnostic reasoning [4]. However, the complexities inherent in medical texts—such as domain-specific jargon, incomplete information, and specialized terminologies—pose challenges that are not always evident in more generic evaluation datasets. In addition, the multifaceted nature of patient data, which can contain laboratory results,

imaging findings, or historical context, introduces a further layer of complexity for models largely trained on general textual information.

Recent studies have highlighted the limitations of large language models in contexts that require robust interpretability [5]. Interpretability is particularly crucial in healthcare, where the precise reasoning behind a diagnosis or treatment recommendation can significantly affect patient outcomes. When large language models generate results that lack transparency, medical experts may be hesitant to incorporate those results into decision-making processes, especially in high-stakes scenarios. Consequently, substantial effort has been dedicated to the design of evaluation protocols aimed at probing whether models align with clinical reasoning principles [6]. Instead of asking whether a model can label a text, researchers increasingly ask whether the model can justify its decisions in a manner consistent with expert-level understanding.

Alongside concerns about interpretability, there is an ongoing debate over the generalizability of large language models within specialized fields [7, 8]. Although pretraining on massive and diverse text corpora endows models with broad linguistic coverage, the specific terminologies and contextual nuances of specialized disciplines often require additional fine-tuning. In medical applications, the differential diagnosis process alone can involve hundreds of nuanced terms, each with a particular set of associated risk factors, comorbidities, and treatment protocols. Overlooking this complexity can lead to an underestimation of the true difficulty of diagnostic inference tasks [9]. Accordingly, the benchmarks employed must be sufficiently rigorous and reflect real-world data to ensure meaningful results.

In parallel, there is also a growing interest in advanced machine learning architectures that combine the strengths of large language models with external knowledge sources [10]. This can take the form of knowledge graphs, ontologies, or even rule-based expert systems designed to guide reasoning steps. One line of inquiry involves integrating medical knowledge bases to complement the model's learned representations, potentially enabling more accurate or interpretable predictions. The synergy between latent semantic representations of text and structured domain knowledge promises improvements in both specificity and recall, though integrating these approaches presents technical and conceptual hurdles [11]. For instance, bridging the gap between unstructured textual embeddings and structured knowledge constraints often requires sophisticated alignment techniques and logic-based rule matching.

Another issue central to diagnostic inference is the role of uncertainty in medical texts. Clinicians frequently use vague or hedging language, indicating potential diagnoses or signaling the need for further testing [12]. A model that incorrectly interprets these nuances might provide overly confident, yet potentially inaccurate, recommendations. Consequently, an ideal benchmark must include not only typical cases but also ambiguous examples reflecting realistic clinical ambiguities [13]. In this vein, probabilistic modeling frameworks can help quantify uncertainty, thereby reflecting a more accurate portrayal of clinical reasoning processes. The capacity to encode and propagate uncertainty is thus a critical dimension in the evaluation and comparison of large language models.

Moreover, the structure of medical documentation itself can influence diagnostic reasoning [14]. Clinicians commonly rely on standardized reporting forms, annotated reports, or integrated information from laboratory results. The extent to which a language model can parse this diversity of data formats directly affects its capacity to draw accurate inferences. Data curation and preprocessing strategies, therefore, become essential components of any benchmarking effort [15]. If the training or evaluation data poorly reflect authentic clinical settings, the resulting performance estimates might mislead researchers or clinicians about how well these models handle real-world complexities.

Ethical and regulatory frameworks also play a pivotal role in shaping the evaluation of large language models for medical use [16]. Patient privacy concerns often limit the availability of robust datasets, leading to an overreliance on synthetic or de-identified records. While these methods provide a stopgap to facilitate experimentation, they may not capture the full intricacies of authentic patient narratives, particularly those that hinge on sensitive socio-demographic factors. Additionally, the high stakes involved in medical practice necessitate an additional layer of scrutiny that goes beyond standard benchmarks [17]. This includes external validation by clinical experts, prospective testing in real clinical workflows, and ongoing monitoring of model outputs for potential biases or oversights.

Even as we acknowledge these challenges, there remains reason for optimism [18]. Advancements in model architectures, computational resources, and data availability can collectively drive continued improvements in the predictive and inferential power of large language models. By subjecting these models to stringent and contextually relevant benchmarks, the field can systematically identify shortcomings and innovations alike. This paper takes a step in that direction by focusing on diagnostic inference tasks in medical texts, providing an extensive comparative analysis of models, while also highlighting areas where future developments could lead to more reliable and actionable automated reasoning systems. [19]

The structure of this work encompasses the application of rigorous evaluation metrics tailored to diagnostic reasoning, consideration of domain-specific constraints such as specialized nomenclature and uncertain evidence, and an exploration of potential synergistic approaches that marry data-driven learning with structured medical knowledge. By doing so, we aim not only to gauge current performance levels but also to chart a course for future research that addresses the unique demands of clinical practice. The following sections detail our methodology, present our experimental results, delve into interpretive aspects of model behavior, and ultimately culminate in conclusions that outline both the limitations and the promise of using large language models for diagnostic tasks. [20]

2. Related Work

Research into applying advanced language models to the medical domain has evolved significantly over the past decade. Early endeavors primarily revolved around shallow machine learning approaches that relied on carefully engineered feature sets derived from textual cues like n-grams, part-of-speech tags, and domain-specific lexica [21]. These methods showed moderate success in tasks such as detecting specific diseases from unstructured clinical notes. However, the transition to deep learning and, more recently, to large language models has drastically altered the landscape. Models such as those built on the Transformer architecture have proven exceptionally adept at capturing contextual and semantic relationships, thereby surpassing conventional machine learning techniques in benchmark comparisons. [22]

Recent attempts to measure progress in this field have often centered on curated benchmark datasets that represent a fraction of real clinical scenarios. Despite their utility, many of these datasets fail to reflect the full spectrum of variability found in genuine patient narratives. Additionally, the inherent complexity of diagnostic inference, which combines textual pattern recognition with domain-specific knowledge, remains a significant challenge for these models [23]. Some researchers have introduced synthetic datasets with controlled variability to isolate specific phenomena or linguistic patterns. While these can yield insights into model behaviors, their ecological validity is sometimes questioned, especially when generalizing to broader clinical contexts. [24]

A growing body of work explores the incorporation of external knowledge sources to enhance performance. Efforts include the integration of Unified Medical Language System (UMLS) ontologies and other medical knowledge bases that encode hierarchical relationships among diseases, symptoms, and treatments. By aligning the model's latent representations with structured concept embeddings, researchers aim to achieve a more robust form of semantic understanding [25]. In parallel, the introduction of logic-based rules, which specify constraints such as "if symptom A and symptom B are present, then condition C is more likely," has also been explored. Such hybrid approaches have shown promise in improving both model performance and interpretability. [26]

Investigators have also probed the generalizability of models trained on publicly available medical corpora, such as PubMed abstracts and clinical guidelines. While these sources can enrich the textual understanding of rare conditions, the practical benefit for diagnostic inference tasks remains mixed. The specialized jargon found in academic publications does not always map directly onto the descriptive, and sometimes incomplete, language common in patient reports [27]. Additionally, domain adaptation techniques, which involve fine-tuning a general-purpose large language model on domain-specific text, have been proposed to mitigate some of these gaps. In practice, the effectiveness of domain adaptation can vary, depending on factors such as dataset size, diversity, and annotation quality.

Parallel to methodological work, discussions have increasingly focused on ethical and regulatory aspects [28]. Scholars note that large language models, despite high performance metrics, can still reproduce biases from their training data, a concern amplified in healthcare settings where such biases can propagate into life-affecting decisions. The notion of "algorithmic accountability" demands rigorous evaluation protocols that delve into model outputs, ensuring they do not perpetuate health disparities or misrepresent demographic groups [29]. Mechanisms for continuous monitoring and feedback loops from clinical experts are often proposed as partial remedies. Nevertheless, robust solutions that fully mitigate bias remain an active area of research.

Another relevant direction pertains to interpretability [30]. Several studies highlight the tension between the complexity of large language models and the need for transparent, justifiable recommendations in clinical practice. Researchers have proposed attention-based visualization methods, gradient-based saliency maps, and post-hoc explanation techniques to offer clinicians some glimpse into the model's decision pathways. Yet, these methods can sometimes present oversimplified views of internal reasoning processes [31]. In diagnostic inference tasks, providing a merely plausible rationale may be insufficient if the reasoning is not genuinely grounded in medical logic. This growing realization pushes the field to seek interpretability solutions that combine transparency with genuine alignment to clinical reasoning standards. [32]

Moreover, current evaluations frequently rely on standard classification metrics like F1 score, accuracy, and Area Under the Curve (AUC). While these are valuable, they may not fully capture the practical intricacies of diagnostic inference. For example, a differential diagnosis task may involve multiple partially correct answers, and the relative ranking of potential conditions can be as important as a top-1 prediction [33]. Consequently, recent studies have proposed more nuanced metrics, including coverage error, ranking loss, and stepwise logical consistency. Some efforts even integrate cost-sensitive evaluations, reflecting the real-world implications of missed diagnoses versus false positives [34]. These refinements highlight the increasing sophistication in how researchers conceptualize and measure the impact of model outputs.

Certain lines of research investigate long-context models that can handle extensive narrative inputs, such as an entire patient file spanning multiple visits. These models aim to capture the evolving clinical picture over time, tracking changes in symptoms, treatments, and test results [35]. In doing so, they open the possibility for more dynamic forms of inference, resembling the iterative reasoning clinicians conduct as they gather more evidence. Yet, the computational demands for processing long sequences remain significant, and effective truncation or summarization strategies must be developed so as not to lose critical information.

In summary, the body of related work reflects a multifaceted exploration of how large language models can be optimized or adapted for diagnostic inference tasks in medical texts [36]. By synthesizing domain-specific knowledge, interpretability techniques, advanced evaluation metrics, and ethically oriented frameworks, this research trajectory is steadily moving toward models that are better aligned with the real-world demands of clinical decision-making. The present study builds upon these efforts by offering a holistic benchmark suite, incorporating both curated and near-real-world datasets, as well as evaluating whether logic-based constraints and external knowledge integration can further enhance performance [37]. This work aims to provide an up-to-date perspective on the strengths and limitations of cutting-edge models, bridging methodological gaps and highlighting areas for future research in the quest for robust, reliable diagnostic inference.

3. Methodology

The core of our methodology resides in establishing a comprehensive framework to evaluate large language models on diagnostic inference tasks in medical texts. We begin with a formal definition of the problem domain [38]. Let the input space be represented by strings denoted as S, where each element $s \in S$ may correspond to a patient case description, a clinical vignette, or any relevant textual record

containing diagnostic information. Our objective is to learn a function $f : S \to D$, mapping each s to a set D of potential diagnoses. [39]

More precisely, we define a structured representation of a patient record as r = (p, c), where p captures patient demographics and history, and c includes current symptoms, lab findings, and any available imaging data. The model aims to output the correct diagnosis or a ranked list of likely diagnoses, denoted by $d \in D$.

Given a set of training examples $\{(r_i, d_i)\}_{i=1}^N$, each pair (r_i, d_i) is assumed to be drawn from an unknown distribution consistent with real clinical scenarios. Our methodology accounts for the possibility that a record r_i can have multiple correct diagnoses [40]. We formalize this multi-label scenario using an indicator function $I(d_i)$, which takes the value 1 if diagnosis d_i is clinically valid for r_i , and 0 otherwise. The model is penalized both for failing to retrieve correct diagnoses and for suggesting diagnoses that are irrelevant to r_i .

A key aspect of our framework involves logic-based constraints that encode clinical knowledge [41]. Specifically, we introduce constraints of the form:

 $\forall x \in R$, Symptom(x) \land RiskFactor(x) \rightarrow HighProbability(Disease(x)),

which indicate the conditions under which certain diagnoses become highly probable [42]. These constraints are integrated during training or inference to guide the model toward outputs that are consistent with domain expertise.

In practice, we incorporate these constraints via an additional loss term, denoted L_{logic} , which imposes a penalty whenever the model's predictions violate established medical rules. To balance data-driven learning and logical reasoning, we define a combined objective function:

$$\mathcal{L} = \alpha L_{\text{data}} + \beta L_{\text{logic}},$$

where L_{data} is the standard cross-entropy loss for classification or ranking tasks, and α , β are hyperparameters controlling the influence of each component. Optimization proceeds via gradient-based methods, with α and β selected through cross-validation. [43]

We evaluate a diverse set of large language models, ranging from those trained on general-domain corpora to models fine-tuned on biomedical literature. Let M_{θ} denote a parameterized model, where θ represents the model parameters. We consider specific instances $M_{\theta}^1, M_{\theta}^2, \ldots, M_{\theta}^k$, each corresponding to a distinct pretraining or fine-tuning scheme. [44]

Our experimental pipeline includes generating tokenized representations for each clinical record, using either subword tokenization or domain-specific vocabularies to preserve semantic granularity. Additionally, we introduce a positional encoding scheme designed to highlight the importance of clinical keywords such as "pain," "fever," and "imaging findings." This augmented representation enables the model to better capture the contextual nuances of medical language. [45]

For the linear algebraic foundation, let X be an embedding matrix of dimension $m \times n$, where m corresponds to the sequence length of the tokenized record and n is the embedding dimension. A standard Transformer-based model projects this embedding matrix into multiple attention heads, generating context-aware representations. To incorporate structured knowledge, we extend each token embedding with a knowledge embedding vector $k_i drawn from an external resource (e.g., a conceptembed ding trained on a medical ontology) [46]. Hence, the$ $<math>x_i k_i$, where denote svector concatenation. The resulting matrix X^c ombined can be dimensionally reduced usin [47]

$$X_{\text{proj}} = X_{\text{combined}}W,$$

which the model then processes through a series of self-attention layers. The outcome is a final state representation that aims to encode both linguistic context and domain-specific knowledge. Subsequent feed-forward and classification layers translate these representations into probabilities over possible diagnoses. [48]

Finally, we address the practical evaluation of uncertainty in diagnostic inference. We propose а Bayesian approximation for the model's output, where ${
m M}_{isreplacedby M_{*},with sampled from a distribution reflecting parameter uncertainty. We can then compute a posterior predictive terms of the second seco$

$$p(d \mid r) = \int p(d \mid r, +)p() d.[49]$$

In practice, we approximate this integral using Monte Carlo dropout or alternative variational inference techniques. This step allows us to derive measures of confidence in the predicted diagnoses, directly reflecting the inherent ambiguity in many medical cases. [50]

Our methodology thus integrates data-driven learning, knowledge-constrained optimization, and uncertainty quantification to create a holistic approach to benchmarking large language models on diagnostic inference tasks. By combining these elements, we aim to shed light on both the potential and limitations of current state-of-the-art models. The next section will outline the experimental design used to implement and test these methodological innovations, followed by a quantitative and qualitative analysis of the results. [51]

4. Experimental Setup and Results

The experimental setup is engineered to provide a thorough evaluation of how well large language models perform in diagnostic inference tasks. We compile multiple medical datasets to capture the varied nature of real-world clinical documentation. These include publicly available collections of de-identified clinical notes, specialized corpora covering specific pathologies, and synthetic data generated to focus on particularly challenging linguistic constructs [52]. We implement a standardized preprocessing pipeline, which includes entity recognition for patient demographics, standardization of vital signs, and detection of negations in textual descriptions. Each dataset is partitioned into training, validation, and test sets, maintaining realistic distributions across conditions and demographics. [53]

We train multiple models ranging from generic large language models pretrained on web-scale data to domain-focused variants finetuned on biomedical text. For each model variant, we fix a maximum sequence length of 512 tokens, reflecting the typical length of a clinical vignette or patient note. Longer documents are segmented, ensuring that clinically relevant context remains largely intact [54]. Hyperparameters like learning rate and batch size are optimized through random search, with separate runs conducted for each model to accommodate differences in parameter space. To mitigate overfitting, we employ early stopping criteria tied to the validation set's performance, specifically monitoring improvements in F1 score for multi-label classification [55]. In the fine-tuning phase, each model typically converges within 5 to 10 epochs, depending on dataset complexity and size.

Evaluation involves multiple metrics to capture different facets of diagnostic accuracy. First, we measure precision, recall, and F1 score, treating each diagnosis as an independent label [56]. This standard approach is supplemented by metrics like the Jaccard index to quantify the overlap in multi-label outputs. Additionally, we compute a ranking-based measure, mean reciprocal rank (MRR), which becomes relevant when the output is a ranked list of potential diagnoses. We also implement a cost-sensitive metric that penalizes missed critical diagnoses more than less severe misclassifications, reflecting the real-world consequences of diagnostic errors [57]. For example, missing a diagnosis of acute myocardial infarction should incur a higher penalty than overlooking a benign condition.

Our experimental findings provide insights into how different model architectures fare in the face of diverse clinical inputs [58]. The general-purpose models often exhibit strong language comprehension for non-technical portions of the text but falter when confronted with highly specialized medical jargon or obscure pathophysiological conditions. In contrast, the domain-focused models demonstrate greater proficiency in interpreting complex medical narratives, particularly those involving comorbidities.

Notably, we observe a performance gap in scenarios where the text contains multiple potential diagnoses [59]. Models lacking explicit logic or knowledge integration occasionally produce contradictory or semantically inconsistent outputs, such as simultaneously predicting both "acute appendicitis" and "resolved appendicitis" for the same patient case.

When we incorporate the logic-based constraints described in the previous section, we record a measurable improvement in both F1 scores and interpretive consistency. The penalization of outputs that contradict well-established medical facts appears to help the models maintain logical coherence across multi-diagnosis tasks [60]. We also evaluate the impact of structured knowledge embeddings by comparing two model variants: one that processes text data exclusively and another that appends ontology-based vectors to each token representation. The latter consistently outperforms the former on most benchmarks, suggesting that domain knowledge provides valuable context clues [61]. These clues can disambiguate certain conditions, such as distinguishing "Type 1 diabetes" from "Type 2 diabetes" based on risk factors and comorbidities present in the text.

In terms of confidence calibration, our Bayesian approximation approach yields probabilities that more closely match empirical frequencies. We measure the calibration error by comparing predicted probabilities of correct diagnoses with observed frequencies in the test set [62]. Models employing Monte Carlo dropout during inference typically achieve lower calibration error than deterministic variants. This result suggests that the representation of parameter uncertainty reduces overconfidence, a particularly desirable feature in medical applications [63]. Indeed, being able to identify ambiguous cases, where the model is less certain, can guide further clinical investigations or additional diagnostic tests.

A noteworthy finding is the variability in performance across diverse subgroups in the data. For instance, performance on pediatric cases often lagged behind that on adult cases, partly due to differences in language and clinical parameters [64]. Similarly, rare diseases posed difficulties for all models, even those augmented with external medical knowledge. This phenomenon highlights the limitations of datadriven learning, as such conditions often appear infrequently in training sets, making it challenging for models to develop a robust understanding of their textual patterns. We quantify these differences by stratifying performance metrics by subgroup, revealing areas where models may require specialized data augmentation or more nuanced handling of domain knowledge. [65]

Beyond quantitative measures, we also conduct qualitative analyses of model outputs. We inspect cases where the model assigned a high probability to diagnoses that domain experts considered unlikely [66]. Manual examination often reveals that the model latched onto misleading textual cues, such as the presence of a medication typically used for a specific disease, without recognizing that it was used off-label or had been discontinued for reasons unrelated to the patient's current complaint. These findings emphasize the importance of context in clinical text understanding and the potential value of robust narrative reasoning mechanisms that track temporal or causal relationships among medical events.

Collectively, these results underscore both the substantial progress large language models have made and the complexities that remain [67]. Models enhanced with domain-specific knowledge and logic constraints exhibit meaningful improvements, yet they still struggle with ambiguous or rare scenarios. Confidence calibration techniques yield practical benefits in identifying uncertain cases, an essential function for clinical use. Our next section explores these outcomes in depth, discussing how the interplay of data-driven methods, knowledge integration, and logic constraints shapes the emergent behaviors of large language models in diagnostic inference [68]. We place particular emphasis on interpretability and real-world applicability, aiming to inform future research directions in automated medical reasoning.

5. Discussion

The experimental results provide a layered view of the challenges and potential solutions associated with large language models in diagnostic inference tasks for medical texts [69]. Several themes emerge that warrant deeper investigation. First, the incorporation of structured knowledge into the neural architecture appears to substantially improve interpretive accuracy. By leveraging ontology-based embeddings, models gain an additional semantic dimension, aiding in the resolution of ambiguities frequently

encountered in clinical narratives [70]. These findings align with earlier work suggesting that domainspecific constraints can complement purely data-driven strategies. However, questions remain regarding the optimal strategies for embedding such knowledge, as naive concatenation can introduce redundant or extraneous information that might inflate computational overhead [71]

Another critical observation pertains to logical consistency in model predictions. Without explicit constraints, large language models may propose diagnoses that conflict with known medical facts, reflecting a purely associative approach rather than genuine reasoning. The logic-based penalty in our experiments proved effective in mitigating such contradictions, but the development of more granular and adaptable rule sets remains an open problem [72]. In practical clinical scenarios, numerous exceptions to general rules exist, and a rigid set of constraints may either over-penalize valid inferences or fail to capture important nuances. Balancing these trade-offs demands sophisticated mechanisms for dynamically adjusting constraint sets, potentially requiring more advanced forms of logic programming or knowledge graph traversal.

Confidence calibration emerged as another important factor, as accurate probabilistic estimates can be vital for high-stakes medical decisions [73]. Our Bayesian approximation approach, incorporating Monte Carlo dropout, demonstrated improved alignment between predicted probabilities and observed outcomes. However, this technique increases computational costs during inference, which may not be feasible in every clinical setting [74, 75]. Future research might explore more efficient variational inference techniques or specialized hardware accelerations to maintain real-time or near-real-time performance. Alternatively, approximate calibration methods that reduce the computational load without sacrificing too much accuracy could prove beneficial.

Despite these advances, the consistent underperformance on rare diseases reveals the data limitation challenges [76]. Large language models excel in identifying patterns that appear frequently in training data, but their inference for less common conditions remains error-prone. Oversampling strategies, synthetic data generation, or domain-adaptive pretraining are potential avenues to address this bottleneck [77]. However, each approach carries trade-offs. Synthetic data may inadvertently introduce artifacts that skew model behavior, while domain adaptation requires carefully curated datasets that still might not encompass every rare condition. Moreover, the ethical and regulatory constraints on sharing medical data limit the volume of diverse training sets [78]. Collaborative networks that facilitate secure, multi-institutional data sharing may alleviate this issue, although such collaborations necessitate robust privacy-preserving methods.

Interpretability is a recurring concern in real-world deployments. While attention visualization or gradient-based saliency can offer partial insights into model predictions, they do not always align with genuine medical reasoning [79]. Large language models may highlight relevant fragments in the text without demonstrating causal or inferential understanding. Development of more advanced explanation techniques that can articulate logical chains of thought, potentially by integrating formal logic representations, could yield more trust in clinical environments [80]. Nonetheless, building and validating such methods remains non-trivial, as they must not only demonstrate plausible reasoning paths but also adhere to medical best practices.

The implications of biases in model performance also merit attention. If a model consistently underperforms for certain demographic groups, it risks perpetuating health disparities [81]. The cause of such biases may range from imbalanced training data to underlying sociocultural factors affecting clinical reporting. Addressing bias requires systematic approaches to dataset composition, model auditing, and performance stratification across patient subpopulations. Ongoing dialogue between technologists, clinicians, and ethicists is crucial to ensure that improvements in automated medical reasoning do not come at the expense of equitable care. [82]

Scalability and integration into clinical workflows represent further frontiers. Even with state-of-theart hardware, large models can be computationally expensive, slowing down inference [83]. Local or on-device solutions with smaller model architectures may be necessary for resource-constrained healthcare facilities. Additionally, embedding these models into electronic health record systems involves robust interoperability standards. The design of application programming interfaces for model inference, data pre-processing modules, and real-time updates from various clinical data streams all require careful engineering and compliance with healthcare data regulations [84]. Achieving seamless integration will likely involve close collaboration between model developers, health information technology professionals, and clinicians.

Finally, while our study provides a benchmark-focused perspective, real-world application must also consider clinical trial validations [85]. A model that excels in controlled evaluations might still yield unpredictable behavior in live medical settings, where incomplete data, human error in data entry, and context-specific nuances abound. Ongoing monitoring and iterative improvement cycles become necessary. As part of this process, direct feedback from clinicians who use the model's outputs can inform targeted refinements, ensuring that the technology evolves in tandem with professional experience and ethical standards. [86]

In summary, the discussion reinforces that while large language models demonstrate considerable promise for diagnostic inference tasks, significant methodological, ethical, and operational challenges remain. The interplay of data-driven representation learning, structured domain knowledge, logical constraints, and robust evaluation metrics defines a rich space for future research. By addressing these challenges in a systematic and transparent manner, the field can move closer to reliable, interpretable, and equitable automated medical reasoning systems [87]. In the concluding section, we consolidate our findings and propose avenues for subsequent investigations, emphasizing the collaborative nature of progress in this domain.

6. Conclusion

The collective insights from this research underscore both the promise and the complexity inherent in deploying large language models for diagnostic inference tasks in medical texts [88]. A principal takeaway is the significance of domain-specific knowledge in augmenting raw neural representations. Models that integrate structured ontologies or logic-based constraints consistently show superior performance compared to their purely data-driven counterparts. This enhancement is particularly evident in contexts where the text contains ambiguous or overlapping symptoms, illustrating how supplemental medical expertise can steer model predictions toward clinically coherent outcomes. [89]

Yet, the attainment of robust performance across a broad spectrum of conditions—ranging from common to extremely rare—remains an elusive goal. Data scarcity, particularly in rare disease cases, continues to hamper generalization. Proposed strategies, such as generating synthetic examples or orchestrating large-scale collaborations for data pooling, highlight the scope for further innovation [90]. While these interventions can mitigate data constraints, they also introduce new considerations. Synthetic data risk introducing artifacts that might mislead the model, while complex collaborations necessitate stringent protocols to preserve patient privacy [91]. Addressing this dual challenge will likely demand an amalgamation of innovative data engineering, ethical oversight, and clinical validation.

Another critical dimension is interpretability. Although attention-based heatmaps and gradient analysis provide some transparency, they do not necessarily equate to genuine reasoning in the medical sense [92]. The potential integration of formal logic into the model's decision-making process holds promise for more trustworthy explanations. Nevertheless, even logic-based approaches face the possibility of oversimplifying the complexities that underlie clinical judgment [93]. As progress unfolds, the onus lies on researchers to develop explanation frameworks that neither compromise the nuanced nature of clinical care nor obscure the computational intricacies of deep neural architectures.

Real-world deployment considerations also surface prominently. While our benchmarks are intentionally designed to capture a wide range of diagnostic challenges, genuine healthcare environments present dynamic variables, including incomplete data entry, evolving patient statuses, and concurrent medical interventions [94]. Coupled with variations in clinical documentation practices among different healthcare providers, these factors demand that any automated solution be adaptable, continuously monitored, and rigorously updated. Feedback loops, wherein clinicians can annotate or correct a model's suggestions, could feed directly into continual learning paradigms, thereby refining performance over time. However, this iterative process must be carefully balanced with regulatory standards governing software as a medical device, clinical trial validations, and institutional guidelines. [95]

An associated dimension is computational feasibility. The fine-tuning and inference steps for large language models can demand substantial computational resources, which may not be accessible in certain healthcare settings [96]. The optimization of model architectures for efficiency, possibly through parameter pruning, knowledge distillation, or specialized hardware acceleration, stands as a vital area of research. Sustainable deployments that accommodate different resource environments can make automated diagnostic inference tools more universally available, potentially democratizing access to advanced clinical decision support.

The inclusion of uncertainty estimates via Bayesian approximation or other calibration techniques is yet another valuable aspect of our findings [97]. Clinicians commonly encounter situations with incomplete or conflicting data, making it critical for computational systems to signal their own uncertainty. While our experiments demonstrate the feasibility of incorporating such techniques, their computational overhead and the complexity of interpreting probabilistic outputs in a clinical context must be addressed. Moreover, translating model confidence scores into practical recommendations for additional tests or referrals will require collaboration among domain experts, statisticians, and user-interface designers, ensuring that these probabilistic signals are both actionable and comprehensible. [98]

Additionally, the ethical and social considerations discussed throughout this paper remain vital to future progress. Biases intrinsic to training data can lead to unequal performance across patient demographics, potentially entrenching healthcare disparities [99]. Systematic audits, performance stratification, and ongoing refinement of data collection protocols can collectively mitigate these risks. Equitable representation in data, together with continuous vigilance from a multidisciplinary research community, holds the key to preventing harmful biases from becoming entrenched in diagnostic tools.

The present study lays a foundation for future research directions that can deepen and broaden the insights gained here [100]. One avenue lies in developing richer frameworks for real-time data integration, enabling models to update diagnostic suggestions as new information surfaces during patient care. Another potential path focuses on multimodal data, combining text with images, lab results, and genetic information [101]. The synergy of these sources has the potential to transform diagnostic accuracy, but it also compounds the technical and interpretive challenges. Finally, frameworks that incorporate continual learning while preserving patient confidentiality offer an exciting domain where large language models can adapt to evolving clinical knowledge over extended periods.

While considerable hurdles remain, this work illuminates the evolving capabilities of large language models to meet the stringent demands of diagnostic reasoning [102]. By methodically combining domain-specific knowledge, robust evaluation metrics, and interpretability features, our investigation illustrates a viable path for pushing the boundaries of automated medical inference. In so doing, it underscores the importance of collaboration between machine learning researchers, clinicians, ethicists, and policymakers. Only through a concerted, interdisciplinary effort can we harness the full potential of these powerful computational engines, bringing them closer to safe, equitable, and efficacious deployment in healthcare systems worldwide. [103]

The results and analyses presented here contribute to an ongoing discourse on the future of artificial intelligence in medicine, highlighting both the remarkable progress made and the complexity still to be unraveled. As large language models continue to evolve in sophistication, their utility for diagnostic inference tasks will likely expand, provided that remaining gaps in data availability, interpretability, and unbiased performance are addressed with due diligence [104]. By tracing a path forward that recognizes technical innovations, real-world viability, and ethical imperatives, we hope this work serves as a constructive reference point for researchers and practitioners seeking to refine and responsibly apply automated reasoning systems in clinical practice.

Closing this discussion, it is evident that the domain of medical text processing for diagnostic inference stands at a pivotal juncture. The synergy of advanced model architectures, formal logic constraints, and carefully curated datasets has enabled significant strides in accuracy and consistency

[105]. Yet, these achievements are merely precursors to a more profound shift in how clinicians interact with computational intelligence. As emerging techniques overcome current limitations, large language models will hold increasing relevance for the next generation of diagnostic decision support systems. By situating our findings within this broader trajectory, we invite further inquiry into strategies that can systematically integrate knowledge, transparency, and adaptability into the digital frameworks of modern healthcare, ultimately fostering improvements in patient outcomes and clinical efficiency on a global scale. [106]

References

- A. Mahmoud and N. Niu, "On the role of semantics in automated requirements tracing," *Requirements Engineering*, vol. 20, pp. 281–300, 1 2014.
- [2] K. B. Wagholikar, M. Torii, S. Jonnalagadda, and H. Liu, "Pooling annotated corpora for clinical concept extraction," *Journal of biomedical semantics*, vol. 4, pp. 3–3, 1 2013.
- [3] S. Lee, J. Han, R. W. Park, G. J. Kim, J. H. Rim, J. Cho, K. H. Lee, J. Lee, S. Kim, and J. H. Kim, "Development of a controlled vocabulary-based adverse drug reaction signal dictionary for multicenter electronic health record-based pharmacovigilance," *Drug safety*, vol. 42, pp. 657–670, 1 2019.
- [4] V. Garla, V. L. Re, Z. Dorey-Stein, F. Kidwai, M. Scotch, J. A. Womack, A. C. Justice, and C. Brandt, "The yale ctakes extensions for document classification: architecture and application," *Journal of the American Medical Informatics Association : JAMIA*, vol. 18, pp. 614–620, 5 2011.
- [5] S. M. Castro, E. Tseytlin, O. Medvedeva, K. J. Mitchell, S. Visweswaran, T. Bekhuis, and R. S. Jacobson, "Automated annotation and classification of bi-rads assessment from radiology reports," *Journal of biomedical informatics*, vol. 69, pp. 177–187, 4 2017.
- [6] S. Ramgopal, L. N. Sanchez-Pinto, C. M. Horvat, M. S. Carroll, Y. Luo, and T. A. Florin, "Artificial intelligence-based clinical decision support in pediatrics.," *Pediatric research*, vol. 93, pp. 334–341, 7 2022.
- [7] R. M. Shah, C. Wong, N. C. Arpey, A. A. Patel, and S. N. Divi, "A surgeon's guide to understanding artificial intelligence and machine learning studies in orthopaedic surgery." *Current reviews in musculoskeletal medicine*, vol. 15, pp. 121–132, 2 2022.
- [8] J. R. Machireddy, "Automation in healthcare claims processing: Enhancing efficiency and accuracy," *International Journal of Science and Research Archive*, vol. 09, no. 01, pp. 825–834, 2023.
- [9] C. Shorten, T. M. Khoshgoftaar, and B. Furht, "Deep learning applications for covid-19," *Journal of big data*, vol. 8, pp. 1–54, 1 2021.
- [10] C. Li and W. Xing, "Natural language generation using deep learning to support mooc learners," International Journal of Artificial Intelligence in Education, vol. 31, pp. 186–214, 1 2021.
- [11] P. H. Yi, T. K. Kim, and C. T. Lin, "Comparison of radiologist versus natural language processing-based image annotations for deep learning system for tuberculosis screening on chest radiographs.," *Clinical imaging*, vol. 87, pp. 34–37, 4 2022.
- [12] S. H. Fenton, D. T. Marc, A. Kennedy, D. Hamada, R. Hoyt, K. Lalani, C. Renda, and R. B. Reynolds, "Aligning the american health information management association entry-level curricula competencies and career map with industry job postings: Cross-sectional study.," *JMIR medical education*, vol. 8, pp. e38004–e38004, 7 2022.
- [13] C. Zheng, B. C. Sun, Y.-L. Wu, M. Ferencik, M.-S. Lee, R. F. Redberg, A. A. Kawatkar, V. V. Musigdilok, and A. L. Sharp, "Automated abstraction of myocardial perfusion imaging reports using natural language processing.," *Journal of nuclear cardiology : official publication of the American Society of Nuclear Cardiology*, vol. 29, pp. 1–10, 11 2020.
- [14] Z. Zeng, L. Yao, A. Roy, X. Li, S. Espino, S. E. Clare, S. A. Khan, and Y. Luo, "Identifying breast cancer distant recurrences from electronic health records using machine learning," *Journal of healthcare informatics research*, vol. 3, pp. 283–299, 4 2019.
- [15] J. Joseph, C. Liu, Q. Hui, K. Aragam, Z. Wang, B. Charest, J. E. Huffman, J. M. Keaton, T. L. Edwards, S. Demissie, L. Djousse, J. P. Casas, J. M. Gaziano, K. Cho, P. W. F. Wilson, L. S. Phillips, null null, C. J. O'Donnell, and Y. V. Sun, "Genetic architecture of heart failure with preserved versus reduced ejection fraction.," *Nature communications*, vol. 13, pp. 7753–, 12 2022.

- [16] R. McLaughlin, J. D. Nguyen, L. S. Brady, S. Eremenc, R. S. Keefe, M. H. Trivedi, M. Sand, K. Koblan, E. Stafford, R. Streiff, S. H. Lisanby, and C. A. Denny, "Acnp 61st annual meeting: Panels, mini-panels and study groups,," *Neuropsychopharmacology : official publication of the American College of Neuropsychopharmacology*, vol. 47, pp. 1–62, 12 2022.
- [17] L. Cao, Q. Yang, and P. S. Yu, "Data science and ai in fintech: an overview," International Journal of Data Science and Analytics, vol. 12, pp. 81–99, 8 2021.
- [18] V. X. Liu, M. P. Clark, M. Mendoza, R. R. Saket, M. N. Gardner, B. J. Turk, and G. J. Escobar, "Automated identification of pneumonia in chest radiograph reports in critically ill patients," *BMC medical informatics and decision making*, vol. 13, pp. 90–90, 8 2013.
- [19] S. A. Johnson, E. Signor, K. Lappe, J. Shi, S. L. Jenkins, S. W. Wikstrom, R. D. Kroencke, D. Hallowell, A. E. Jones, and D. M. Witt, "A comparison of natural language processing to icd-10 codes for identification and characterization of pulmonary embolism.," *Thrombosis research*, vol. 203, pp. 190–195, 5 2021.
- [20] N. Hong, A. Wen, D. J. Stone, S. Tsuji, P. R. Kingsbury, L. V. Rasmussen, J. A. Pacheco, P. Adekkanattu, F. Wang, Y. Luo, J. Pathak, H. Liu, and G. Jiang, "Developing a fhir-based ehr phenotyping framework: A case study for identification of patients with obesity and multiple comorbidities from discharge summaries," *Journal of biomedical informatics*, vol. 99, pp. 103310–103310, 10 2019.
- [21] N. G. Nia, E. Kaplanoglu, and A. Nasab, "Evaluation of artificial intelligence techniques in disease diagnosis and prediction," *Discover Artificial Intelligence*, vol. 3, 1 2023.
- [22] A. Porcalla, N. Barshteyn, S. Snyder, and M. Bhattacharya, "An innovative, collaborative, and strategic approach to proactively evaluate and update drug interactions based on prescribing information of newly approved medicinal products.," *Therapeutic innovation & regulatory science*, vol. 51, pp. 780–786, 5 2017.
- [23] M. Delsoz, H. Raja, Y. Madadi, A. A. Tang, B. M. Wirostko, M. Y. Kahook, and S. Yousefi, "The use of chatgpt to assist in diagnosing glaucoma based on clinical case reports.," *Ophthalmology and therapy*, vol. 12, pp. 3121–3132, 9 2023.
- [24] A. Javaid, M. A. Siddique, A. A. Reshi, null Mui-zzud din, F. Rustam, E. Lee, and V. Rupapara, "Coal mining accident causes classification using voting-based hybrid classifier (vhc)," *Journal of Ambient Intelligence and Humanized Computing*, vol. 14, pp. 13211–13221, 3 2022.
- [25] J. Whelan, M. Ghoniem, N. Médoc, M. Apicella, and E. Beck, "Applying a novel approach to scoping review incorporating artificial intelligence: mapping the natural history of gonorrhoea.," *BMC medical research methodology*, vol. 21, pp. 183–, 9 2021.
- [26] P. A. Heidenreich, "Can natural language processing fulfill the promise of electronic medical records," *Journal of cardiac failure*, vol. 20, pp. 465–466, 5 2014.
- [27] I. Ghanzouri, S. Amal, V. Ho, L. Safarnejad, J. Cabot, C. G. Brown-Johnson, N. Leeper, S. Asch, N. H. Shah, and E. G. Ross, "Performance and usability testing of an automated tool for detection of peripheral artery disease using electronic health records.," *Scientific reports*, vol. 12, pp. 13364–, 8 2022.
- [28] A. Sharma, Structural and network-based methods for knowledge-based systems. PhD thesis, Northwestern University, 2011.
- [29] P. Mathur, J. B. Cywinski, K. Maheshwari, J. Niezgoda, J. Mathew, C. C. do Nascimento, B. Abdelmalak, and F. Papay, "Automated analysis of ambulatory surgery patient experience comments using artificial intelligence for quality improvement: A patient centered approach," *Intelligence-Based Medicine*, vol. 5, pp. 100043–, 2021.
- [30] C. Jujjavarapu, V. Pejaver, T. Cohen, S. D. Mooney, P. J. Heagerty, and J. G. Jarvik, "A comparison of natural language processing methods for the classification of lumbar spine imaging findings related to lower back pain.," *Academic radiology*, vol. 29 Suppl 3, pp. S188–S200, 12 2021.
- [31] S. Yu, K. K. Kumamaru, E. George, R. M. Dunne, A. Bedayat, M. Neykov, A. R. Hunsaker, K. E. Dill, T. Cai, and F. J. Rybicki, "Classification of ct pulmonary angiography reports by presence, chronicity, and location of pulmonary embolism with natural language processing," *Journal of biomedical informatics*, vol. 52, pp. 386–393, 8 2014.
- [32] A. A. Allam, M. L. Nagy, G. R. Thoma, and M. Krauthammer, "Neural networks versus logistic regression for 30 days all-cause readmission prediction," *Scientific reports*, vol. 9, pp. 9277–9277, 6 2019.

- [33] J. C. Denny, A. Spickard, P. Speltz, R. Porier, D. Rosenstiel, and J. S. Powers, "Using natural language processing to provide personalized learning opportunities from trainee clinical notes," *Journal of biomedical informatics*, vol. 56, pp. 292–299, 6 2015.
- [34] M. C. Tremblay, D. J. Berndt, S. L. Luther, P. Foulis, and D. D. French, "Identifying fall-related injuries: Text mining the electronic medical record," *Information Technology and Management*, vol. 10, pp. 253–265, 11 2009.
- [35] R. T. Drumright, K. A. Regan, A. L. Lin, M. G. Moroux, and S. S. R. Iyer, "Utility of wound cultures in the management of open globe injuries: a 5-year retrospective review," *Journal of ophthalmic inflammation and infection*, vol. 10, pp. 1–5, 2 2020.
- [36] K. S. Allen, D. R. Hood, J. Cummins, S. Kasturi, E. A. Mendonca, and J. R. Vest, "Natural language processing-driven state machines to extract social factors from unstructured clinical documentation.," *JAMIA open*, vol. 6, pp. 00ad024–, 4 2023.
- [37] R. A. Denu and M. E. Burkard, "Analysis of the "centrosome-ome" identifies mcph1 deletion as a cause of centrosome amplification in human cancer," *Scientific reports*, vol. 10, pp. 11921–11921, 7 2020.
- [38] B. Shiner, L. W. D'Avolio, T. M. Nguyen, M. H. Zayed, Y. Young-Xu, R. A. Desai, P. P. Schnurr, L. D. Fiore, and B. V. Watts, "Measuring use of evidence based psychotherapy for posttraumatic stress disorder.," *Administration and policy in mental health*, vol. 40, pp. 311–318, 4 2012.
- [39] J. P. Ferraro, Y. Ye, P. H. Gesteland, P. J. Haug, F. R. Tsui, G. F. Cooper, R. E. V. Bree, T. Ginter, A. J. Nowalk, and M. M. Wagner, "The effects of natural language processing on cross-institutional portability of influenza case detection for disease surveillance," *Applied clinical informatics*, vol. 8, pp. 560–580, 5 2017.
- [40] L. Y. Jiang, X. C. Liu, N. P. Nejatian, M. Nasir-Moin, D. Wang, A. Abidin, K. Eaton, H. A. Riina, I. Laufer, P. Punjabi, M. Miceli, N. C. Kim, C. Orillac, Z. Schnurman, C. Livia, H. Weiss, D. Kurland, S. Neifert, Y. Dastagirzada, D. Kondziolka, A. T. M. Cheung, G. Yang, M. Cao, M. Flores, A. B. Costa, Y. Aphinyanaphongs, K. Cho, and E. K. Oermann, "Health system-scale language models are all-purpose prediction engines.," *Nature*, vol. 619, pp. 357–362, 6 2023.
- [41] C. Shivade, P. Raghavan, E. Fosler-Lussier, P. J. Embi, N. Elhadad, S. B. Johnson, and A. Lai, "A review of approaches to identifying patient phenotype cohorts using electronic health records," *Journal of the American Medical Informatics Association : JAMIA*, vol. 21, pp. 221–230, 11 2013.
- [42] M. Guo, Y. Ma, E. Eworuke, M. Khashei, J. Song, Y. Zhao, and F. Jin, "Identifying covid-19 cases and extracting patient reported symptoms from reddit using natural language processing.," *Scientific reports*, vol. 13, pp. 13721–, 8 2023.
- [43] A. Bansal, R. P. Padappayil, C. Garg, A. Singal, M. Gupta, and A. L. Klein, "Utility of artificial intelligence amidst the covid 19 pandemic: A review.," *Journal of medical systems*, vol. 44, pp. 156–156, 8 2020.
- [44] B. Hunter, S. Reis, D. Campbell, S. Matharu, P. Ratnakumar, L. Mercuri, S. Hindocha, H. Kalsi, E. Mayer, B. Glampson, E. J. Robinson, B. Al-Lazikani, L. Scerri, S. Bloch, and R. C. T. Lee, "Development of a structured query language and natural language processing algorithm to identify lung nodules in a cancer centre.," *Frontiers in medicine*, vol. 8, pp. 748168–, 11 2021.
- [45] R. Avula, "Strategies for minimizing delays and enhancing workflow efficiency by managing data dependencies in healthcare pipelines," *Eigenpub Review of Science and Technology*, vol. 4, no. 1, pp. 38–57, 2020.
- [46] L. Wang, L. Luo, Y. Wang, J. A. Wampfler, P. Yang, and H. Liu, "Ichi information extraction for populating lung cancer clinical research data," *Proceedings. IEEE International Conference on Healthcare Informatics*, vol. 2019, pp. 8904601–2, 11 2019.
- [47] R. Avula, "Healthcare data pipeline architectures for ehr integration, clinical trials management, and real-time patient monitoring," *Quarterly Journal of Emerging Technologies and Innovations*, vol. 8, no. 3, pp. 119–131, 2023.
- [48] M. Muniswamaiah, T. Agerwala, and C. Tappert, "Big data in cloud computing review and opportunities," arXiv preprint arXiv:1912.10821, 2019.
- [49] M. C. R. Melo, J. R. M. A. Maasch, and C. de la Fuente-Nunez, "Accelerating antibiotic discovery through artificial intelligence.," *Communications biology*, vol. 4, pp. 1050–1050, 9 2021.
- [50] K. D. Forbus, C. Riesbeck, L. Birnbaum, K. Livingston, A. Sharma, and L. Ureel, "Integrating natural language, knowledge representation and reasoning, and analogical processing to learn by reading," in *Proceedings of the National Conference on Artificial Intelligence*, vol. 22, p. 1542, Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 2007.

- [51] R. Bellazzi, M. Masseroli, S. N. Murphy, A. Shabo, and P. Romano, "Clinical bioinformatics: challenges and opportunities.," BMC bioinformatics, vol. 13, pp. 1–8, 9 2012.
- [52] M. Chaudhary, K. Kosyluk, S. Thomas, and T. Neal, "On the use of aspect-based sentiment analysis of twitter data to explore the experiences of african americans during covid-19,," *Scientific reports*, vol. 13, pp. 10694–, 7 2023.
- [53] N. B. Link, S. Huang, T. Cai, J. Sun, K. Dahal, L. Costa, K. Cho, K. Liao, T. Cai, C. Hong, and null null, "Binary acronym disambiguation in clinical notes from electronic health records with an application in computational phenotyping.," *International journal of medical informatics*, vol. 162, pp. 104753–104753, 4 2022.
- [54] B. Humm, H. Bense, J. Bock, M. Classen, O. Halvani, C. Herta, T. Hoppe, O. Juwig, and M. Siegel, "Applying machine intelligence in practice," *Informatik Spektrum*, vol. 43, pp. 137–144, 3 2020.
- [55] H. A. Piwowar and W. W. Chapman, "Amia identifying data sharing in biomedical literature.," AMIA ... Annual Symposium proceedings. AMIA Symposium, vol. 2008, pp. 596–600, 11 2008.
- [56] T. M. Ko, V. A. Kalesnikava, D. Jurgens, and B. Mezuk, "A data science approach to estimating the frequency of driving cessation associated suicide in the us: Evidence from the national violent death reporting system.," *Frontiers in public health*, vol. 9, pp. 689967–, 8 2021.
- [57] F. M. D. L. Vega, S. Chowdhury, B. Moore, E. Frise, J. McCarthy, E. J. Hernandez, T. C. Wong, K. N. James, L. Guidugli, P. B. Agrawal, C. A. Genetti, C. A. Brownstein, A. H. Beggs, B. S. Löscher, A. Franke, B. E. Boone, S. Levy, K. Õunap, S. Pajusalu, M. J. Huentelman, K. Ramsey, M. Naymik, V. Narayanan, N. Veeraraghavan, P. Billings, M. G. Reese, M. Yandell, and S. F. Kingsmore, "Artificial intelligence enables comprehensive genome interpretation and nomination of candidate diagnoses for rare genetic diseases," *Genome medicine*, vol. 13, pp. 153–, 10 2021.
- [58] T. K. Mackey, J. Li, V. Purushothaman, M. Nali, N. Shah, C. Bardier, M. Cai, and B. A. Liang, "Big data, natural language processing, and deep learning to detect and characterize illicit covid-19 product sales: Infoveillance study on twitter and instagram.," *JMIR public health and surveillance*, vol. 6, pp. e20794–, 8 2020.
- [59] C. Roth, R. E. Foraker, P. R. O. Payne, and P. J. Embi, "Community-level determinants of obesity: harnessing the power of electronic health records for retrospective data analysis," *BMC medical informatics and decision making*, vol. 14, pp. 36–36, 5 2014.
- [60] M. Muniswamaiah, T. Agerwala, and C. C. Tappert, "Context-aware query performance optimization for big data analytics in healthcare," in 2019 IEEE High Performance Extreme Computing Conference (HPEC-2019), pp. 1–7, 2019.
- [61] S. H. Chu, E. S. Wan, M. H. Cho, S. Goryachev, V. S. Gainer, J. G. Linneman, E. J. Scotty, S. J. Hebbring, S. N. Murphy, J. Lasky-Su, S. T. Weiss, J. W. Smoller, and E. W. Karlson, "An independently validated, portable algorithm for the rapid identification of copd patients using electronic health records.," *Scientific reports*, vol. 11, pp. 1–9, 10 2021.
- [62] L. Zhou, Y. Lu, C. J. Vitale, P. Mar, F. Y. Chang, N. Dhopeshwarkar, and R. A. Rocha, "Representation of information about family relatives as structured data in electronic health records," *Applied clinical informatics*, vol. 5, pp. 349–367, 4 2014.
- [63] A. Sharma, M. Witbrock, and K. Goolsbey, "Controlling search in very large commonsense knowledge bases: a machine learning approach," arXiv preprint arXiv:1603.04402, 2016.
- [64] H. Kang and Y. Gong, "Developing a similarity searching module for patient safety event reporting system using semantic similarity measures.," *BMC medical informatics and decision making*, vol. 17, pp. 75–75, 7 2017.
- [65] P. Suryanarayanan, C.-H. Tsou, A. Poddar, D. Mahajan, B. Dandala, P. Madan, A. Agrawal, C. Wachira, O. M. Samuel, O. Bar-Shira, C. Kipchirchir, S. Okwako, W. Ogallo, F. Otieno, T. Nyota, F. Matu, V. R. Barros, D. Shats, O. Kagan, S. L. Remy, O. Bent, P. Guhan, S. Mahatma, A. Walcott-Bryant, D. Pathak, and M. Rosen-Zvi, "Ai-assisted tracking of worldwide non-pharmaceutical interventions for covid-19.," *Scientific data*, vol. 8, pp. 94–94, 3 2021.
- [66] K. D. Forbus, K. Lockwood, A. B. Sharma, and E. Tomai, "Steps towards a second generation learning by reading system.," in AAAI Spring Symposium: Learning by Reading and Learning to Read, pp. 36–43, 2009.
- [67] A. Nashed, S. Zhang, C.-W. Chiang, M. Zitu, G. A. Otterson, C. J. Presley, K. Kendra, S. H. Patel, A. Johns, M. Li, M. Grogan, G. Lopez, D. H. Owen, and L. Li, "Comparative assessment of manual chart review and icd claims data in evaluating immunotherapy-related adverse events," *Cancer immunology, immunotherapy : CII*, vol. 70, pp. 2761–2769, 2 2021.

- [68] J. He, F. Li, X. Hu, J. Li, Y. Nian, J. Wang, Y. Xiang, Q. Wei, H. Xu, and C. Tao, "Chemical-protein relation extraction with pre-trained prompt tuning," *Proceedings. IEEE International Conference on Healthcare Informatics*, vol. 2022, pp. 608–609, 9 2022.
- [69] J. Zhang, S. Whebell, J. Gallifant, S. Budhdeo, H. Mattie, P. Lertvittayakumjorn, M. D. P. A. Lopez, B. J. Tiangco, J. W. Gichoya, H. Ashrafian, L. A. Celi, and J. T. Teo, "An interactive dashboard to track themes, development maturity, and global equity in clinical artificial intelligence research.," *The Lancet. Digital health*, vol. 4, no. 4, pp. e212–e213, 2022.
- [70] V. Chittajallu, E. Mansoor, J. Perez, Y. A. Omar, S. A. Firkins, H. Yoo, B. Baggott, and R. Simons-Linares, "Association of bariatric surgery with risk of incident obesity-associated malignancies: a multi-center population-based study.," *Obesity surgery*, vol. 33, pp. 4065–4069, 11 2023.
- [71] M. Muniswamaiah, T. Agerwala, and C. C. Tappert, "Federated query processing for big data in data science," in 2019 IEEE International Conference on Big Data (Big Data), pp. 6145–6147, IEEE, 2019.
- [72] H. Horng, J. Steinkamp, C. E. Kahn, and T. S. Cook, "Ensemble approaches to recognize protected health information in radiology reports," *Journal of digital imaging*, vol. 35, pp. 1694–1698, 6 2022.
- [73] E. Jung, H. Jain, A. P. Sinha, and C. Gaudioso, "Building a specialized lexicon for breast cancer clinical trial subject eligibility analysis," *Health informatics journal*, vol. 27, pp. 1460458221989392–1460458221989392, 2 2021.
- [74] M. Bhattacharya, S. Snyder, M. Malin, M. M. Truffa, S. Marinic, R. Engelmann, and R. R. Raheja, "Using social media data in routine pharmacovigilance: A pilot study to identify safety signals and patient perspectives," *Pharmaceutical Medicine*, vol. 31, pp. 167–174, 4 2017.
- [75] J. R. Machireddy, "Harnessing ai and data analytics for smarter healthcare solutions," *International Journal of Science and Research Archive*, vol. 08, no. 02, pp. 785–798, 2023.
- [76] M. Abouelyazid and C. Xiang, "Machine learning-assisted approach for fetal health status prediction using cardiotocogram data," *International Journal of Applied Health Care Analytics*, vol. 6, no. 4, pp. 1–22, 2021.
- [77] P. C. Johnson, N. Markovitz, T. F. Gray, S. Bhatt, R. D. Nipp, N. N. Ufere, J. Rice, M. J. Reynolds, M. W. Lavoie, C. E. Topping, M. A. Clay, C. Lindvall, and A. El-Jawahri, "Association of social support with overall survival and healthcare utilization in patients with aggressive hematologic malignancies.," *Journal of the National Comprehensive Cancer Network* : JNCCN, vol. -1, pp. 1–7, 10 2021.
- [78] R. Avula, "Applications of bayesian statistics in healthcare for improving predictive modeling, decision-making, and adaptive personalized medicine," *International Journal of Applied Health Care Analytics*, vol. 7, no. 11, pp. 29–43, 2022.
- [79] S. Singhal, B. Hegde, P. Karmalkar, J. Muhith, and H. Gurulingappa, "Weakly supervised learning for categorization of medical inquiries for customer service effectiveness.," *Frontiers in research metrics and analytics*, vol. 6, pp. 683400–, 8 2021.
- [80] E. Sholle, L. C. Pinheiro, P. Adekkanattu, M. A. Davila, S. B. Johnson, J. Pathak, S. Sinha, C. Li, S. A. Lubansky, M. M. Safford, and T. R. Campion, "Underserved populations with missing race ethnicity data differ significantly from those with structured race/ethnicity documentation.," *Journal of the American Medical Informatics Association : JAMIA*, vol. 26, pp. 722–729, 4 2019.
- [81] E. L. Palmer, J. H. Higgins, S. Hassanpour, J. D. Sargent, C. M. Robinson, J. A. Doherty, and T. Onega, "Assessing data availability and quality within an electronic health record system through external validation against an external clinical data source," *BMC medical informatics and decision making*, vol. 19, pp. 1–9, 7 2019.
- [82] M. Wu, Y. Zhang, M. Markley, C. Cassidy, N. Newman, and A. Porter, "Covid-19 knowledge deconstruction and retrieval: an intelligent bibliometric solution.," *Scientometrics*, vol. 129, pp. 1–7259, 5 2023.
- [83] A. Sharma, K. M. Goolsbey, and D. Schneider, "Disambiguation for semi-supervised extraction of complex relations in large commonsense knowledge bases," in 7th Annual Conference on Advances in Cognitive Systems, 2019.
- [84] Y. Wang, E. Willis, V. K. Yeruva, D. Ho, and Y. Lee, "A case study of using natural language processing to extract consumer insights from tweets in american cities for public health crises.," *BMC public health*, vol. 23, pp. 935–, 5 2023.
- [85] O. V. Patterson, M. S. Freiberg, M. Skanderson, S. J. Fodeh, C. Brandt, and S. L. DuVall, "Unlocking echocardiogram measurements for heart disease research through natural language processing.," *BMC cardiovascular disorders*, vol. 17, pp. 151–151, 6 2017.

- [86] M. A. Woodward, N. Maganti, L. M. Niziol, S. Amin, A. Hou, and K. Singh, "Development and validation of a natural language processing algorithm to extract descriptors of microbial keratitis from the electronic health record.," *Cornea*, vol. 40, pp. 1548–1553, 5 2021.
- [87] W. Chen, C. Durkin, Y. Huang, B. Adler, S. Rust, and S. Lin, "Simplified readability metric drives improvement of radiology reports: an experiment on ultrasound reports at a pediatric hospital.," *Journal of digital imaging*, vol. 30, pp. 710–717, 5 2017.
- [88] D. Newman-Griffis, J. Porcino, A. Zirikly, T. Thieu, J. C. Maldonado, P.-S. Ho, M. Ding, L. Chan, and E. K. Rasch, "Broadening horizons: the case for capturing function and the role of health informatics in its use," *BMC public health*, vol. 19, pp. 1–13, 10 2019.
- [89] R. W. Grundmeier, A. J. Masino, T. C. Casper, J. M. Dean, J. J. Bell, F. R. Enriquez, S. J. Deakyne, J. M. Chamberlain, and E. R. Alpern, "Identification of long bone fractures in radiology reports using natural language processing to support healthcare quality improvement," *Applied clinical informatics*, vol. 7, pp. 1051–1068, 11 2016.
- [90] A. Sorbello, S. A. Haque, R. Hasan, R. Jermyn, A. Hussein, A. Vega, K. Zembrzuski, A. Ripple, and M. Ahadpour, "Artificial intelligence-enabled software prototype to inform opioid pharmacovigilance from electronic health records: Development and usability study.," *JMIR AI*, vol. 2, pp. e45000–e45000, 7 2023.
- [91] Z. Chen, X. Cheng, S. Dong, Z. Dou, J. Guo, X. Huang, Y. Lan, C. Li, R. Li, T.-Y. Liu, Y. Liu, J. Ma, B. Qin, M. Wang, J.-R. Wen, J. Xu, M. Zhang, P. Zhang, and Q. Zhang, "Information retrieval: a view from the chinese ir community," *Frontiers of Computer Science*, vol. 15, pp. 151601–, 9 2020.
- [92] Y. Li, A. J. Yepes, and C. Xiao, "Combining social media and fda adverse event reporting system to detect adverse drug reactions," *Drug safety*, vol. 43, pp. 893–903, 5 2020.
- [93] A. K. Saxena, "Evaluating the regulatory and policy recommendations for promoting information diversity in the digital age," *International Journal of Responsible Artificial Intelligence*, vol. 11, no. 8, pp. 33–42, 2021.
- [94] N. Mehta, K. Born, and B. Fine, "How artificial intelligence can help us 'choose wisely'," *Bioelectronic medicine*, vol. 7, pp. 5–5, 4 2021.
- [95] J. R. Curtis, L. Chen, P. Higginbotham, W. B. Nowell, R. Gal-Levy, J. H. Willig, M. M. Safford, J. Coe, K. O'Hara, and R. Sa'adon, "Social media for arthritis-related comparative effectiveness and safety research and the impact of direct-to-consumer advertising.," *Arthritis research & therapy*, vol. 19, pp. 48–48, 3 2017.
- [96] A. N. Garman, M. P. Standish, and D. H. Kim, "Enhancing efficiency, reliability, and rigor in competency model analysis using natural language processing," *The Journal of Competency-Based Education*, vol. 3, 6 2018.
- [97] S. V. Ramanan, K. Radhakrishna, A. Waghmare, T. Raj, S. P. Nathan, S. M. Sreerama, and S. Sampath, "Dense annotation of free-text critical care discharge summaries from an indian hospital and associated performance of a clinical nlp annotator," *Journal of medical systems*, vol. 40, pp. 1–9, 6 2016.
- [98] N. Hong, A. Wen, F. Shen, S. Sohn, C. Wang, H. Liu, and G. Jiang, "Developing a scalable fhir-based clinical data normalization pipeline for standardizing and integrating unstructured and structured electronic health record data.," *JAMIA* open, vol. 2, pp. 570–579, 10 2019.
- [99] S. N. Murphy, C. Herrick, Y. Wang, T. D. Wang, D. Sack, K. P. Andriole, J. Wei, N. Reynolds, W. J. Plesniak, B. R. Rosen, S. D. Pieper, and R. L. Gollub, "High throughput tools to access images from clinical archives for research," *Journal of digital imaging*, vol. 28, pp. 194–204, 10 2014.
- [100] Y. Zhang, A. Nie, A. M. Zehnder, R. L. Page, and J. Zou, "Vettag: improving automated veterinary diagnosis coding via large-scale language modeling.," *NPJ digital medicine*, vol. 2, pp. 35–35, 5 2019.
- [101] D. Sharma, S. Purushotham, and C. K. Reddy, "Medfusenet: An attention-based multimodal deep learning model for visual question answering in the medical domain.," *Scientific reports*, vol. 11, pp. 19826–, 10 2021.
- [102] N. C. Ernecoff, K. L. Wessell, L. C. Hanson, A. M. Lee, C. M. Shea, S. B. Dusetzina, M. Weinberger, and A. V. Bennett, "Electronic health record phenotypes for identifying patients with late-stage disease: a method for research and clinical application," *Journal of general internal medicine*, vol. 34, pp. 2818–2823, 8 2019.
- [103] C. Xiang and M. Abouelyazid, "The impact of generational cohorts and visit environment on telemedicine satisfaction: A novel investigation," 2020.

- [104] A. Syrowatka, M. Kuznetsova, A. Alsubai, A. L. Beckman, P. A. Bain, K. J. T. Craig, J. Hu, G. P. Jackson, K. Rhee, and D. W. Bates, "Leveraging artificial intelligence for pandemic preparedness and response: a scoping review to identify key use cases," *NPJ digital medicine*, vol. 4, pp. 96–96, 6 2021.
- [105] null Muhammad Bilal Ahmad Jamil and null Duryab Shahzadi, "A systematic review a conversational interface agent for the export business acceleration," *Lahore Garrison University Research Journal of Computer Science and Information Technology*, vol. 7, pp. 37–49, 8 2023.
- [106] J. Zhang, A. Can, P. M. R. Lai, S. Mukundan, V. M. Castro, D. Dligach, S. Finan, V. S. Gainer, N. A. Shadick, G. Savova, S. N. Murphy, T. Cai, S. T. Weiss, and R. Du, "Surrounding vascular geometry associated with basilar tip aneurysm formation.," *Scientific reports*, vol. 10, pp. 17928–17928, 10 2020.