



Original Research

Semantic Embedding Alignment for Cross-Institutional Clinical Text Mining

Santosh Bhandari¹

¹Purbanchal University, Department of Computer Science, Biratnagar-12, Morang, Nepal.

Abstract

This paper explores a unified framework for semantic embedding alignment in cross-institutional clinical text mining contexts. The goal is to ensure that concept representations across different medical institutions retain consistent semantic information despite discrepancies in data collection procedures, annotation guidelines, and linguistic variations in documentation style. The approach focuses on leveraging geometry-aware transformations to map institution-specific embeddings into a common latent space, allowing domain-invariant representations that facilitate downstream clinical tasks. The method addresses challenges arising from heterogeneous data distributions and multi-scale contextual embeddings, providing a foundation for federated clinical studies that respect local privacy constraints. Central to this framework is the concept of establishing topological homomorphisms between embedding spaces, captured through advanced linear algebraic and logical formulations that enable robust alignment even under partial supervision or noisy label constraints. The significance of this methodology lies in its ability to harmonize terminological discrepancies and to reduce the risk of model miscalibration when applying machine learning techniques to sensitive, domain-specific corpora. Through the integration of manifold projections and learned semantic correspondences, the framework promises to facilitate tasks such as clinical named entity recognition, phenotype identification, and automated diagnostic coding, thereby enhancing interoperability and reducing diagnostic gaps across multiple healthcare institutions.

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 general-purpose language tasks. 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 [2]. 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. For instance, tasks may require a model to read a patient report describing a constellation of symptoms and then deduce the most probable diagnosis [3]. 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. 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. [4]

Recent studies have highlighted the limitations of large language models in contexts that require robust interpretability. 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 [5]. 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. 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 [6]. 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. Accordingly, the benchmarks employed must be sufficiently rigorous and reflect real-world data to ensure meaningful results. [7]

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. 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 [8]. 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. 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 [9]. 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. In this vein, probabilistic modeling frameworks can help quantify uncertainty, thereby reflecting a more accurate portrayal of clinical reasoning processes [10]. 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. Clinicians commonly rely on standardized reporting forms, annotated reports, or integrated information from laboratory results [11, 12]. 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. 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. [13]

Ethical and regulatory frameworks also play a pivotal role in shaping the evaluation of large language models for medical use. 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 [14]. 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. Advancements in model architectures, computational resources, and data availability can collectively drive continued improvements in the predictive and inferential power of large language models [15]. 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.

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 [16]. 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.

2. Related Work

Research into applying advanced language models to the medical domain has evolved significantly over the past decade [17]. 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. 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. [18]

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 [19]. 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.

A growing body of work explores the incorporation of external knowledge sources to enhance performance [20]. 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. 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 [21]. Such hybrid approaches have shown promise in improving both model performance and interpretability.

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 [22]. 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 [23]. 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. Mechanisms for continuous monitoring and feedback loops from clinical experts are often proposed as partial remedies [24]. Nevertheless, robust solutions that fully mitigate bias remain an active area of research.

Another relevant direction pertains to interpretability. Several studies highlight the tension between the complexity of large language models and the need for transparent, justifiable recommendations in clinical practice [25]. 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. 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. [26]

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 [27]. 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. These refinements highlight the increasing sophistication in how researchers conceptualize and measure the impact of model outputs. [28]

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. In doing so, they open the possibility for more dynamic forms of inference, resembling the iterative reasoning clinicians conduct as they gather more evidence [29]. 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. 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 [30]. 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 [31]. 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 $f : S \rightarrow D$, mapping each s to a set D of potential diagnoses.

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 [32]. 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. 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 [33]. 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. Specifically, we introduce constraints of the form:

$$\forall x \in R, \text{Symptom}(x) \wedge \text{RiskFactor}(x) \rightarrow \text{HighProbability}(\text{Disease}(x)),$$

which indicate the conditions under which certain diagnoses become highly probable [34]. 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. [35]

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 $M_\theta^1, M_\theta^2, \dots, M_\theta^k$, each corresponding to a distinct pretraining or fine-tuning scheme. [36]

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.

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 [37]. 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 concept embedding trained on a medical ontology). Hence, the combined matrix X^c can be dimensionally reduced using [38]

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

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. [39]

Finally, we address the practical evaluation of uncertainty in diagnostic inference. We propose a Bayesian approximation for the model’s output, where M is replaced by M_+ , with sampled from a distribution reflecting parameter uncertainty. We can then compute a posterior predictive

$$p(d | r) = \int p(d | r, +) p() d. [40]$$

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.

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 [41]. 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.

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 [42]. 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. 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. [43]

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 [44]. 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. In the fine-tuning phase, each model typically converges within 5 to 10 epochs, depending on dataset complexity and size. [45]

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. 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 [46]. 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. 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 [47]. 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 [48]. 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 [49]. 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. 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 [50]. We measure the calibration error by comparing predicted probabilities of correct diagnoses with observed frequencies in the test set. 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 [51]. 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 [52]. Similarly, rare diseases posed difficulties for all models, even those augmented with external medical knowledge. This phenomenon highlights the limitations of data-driven 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. [53]

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. 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. [54]

Collectively, these results underscore both the substantial progress large language models have made and the complexities that remain. 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 [55]. 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. 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 [56]. 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 [57]. These findings align with earlier work suggesting that domain-specific 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.

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 [58]. 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. 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. [59]

Confidence calibration emerged as another important factor, as accurate probabilistic estimates can be vital for high-stakes medical decisions. 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 [60]. 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 [61]. 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. 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 [62]. Moreover, the ethical and regulatory constraints on sharing medical data limit the volume of diverse training sets. 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 [63]. While attention visualization or gradient-based saliency can offer partial insights into model predictions, they do not always align with genuine medical reasoning. 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 [64]. 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 [65]. 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.

Scalability and integration into clinical workflows represent further frontiers [66]. Even with state-of-the-art hardware, large models can be computationally expensive, slowing down inference. 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 [67]. 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. 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 [68]. 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. [69]

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. In the concluding section, we consolidate our findings and propose avenues for subsequent investigations, emphasizing the collaborative nature of progress in this domain. [70]

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. 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 [71]. 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.

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 [72]. Proposed strategies, such as generating synthetic examples or orchestrating large-scale collaborations for data pooling, highlight the scope for further innovation. 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 [73, 74]. 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. The potential integration of formal logic into the model's decision-making process holds promise for more trustworthy explanations [75]. Nevertheless, even logic-based approaches face the possibility of oversimplifying the complexities that underlie clinical judgment. 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 [76]. 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. 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 [77]. However, this iterative process must be carefully balanced with regulatory standards governing software as a medical device, clinical trial validations, and institutional guidelines.

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 [78]. 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. Clinicians commonly encounter situations with incomplete or conflicting data, making it critical for computational systems to signal their own uncertainty [79]. 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.

Additionally, the ethical and social considerations discussed throughout this paper remain vital to future progress [80]. Biases intrinsic to training data can lead to unequal performance across patient demographics, potentially entrenching healthcare disparities. 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. [81]

The present study lays a foundation for future research directions that can deepen and broaden the insights gained here. 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 [82]. 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. 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 [83]. 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.

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 [84]. 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. 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 [85]. The synergy of advanced model architectures, formal logic

constraints, and carefully curated datasets has enabled significant strides in accuracy and consistency. 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. [86]

References

- [1] A. R. Pah, L. J. Rasmussen-Torvik, S. Goel, P. Greenland, and A. N. Kho, “Big data: What is it and what does it mean for cardiovascular research and prevention policy,” *Current Cardiovascular Risk Reports*, vol. 9, pp. 1–9, 11 2014.
- [2] J. M. Harrison, A. Yala, P. Mikhael, J. Roldan, D. Ciprani, T. Michelakos, L. Bolm, M. Qadan, C. Ferrone, C. F.-D. Castillo, K. D. Lillemoe, E. Santus, and K. Hughes, “Successful development of a natural language processing algorithm for pancreatic neoplasms and associated histologic features.,” *Pancreas*, vol. 52, pp. e219–e223, 4 2023.
- [3] A. Arno, J. Elliott, B. C. Wallace, T. Turner, and J. Thomas, “The views of health guideline developers on the use of automation in health evidence synthesis,” *Systematic reviews*, vol. 10, pp. 1–10, 1 2021.
- [4] 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.
- [5] C. Xiang and M. Abouelyazid, “The impact of generational cohorts and visit environment on telemedicine satisfaction: A novel investigation,” 2020.
- [6] R. Avula, “Addressing barriers in data collection, transmission, and security to optimize data availability in healthcare systems for improved clinical decision-making and analytics,” *Applied Research in Artificial Intelligence and Cloud Computing*, vol. 4, no. 1, pp. 78–93, 2021.
- [7] L. Wu, T. Ingle, Z. Liu, A. Zhao-Wong, S. C. Harris, S. Thakkar, G. Zhou, J. Yang, J. Xu, D. Mehta, W. Ge, W. Tong, and H. Fang, “Study of serious adverse drug reactions using fda-approved drug labeling and meddra,” *BMC bioinformatics*, vol. 20, pp. 129–139, 3 2019.
- [8] M. Hu and M. Kejriwal, “Measuring spatio-textual affinities in twitter between two urban metropolises,” *Journal of computational social science*, vol. 5, pp. 1–26, 6 2021.
- [9] Y. Kim, J. H. Garvin, M. K. Goldstein, T. S. Hwang, A. Redd, D. Bolton, P. A. Heidenreich, and S. M. Meystre, “Extraction of left ventricular ejection fraction information from various types of clinical reports,” *Journal of biomedical informatics*, vol. 67, pp. 42–48, 2 2017.
- [10] C. J. Roth, D. A. Clunie, D. Vining, S. A. Berkowitz, A. Berlin, J.-P. Bissonnette, S. D. Clark, T. C. Cornish, M. Eid, C. M. Gaskin, A. K. Goel, G. Jacobs, D. Kwan, D. M. Luviano, M. P. McBee, K. Miller, A. M. Hafiz, C. Obcemea, A. V. Parwani, V. Rotemberg, E. L. Silver, E. S. Storm, J. E. Tcheng, K. S. Thullner, and L. R. Folio, “Multispecialty enterprise imaging workgroup consensus on interactive multimedia reporting current state and road to the future: Himss-siim collaborative white paper,” *Journal of digital imaging*, vol. 34, pp. 495–522, 6 2021.
- [11] M. Tušl, A. Thelen, K. Marcus, A. Peters, E. Shalaeva, B. Scheckel, M. Sykora, S. Elayan, J. A. Naslund, K. Shankardass, S. J. Mooney, M. Fadda, and O. Gruebner, “Opportunities and challenges of using social media big data to assess mental health consequences of the covid-19 crisis and future major events.,” *Discover mental health*, vol. 2, pp. 14–, 6 2022.
- [12] 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.
- [13] S. Ahmed, I. E. Nielsen, A. Tripathi, S. Siddiqui, R. P. Ramachandran, and G. Rasool, “Transformers in time-series analysis: A tutorial,” *Circuits, Systems, and Signal Processing*, vol. 42, pp. 7433–7466, 7 2023.
- [14] L. S. Weiss, X. Zhou, A. M. Walker, A. N. Ananthakrishnan, R. Shen, R. E. Sobel, A. Bate, and R. F. Reynolds, “A case study of the incremental utility for disease identification of natural language processing in electronic medical records,” *Pharmaceutical Medicine*, vol. 32, pp. 31–37, 12 2017.

- [15] Y. Wang, S. Sohn, S. Liu, F. Shen, L. Wang, E. J. Atkinson, S. Amin, and H. Liu, "A clinical text classification paradigm using weak supervision and deep representation.," *BMC medical informatics and decision making*, vol. 19, pp. 1–13, 1 2019.
- [16] Z. Zeng, S. Espino, A. Roy, X. Li, S. A. Khan, S. E. Clare, X. Jiang, R. E. Neapolitan, and Y. Luo, "Using natural language processing and machine learning to identify breast cancer local recurrence," *BMC bioinformatics*, vol. 19, pp. 65–74, 12 2018.
- [17] 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.
- [18] A. Alessa and M. Faezipour, "A review of influenza detection and prediction through social networking sites," *Theoretical biology & medical modelling*, vol. 15, pp. 2–2, 2 2018.
- [19] J. Chen, J. Zheng, and H. Yu, "Finding important terms for patients in their electronic health records: A learning-to-rank approach using expert annotations.," *JMIR medical informatics*, vol. 4, pp. e40–, 11 2016.
- [20] P. Sharedalal, A. Singh, N. Shah, and D. Jain, "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. 1–3, 1 2021.
- [21] M. Alsaidi, M. T. Jan, A. Altaher, H. Zhuang, and X. Zhu, "Tackling the class imbalanced dermoscopic image classification using data augmentation and gan," *Multimedia Tools and Applications*, vol. 83, pp. 49121–49147, 10 2023.
- [22] 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.
- [23] L. E. Young, Y. Nan, E. Jang, and R. Stevens, "Digital epidemiological approaches in hiv research: a scoping methodological review.," *Current HIV/AIDS reports*, vol. 20, pp. 470–480, 11 2023.
- [24] E. Jing and Y.-Y. Ahn, "Characterizing partisan political narrative frameworks about covid-19 on twitter.," *EPJ data science*, vol. 10, pp. 53–53, 10 2021.
- [25] M. Muniswamaiah, T. Agerwala, and C. Tappert, "Big data in cloud computing review and opportunities," *arXiv preprint arXiv:1912.10821*, 2019.
- [26] Q.-Y. Zhong, L. Mittal, M. Nathan, K. M. Brown, D. K. González, T. Cai, S. Finan, B. Gelaye, P. Avillach, J. W. Smoller, E. W. Karlson, T. Cai, and M. A. Williams, "Use of natural language processing in electronic medical records to identify pregnant women with suicidal behavior: towards a solution to the complex classification problem," *European journal of epidemiology*, vol. 34, pp. 153–162, 12 2018.
- [27] A. Abhishek and A. Basu, "A framework for disambiguation in ambiguous iconic environments," in *AI 2004: Advances in Artificial Intelligence: 17th Australian Joint Conference on Artificial Intelligence, Cairns, Australia, December 4-6, 2004. Proceedings 17*, pp. 1135–1140, Springer, 2005.
- [28] S. Bayer, C. Clark, O. Dang, J. S. Aberdeen, S. Brajovic, K. Swank, L. Hirschman, and R. Ball, "Ade eval: An evaluation of text processing systems for adverse event extraction from drug labels for pharmacovigilance," *Drug safety*, vol. 44, pp. 83–94, 10 2020.
- [29] K. Goodman, J. Krueger, and J. Crowley, "The automatic clinical trial: leveraging the electronic medical record in multisite cancer clinical trials.," *Current oncology reports*, vol. 14, pp. 502–508, 8 2012.
- [30] J. LaFleur, R. E. Nelson, Y. Yao, R. A. Adler, and J. R. Nebeker, "Validated risk rule using computerized data to identify males at high risk for fracture," *Osteoporosis international : a journal established as result of cooperation between the European Foundation for Osteoporosis and the National Osteoporosis Foundation of the USA*, vol. 23, pp. 1017–1027, 5 2011.
- [31] L. W. D'Avolio, T. M. Nguyen, S. Goryachev, and L. D. Fiore, "Automated concept-level information extraction to reduce the need for custom software and rules development.," *Journal of the American Medical Informatics Association : JAMIA*, vol. 18, pp. 607–613, 6 2011.
- [32] J. P. Erinjeri, D. Picus, F. W. Prior, D. T. Rubin, and P. Koppel, "Development of a google-based search engine for data mining radiology reports," *Journal of digital imaging*, vol. 22, pp. 348–356, 4 2008.
- [33] H. Edrees, W. Song, A. Syrowatka, A. Simona, M. G. Amato, and D. W. Bates, "Intelligent telehealth in pharmacovigilance: A future perspective.," *Drug safety*, vol. 45, pp. 449–458, 5 2022.

- [34] 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.
- [35] C. S. Rosen, M. M. Matthieu, S. W. Stirman, J. M. Cook, S. J. Landes, N. C. Bernardy, K. M. Chard, J. J. Crowley, A. Eftekhari, E. P. Finley, J. L. Hamblen, J. M. Harik, S. M. Kehle-Forbes, L. A. Meis, E. Osei-Bonsu, A. L. Rodriguez, K. J. Ruggiero, J. I. Ruzek, B. N. Smith, L. R. Trent, and B. V. Watts, "A review of studies on the system-wide implementation of evidence-based psychotherapies for posttraumatic stress disorder in the veterans health administration," *Administration and policy in mental health*, vol. 43, pp. 957–977, 7 2016.
- [36] F. FitzHenry, H. J. Murff, M. E. Matheny, N. Gentry, E. M. Fielstein, S. H. Brown, R. M. Reeves, D. Aronsky, P. L. Elkin, V. P. Messina, and T. Speroff, "Exploring the frontier of electronic health record surveillance: the case of postoperative complications.," *Medical care*, vol. 51, no. 6, pp. 509–516, 2013.
- [37] E. Pierce, N. N. Boytsov, J. Vasey, T. C. Sudaria, X. Liu, K. W. Lavelle, A. Bogdanov, and O. Goldblum, "A qualitative analysis of provider notes of atopic dermatitis-related visits using natural language processing methods," *Dermatology and therapy*, vol. 11, pp. 1305–1318, 5 2021.
- [38] T. Bai, A. K. Chanda, B. L. Egleston, and S. Vucetic, "Ehr phenotyping via jointly embedding medical concepts and words into a unified vector space," *BMC medical informatics and decision making*, vol. 18, pp. 15–25, 12 2018.
- [39] S. Shen, K. Zhang, Y. Zhou, and S. Ci, "Security in edge-assisted internet of things: challenges and solutions," *Science China Information Sciences*, vol. 63, pp. 220302–, 11 2020.
- [40] S. E. Davis, L. Zabolka, R. J. Desai, S. V. Wang, J. C. Maro, K. Coughlin, J. J. Hernández-Muñoz, D. Stojanovic, N. H. Shah, and J. C. Smith, "Use of electronic health record data for drug safety signal identification: A scoping review.," *Drug safety*, vol. 46, pp. 725–742, 6 2023.
- [41] J. Scheibmeir and Y. K. Malaiya, "Social media analytics of the internet of things," *Discover Internet of Things*, vol. 1, pp. 1–15, 7 2021.
- [42] A. Banerji, K. H. Lai, Y. Li, R. R. Saff, C. A. Camargo, K. G. Blumenthal, and L. Zhou, "Natural language processing combined with icd-9-cm codes as a novel method to study the epidemiology of allergic drug reactions.," *The journal of allergy and clinical immunology. In practice*, vol. 8, pp. 1032–1038.e1, 12 2019.
- [43] B. Zaidat, Y. S. Lahoti, A. Yu, K. S. Mohamed, S. K. Cho, and J. S. Kim, "Artificially intelligent billing in spine surgery: An analysis of a large language model.," *Global spine journal*, vol. 15, pp. 1113–1120, 12 2023.
- [44] K. P. Liao, J. Sun, T. Cai, N. Link, C. Hong, J. Huang, J. E. Huffman, J. L. Gronsbell, Y. Zhang, Y. L. Ho, V. M. Castro, V. S. Gainer, S. N. Murphy, C. J. O'Donnell, J. M. Gaziano, K. Cho, P. Szolovits, I. S. Kohane, S. Yu, and T. Cai, "High-throughput multimodal automated phenotyping (map) with application to phewas.," *Journal of the American Medical Informatics Association : JAMIA*, vol. 26, pp. 1255–1262, 8 2019.
- [45] T. T. Ingol, R. Li, R. Ronau, M. A. Klebanoff, R. Oza-Frank, J. Rausch, K. M. Boone, and S. A. Keim, "Underdiagnosis of obesity in pediatric clinical care settings among children born preterm: a retrospective cohort study.," *International journal of obesity (2005)*, vol. 45, pp. 1717–1727, 5 2021.
- [46] S. Golas, T. Shibahara, S. Agboola, H. Otaki, J. Sato, T. Nakae, T. Hisamitsu, G. Kojima, J. Felsted, S. Kakarmath, J. C. Kvedar, and K. Jethwani, "A machine learning model to predict the risk of 30-day readmissions in patients with heart failure: a retrospective analysis of electronic medical records data," *BMC medical informatics and decision making*, vol. 18, pp. 44–44, 6 2018.
- [47] M. Garg, "Mental health analysis in social media posts: A survey.," *Archives of computational methods in engineering : state of the art reviews*, vol. 30, pp. 1819–1842, 1 2023.
- [48] S. Mandal, R. A. Gandhi, and H. Siy, "Modular norm models: practical representation and analysis of contractual rights and obligations," *Requirements Engineering*, vol. 25, pp. 383–412, 8 2019.
- [49] W. T. Kerr and K. N. McFarlane, "Machine learning and artificial intelligence applications to epilepsy: a review for the practicing epileptologist.," *Current neurology and neuroscience reports*, vol. 23, pp. 869–879, 12 2023.
- [50] S. Moon, S. Liu, D. C. Chen, Y. Wang, D. L. Wood, R. Chaudhry, H. Liu, and P. R. Kingsbury, "Salience of medical concepts of inside clinical texts and outside medical records for referred cardiovascular patients," *Journal of healthcare informatics research*, vol. 3, pp. 200–219, 1 2019.

- [51] R. Ramon-Gonen, A. Dori, and S. Shelly, "Towards a practical use of text mining approaches in electrodiagnostic data.," *Scientific reports*, vol. 13, pp. 19483–, 11 2023.
- [52] S. Chan, V. Reddy, B. Myers, Q. Thibodeaux, N. Brownstone, and W. Liao, "Machine learning in dermatology: Current applications, opportunities, and limitations," *Dermatology and therapy*, vol. 10, pp. 365–386, 4 2020.
- [53] Y. Liu, Y. Luo, and A. M. Naidech, "Big data in stroke: How to use big data to make the next management decision.," *Neurotherapeutics : the journal of the American Society for Experimental NeuroTherapeutics*, vol. 20, pp. 744–757, 3 2023.
- [54] C. R. Weir, N. Staggers, B. Gibson, K. Doing-Harris, R. Barrus, and R. Dunlea, "A qualitative evaluation of the crucial attributes of contextual information necessary in ehr design to support patient-centered medical home care," *BMC medical informatics and decision making*, vol. 15, pp. 30–30, 4 2015.
- [55] M. Chary, N. Genes, A. M. McKenzie, and A. F. Manini, "Leveraging social networks for toxicovigilance.," *Journal of medical toxicology : official journal of the American College of Medical Toxicology*, vol. 9, pp. 184–191, 4 2013.
- [56] A. B. Chapman, A. L. Jones, A. T. Kelley, B. E. Jones, L. Gawron, A. E. Montgomery, T. Byrne, Y. Suo, J. Cook, W. B. P. Pettey, K. S. Peterson, M. Jones, and R. E. Nelson, "Rehoused: A novel measurement of veteran housing stability using natural language processing.," *Journal of biomedical informatics*, vol. 122, pp. 103903–103903, 8 2021.
- [57] E. Klang, B. R. Kummer, N. S. Dangayach, A. Zhong, M. A. Kia, P. Timsina, I. Cossentino, A. Costa, M. A. Levin, and E. K. Oermann, "Predicting adult neuroscience intensive care unit admission from emergency department triage using a retrospective, tabular-free text machine learning approach," *Scientific reports*, vol. 11, pp. 1381–, 1 2021.
- [58] A. Sharma and K. Forbus, "Automatic extraction of efficient axiom sets from large knowledge bases," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 27, pp. 1248–1254, 2013.
- [59] A. Sharma and K. D. Forbus, "Graph-based reasoning and reinforcement learning for improving q/a performance in large knowledge-based systems," in *2010 AAAI Fall Symposium Series*, 2010.
- [60] N. Sadat, M. A. Aziz, N. Mohammed, S. V. S. Pakhomov, H. Liu, and X. Jiang, "A privacy-preserving distributed filtering framework for nlp artifacts," *BMC medical informatics and decision making*, vol. 19, pp. 183–, 9 2019.
- [61] R. Avula *et al.*, "Data-driven decision-making in healthcare through advanced data mining techniques: A survey on applications and limitations," *International Journal of Applied Machine Learning and Computational Intelligence*, vol. 12, no. 4, pp. 64–85, 2022.
- [62] J. A. Walsh, S. Pei, G. K. Penmettsa, J. L. Hansen, G. W. Cannon, D. O. Clegg, and B. C. Sauer, "Identification of axial spondyloarthritis patients in a large dataset: The development and validation of novel methods," *The Journal of rheumatology*, vol. 47, pp. 42–49, 3 2019.
- [63] S. Nath, A. Marie, S. Ellershaw, E. Korot, and P. A. Keane, "New meaning for nlp: the trials and tribulations of natural language processing with gpt-3 in ophthalmology.," *The British journal of ophthalmology*, vol. 106, pp. 889–892, 5 2022.
- [64] B. Tang, H. Cao, Y. Wu, M. Jiang, and H. Xu, "Recognizing clinical entities in hospital discharge summaries using structural support vector machines with word representation features," *BMC medical informatics and decision making*, vol. 13, pp. 1–10, 4 2013.
- [65] A. Sharma and K. Forbus, "Graph traversal methods for reasoning in large knowledge-based systems," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 27, pp. 1255–1261, 2013.
- [66] F. Li, Y. Jin, W. Liu, B. P. S. Rawat, P. Cai, and H. Yu, "Fine-tuning bidirectional encoder representations from transformers (bert)-based models on large-scale electronic health record notes: An empirical study," *JMIR medical informatics*, vol. 7, pp. e14830–, 9 2019.
- [67] A. Can, P. M. R. Lai, V. M. Castro, S. Yu, D. Dligach, S. Finan, V. S. Gainer, N. A. Shadick, G. Savova, S. N. Murphy, T. Cai, S. T. Weiss, and R. Du, "Decreased total iron binding capacity may correlate with ruptured intracranial aneurysms.," *Scientific reports*, vol. 9, pp. 6054–6054, 4 2019.
- [68] R. Dias and A. Torkamani, "Artificial intelligence in clinical and genomic diagnostics.," *Genome medicine*, vol. 11, pp. 70–70, 11 2019.
- [69] Y. Si, J. Wang, H. Xu, and K. Roberts, "Enhancing clinical concept extraction with contextual embeddings.," *Journal of the American Medical Informatics Association : JAMIA*, vol. 26, pp. 1297–1304, 7 2019.

- [70] L. Rasmy, Y. Xiang, Z. Xie, C. Tao, and D. Zhi, "Med-bert: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction.," *NPJ digital medicine*, vol. 4, pp. 86–86, 5 2021.
- [71] X. Cai, S. Liu, J. Han, L. Yang, Z. Liu, and T. Liu, "Chestxraybert: A pretrained language model for chest radiology report summarization," *IEEE Transactions on Multimedia*, vol. 25, pp. 845–855, 2023.
- [72] H. Sharma, C. Mao, Y. Zhang, H. Vatani, L. Yao, Y. Zhong, L. V. Rasmussen, G. Jiang, J. Pathak, and Y. Luo, "Developing a portable natural language processing based phenotyping system," *BMC medical informatics and decision making*, vol. 19, pp. 78–78, 4 2019.
- [73] B. Athira, J. Jones, S. M. Idicula, A. Kulanthaivel, and E. Zhang, "Annotating and detecting topics in social media forum and modelling the annotation to derive directions-a case study," *Journal of Big Data*, vol. 8, pp. 1–23, 2 2021.
- [74] 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.
- [75] Z. Xiao, J. Zhu, Y. Wang, P. Zhou, W. H. Lam, M. A. Porter, and Y. Sun, "Detecting political biases of named entities and hashtags on twitter," *EPJ Data Science*, vol. 12, 6 2023.
- [76] 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.
- [77] C. Reich and B. Meder, "The heart and artificial intelligence-how can we improve medicine without causing harm.," *Current heart failure reports*, vol. 20, pp. 271–279, 6 2023.
- [78] C. M. Lineback, R. Garg, E. Oh, A. M. Naidech, J. L. Holl, and S. Prabhakaran, "Prediction of 30-day readmission after stroke using machine learning and natural language processing.," *Frontiers in neurology*, vol. 12, pp. 649521–, 7 2021.
- [79] K. E. Corey, U. Kartoun, H. Zheng, and S. Y. Shaw, "Development and validation of an algorithm to identify nonalcoholic fatty liver disease in the electronic medical record.," *Digestive diseases and sciences*, vol. 61, pp. 913–919, 11 2015.
- [80] T. B. Smith, R. Vacca, L. Mantegazza, and I. Capua, "Discovering new pathways toward integration between health and sustainable development goals with natural language processing and network science.," *Globalization and health*, vol. 19, pp. 44–, 6 2023.
- [81] G. E. Powell, H. A. Seifert, T. Reblin, P. J. Burstein, J. Blowers, J. A. Menius, J. L. Painter, M. Thomas, C. E. Pierce, H. Rodriguez, J. S. Brownstein, C. C. Freifeld, H. Bell, and N. Dasgupta, "Social media listening for routine post-marketing safety surveillance," *Drug safety*, vol. 39, pp. 443–454, 1 2016.
- [82] 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.
- [83] J. Xu, P. Yang, S. Xue, B. Sharma, M. Sanchez-Martin, F. Wang, K. A. Beaty, D. Elinor, and B. Parikh, "Translating cancer genomics into precision medicine with artificial intelligence: applications, challenges and future perspectives.," *Human genetics*, vol. 138, pp. 109–124, 1 2019.
- [84] C.-Y. Chen, P. H. Lee, V. M. Castro, J. Minnier, A. W. Charney, E. A. Stahl, D. M. Ruderfer, S. N. Murphy, V. S. Gainer, T. Cai, I. Jones, C. N. Pato, M. T. Pato, M. Landén, P. Sklar, R. H. Perlis, and J. W. Smoller, "Genetic validation of bipolar disorder identified by automated phenotyping using electronic health records.," *Translational psychiatry*, vol. 8, pp. 86–86, 4 2018.
- [85] L. T.-E. Cheng, J. Zheng, G. Savova, and B. J. Erickson, "Discerning tumor status from unstructured mri reports—completeness of information in existing reports and utility of automated natural language processing.," *Journal of digital imaging*, vol. 23, pp. 119–132, 5 2009.
- [86] M. E. Matheny, F. FitzHenry, T. Speroff, J. Green, M. L. Griffith, E. E. Vasilevskis, E. M. Fielstein, P. L. Elkin, and S. H. Brown, "Detection of infectious symptoms from va emergency department and primary care clinical documentation," *International journal of medical informatics*, vol. 81, pp. 143–156, 1 2012.