

Journal Pre-proof



A comprehensive review of Intelligent Question-Answering Systems in Traditional Chinese Medicine Based on LLMs

Qilan Xu, Tong Wu, Yiwen Wang, Xingyu Li, Heshui Yu, Shixin Cen, Zheng Li

PII: S2095-1779(25)00223-0

DOI: <https://doi.org/10.1016/j.jpha.2025.101406>

Reference: JPHA 101406

To appear in: *Journal of Pharmaceutical Analysis*

Received Date: 25 March 2025

Revised Date: 2 July 2025

Accepted Date: 19 July 2025

Please cite this article as: Q. Xu, T. Wu, Y. Wang, X. Li, H. Yu, S. Cen, Z. Li, A comprehensive review of Intelligent Question-Answering Systems in Traditional Chinese Medicine Based on LLMs, *Journal of Pharmaceutical Analysis*, <https://doi.org/10.1016/j.jpha.2025.101406>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 The Author(s). Published by Elsevier B.V. on behalf of Xi'an Jiaotong University.

A Comprehensive Review of Intelligent Question-Answering Systems in Traditional Chinese Medicine Based on LLMs

Qilan Xu^{a, 1}, Tong Wu^{a, 1}, Yiwen Wang^a, Xingyu Li^a, Heshui Yu^{a, b, c, d}, Shixin Cen<sup>a,
b, c, d, **</sup>, Zheng Li^{a, b, c, d, e_—*}

^a College of Pharmaceutical Engineering of Traditional Chinese Medicine, Tianjin University of Traditional Chinese Medicine, Tianjin, 301617, China

^b State Key Laboratory of Component-Based Chinese Medicine, Tianjin, 301617, China

^c State Key Laboratory of Chinese Medicine Modernization, Tianjin, 301617, China

^d Tianjin Key Laboratory of Intelligent and Green Pharmaceuticals for Traditional Chinese Medicine

^e Haihe Laboratory of Modern Chinese Medicine, Tianjin, 301617, China

¹ Both authors contributed equally to this work.

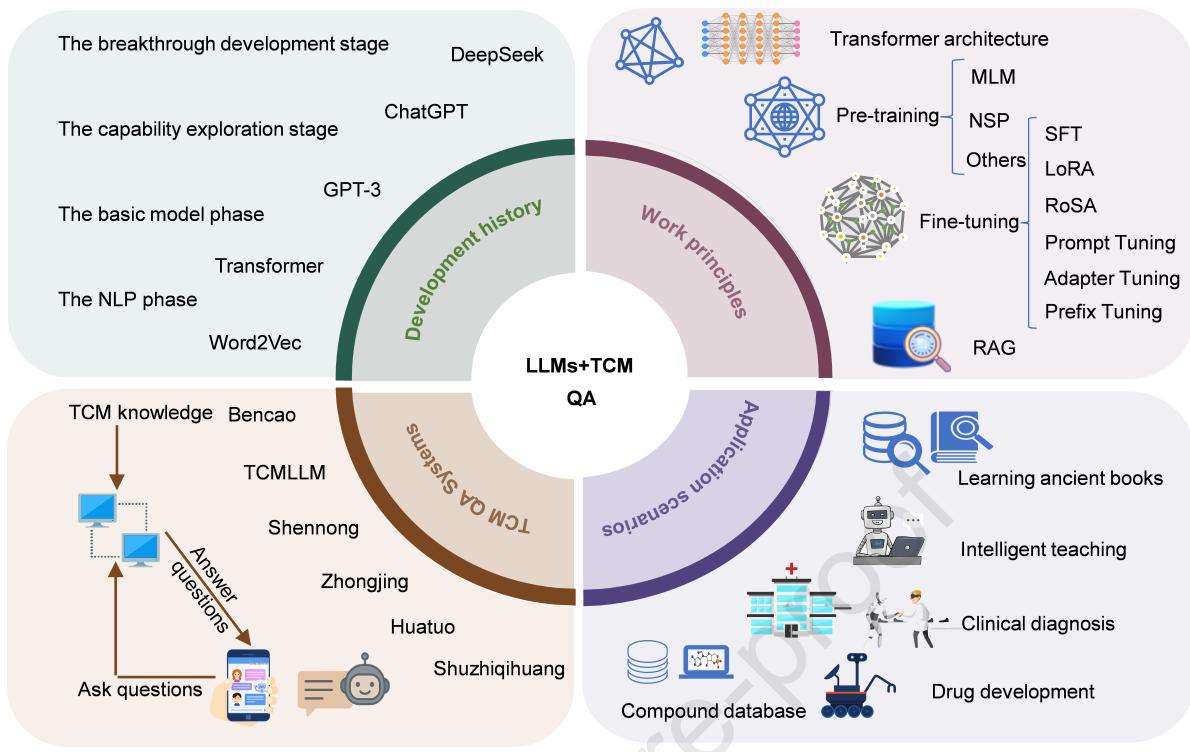
* Corresponding author.

College of Pharmaceutical Engineering of Traditional Chinese Medicine, Tianjin University of Traditional Chinese Medicine, Tianjin, 301617, China.

** Corresponding author.

College of Pharmaceutical Engineering of Traditional Chinese Medicine, Tianjin University of Traditional Chinese Medicine, Tianjin, 301617, China.

E-mail addresses: lizheng@tjutcm.edu.cn (Z. Li), censhixin@tjutcm.edu.cn (S. Cen).



Review paper**A comprehensive review of Intelligent Question-Answering Systems in
Traditional Chinese Medicine Based on LLMs****Abstract:**

Large language models (LLMs) are advanced deep learning models with billions or even trillions of parameters, enabling powerful natural language processing and knowledge reasoning capabilities. Their applications in the medical domain have been rapidly expanding, spanning medical research, clinical diagnosis, drug development, and patient management. As a cornerstone of China's healthcare system, traditional Chinese medicine (TCM) faces significant challenges, including difficulties in knowledge extraction, and lack of standardization. The emergence of TCM-focused LLMs presents a transformative opportunity, offering a novel technological framework to process vast amounts of TCM data, uncover hidden theoretical insights, and enhance both research and clinical applications. Despite the growing interest in AI-driven medical solutions, systematic research on LLMs in the TCM domain remains limited. This article provides a comprehensive review of LLM development, detailing their underlying mechanisms, training methodologies, and key technological advancements. It further explores the unique characteristics and diverse application scenarios of existing TCM-LLMs. Additionally, this study also conducts a horizontal comparison of the differences between intelligent question-answering (QA) systems on general LLMs and QA systems on TCM-LLMs, discusses challenges and potential risks, and offers strategic recommendations for future development. By synthesizing current advancements and addressing critical gaps, this work aims to support the continued modernization and intelligent evolution of TCM, fostering its integration into contemporary healthcare systems.

Keywords: Large language models, Traditional Chinese medicine-large language models, Traditional Chinese medicine question-answering systems

28 **1. Introduction**

29 Large language models (LLMs) are deep neural networks with a vast number of
30 parameters, designed to develop advanced language comprehension and generation
31 capabilities through large-scale pre-training (PT) tasks. These models learn intricate
32 linguistic structures and patterns, enhancing their ability to process and generate
33 human-like text [1]. Through the processing of massive text data, LLMs perform
34 excellently in various natural language processing (NLP) tasks. For example, in
35 machine translation, text generation, question answering (QA) systems, and sentiment
36 analysis, the capabilities of LLMs significantly surpass those of traditional models,
37 demonstrating strong generalization performance [2,3]. Their deep understanding of
38 complex language features enables them to efficiently handle semantic reasoning,
39 knowledge integration, and task adaptation. Consequently, with their powerful
40 language comprehension and information processing capabilities, LLMs have rapidly
41 expanded their applications in the medical field, including new drug design [4],
42 development of personalized medical treatment plans [5], clinical diagnosis of diseases
43 [6], and medical education [7]. LLMs have demonstrated significant value and become
44 an important technical tool for promoting the intelligent development of healthcare.

45 Nowadays, the field of traditional Chinese medicine (TCM) faces numerous
46 challenges, such as insufficient data standardization, and difficulty in data integration.
47 TCM data sources are diverse, covering multiple aspects such as ancient books,
48 literatures, and clinical records. The data types include different fields such as biology,
49 chemistry, and pharmacology. In terms of data format, it presents multi-modal
50 characteristics, including various forms such as text, charts, and images. Therefore, we
51 have systematically sorted out and summarized the data sources in TCM (Fig. 1). The
52 integration of TCM information is confronted with complex and highly heterogeneous
53 issues [8,9]. Traditional manual analysis methods are time-consuming, labor-intensive,
54 and difficult to systematically and efficiently uncover the deep-seated patterns behind
55 the data [10].

56 < **Fig. 1.** Summary of the data in the field of TCM. >

57 Based on previous reports, some studies have investigated the application of LLMs
58 to overcome the challenges of data integration in TCM analysis [11]. The powerful NLP
59 capabilities of LLMs enable them to quickly process and interpret complex TCM data,
60 especially demonstrating unique advantages in handling TCM's multimodal and
61 heterogeneous data. These models can automatically organize classical Chinese
62 medical literature and extract key information. Additionally, they uncover deep
63 correlations between different data sources through comprehensive analysis of clinical
64 cases and TCM knowledge bases. For instance, by integrating classical TCM literature,
65 modern clinical data, and Chinese medicinal materials knowledge, LLMs can identify
66 diagnostic and treatment patterns, optimize treatment plans, and provide personalized
67 treatment recommendations [12]. LLMs also offer critical technical support for the
68 standardization, digitalization, and intelligent development of TCM [13]. Through their
69 automated data processing and pattern recognition capabilities, LLMs can significantly
70 enhance TCM research and clinical diagnosis. Meanwhile, the model's high-
71 performance computing and data mining potential offer strong support for the transition
72 of TCM from traditional experience-based medicine to modern medicine.

73 By using the Web of Sciences database, we have retrieved literature related to
74 LLM research on a global scale. The paper search involves keywords such as LLMs,
75 and intelligent QA systems (Fig. 2). After removing irrelevant and duplicate articles, a
76 total of 212 papers have been published since 2020, with 187 published in 2023.
77 Subsequently, we conducted keyword analysis and found that in recent years, the
78 combination of TCM and LLMs has become one of the research hotspots. As a result,
79 a variety of representative TCM-LLMs have emerged, including the Shuzhiqihuang 2.0,
80 Huatuo, Huangdi, and so on. These models have demonstrated extensive application
81 prospects and significant achievements in the practice of TCM.

82 < **Fig. 2.** Overview analysis diagram of intelligent question-answering (QA) systems
83 based on large language models (LLMs). >

84 As an important research direction of TCM-LLMs, compared with traditional
85 search engines, intelligent QA systems can quickly and efficiently obtain knowledge
86 and information conveniently. With their powerful automated information retrieval and
87 knowledge reasoning capabilities, they can provide accurate consulting services for
88 TCM and further support clinical decision-making and patient management [14].
89 Previous reviews of LLMs have focused on their model architectures, key technologies,
90 and training methods [15]. Particularly, in general fields such as NLP and dialogue
91 systems, a relatively systematic understanding has been established [16]. However,
92 with the continuous acceleration of the digital and intelligent transformation of TCM
93 disciplines, TCM-LLMs have gradually become a new research hotspot. Although
94 existing literature have reported some application explorations of TCM-LLMs, such as
95 prescription recommendation, case analysis and QA on TCM knowledge [17], there is
96 still a lack of systematic review and evaluation of their overall development status and
97 potential challenges.

98 Based on the above background, this work aims to provide a systematic review of
99 relevant research on TCM-LLMs. First, it reviews the evolutionary progression of
100 LLMs as well as the key technologies and core methods, outlines the model
101 architectures and training processes commonly adopted in LLMs. On this basis, this
102 paper further analyzes the unique characteristics and technical advantages of the TCM
103 large language model, deeply explores its diverse application scenarios in practice, and
104 examines the profound impact it may have on the field of TCM. At the same time, this
105 articles also analyzes the core challenges faced by TCM-LLMs in aspects, such as data
106 quality, health management, and ethical norms. These challenges not only determine
107 the technical barriers in this field but also directly affect the sustainable development
108 and wide application of TCM-LLMs. Through comprehensive analysis of relevant
109 research results and limitations, this work summarizes the research status and
110 development trends of TCM-LLMs, puts forward several thoughts and suggestions, and

111 provides a scientific basis and reference significance for subsequent research and
 112 industrial practice.

113 **2. The evolutionary progression of LLMs**

114 The development of LLMs can be roughly divided into four stages: the NLP phase,
 115 the basic model phase, the capability exploration stage, and the breakthrough
 116 development stage [18,19]. To more systematically show the development process of
 117 LLMs, we further sorted and drew its development process diagram (Fig. 3).

118 < Fig. 3. Background of large language models (LLMs) development. >

119 2.1. The NLP phase (2013-2014)

120 During this period, word embedding (WE) techniques such as Word2Vec and
 121 GloVe rapidly advanced and became a significant milestone in NLP. Unlike traditional
 122 methods, such as N-grams or one-hot vectors, WE can map words into a low-
 123 dimensional, dense vector space, addressing data sparsity and capturing rich semantic
 124 and syntactic features. Word2Vec employs the Continuous-Bag-of-Words and Skip-
 125 gram models with shallow neural networks to predict local context, while GloVe
 126 generates high-quality word vectors by optimizing the global word frequency co-
 127 occurrence matrix. The development of WE technology has spurred innovations in
 128 language models and provided efficient, versatile representations for tasks such as
 129 sentiment analysis and machine translation [20–22].

130 2.2. The basic model phase (2017-2018)

131 In 2017, Vaswani et al. [23] introduced the Transformer architecture,
 132 revolutionizing machine translation by overcoming the limitations of traditional
 133 Recurrent Neural Networks (RNNs) and Long Short-Term Memory Networks. The
 134 Transformer's innovative attention mechanism enabled efficient parallel computing,
 135 improving both translation quality and model performance. This breakthrough laid the
 136 foundation for advancements in NLP tasks such as language generation and question
 137 answering, establishing the Transformer as the dominant architecture in deep learning
 138 (DL). In 2018, OpenAI released GPT-1, the first pre-trained transformer-based model

139 with 117 million (M) parameters, followed by Google's BERT model with 340 M
 140 parameters [24,25]. The success of these models marked the rise of transformer
 141 architectures in NLP, advancing the application of NLP technologies and ushering in a
 142 new era in DL dominated by pre-trained models.

143 2.3. The capability exploration stage (2019-2022)

144 During this stage, the deployment of large-scale language models faced challenges
 145 such as high computational demands, long training times, and the risk of overfitting due
 146 to task-specific fine-tuning (FT). Researchers began exploring solutions to eliminate
 147 the need for task-specific FT, enhancing model capabilities [26]. Compared to GPT-1,
 148 GPT-2 saw significant improvements, with its parameter scale increasing from 117 M
 149 to 1.5 billion (B), and training data expanding from 5 GB. These advancements greatly
 150 improved GPT-2's generalization, boosting performance across tasks [24,27]. OpenAI
 151 then released GPT-3, which further increased its parameter scale to 175 B, positioning
 152 it as one of the largest models of its time [24,28]. These advancements in language
 153 models provided technical support for new NLP research, shifting tasks from
 154 specialized to more generalized. At the same time, large models exhibit extraordinary
 155 generality and scalability.

156 2.4. The breakthrough development stage (2022-Present)

157 The release of ChatGPT marked a significant milestone in the development of
 158 large models, highlighting progress in NLP. With a parameter scale in the hundreds of
 159 billions and training data reaching hundreds of terabytes, ChatGPT advanced natural
 160 language generation technology and enhanced the intelligence and interactivity of
 161 dialogue systems [29]. Building on ChatGPT, OpenAI launched the GPT-3.5 and GPT-
 162 4 series, with GPT-4 offering notable improvements in model capabilities and
 163 application range, particularly in language generation accuracy and context
 164 understanding depth. Additionally, GPT-4's multimodal capabilities enable it to
 165 process both text and images, significantly improving user experience [24]. These
 166 models provide new directions for multimodal intelligence and complex data analysis,

167 and offer new tools for the development of TCM.

168 **3. The work principles and key technologies of LLMs**

169 **3.1. LLMs working principle:**

170 LLMs are deep neural networks with hundreds of billions of parameters, built on
 171 various architectures such as Transformer, MAMBA [30], Falcon Mamba 7B [31], and
 172 Receptance Weighted Key Value [32], with transformer being the most fundamental
 173 and widely adopted. The Transformer architecture (Fig. 4A), distinct from traditional
 174 RNNs and CNNs, enables high parallelism, reduces training time. It also excels at
 175 processing large-scale data, making it particularly suitable for tasks like machine
 176 translation and text generation. It consists of two core components, the Encoder, which
 177 reads and understands input text, and the Decoder, which generates output. This
 178 encoder-decoder collaboration underpins its outstanding performance in language
 179 modeling and generation tasks [23,24]. The core architecture of the transformer model
 180 consists of five fundamental components, namely word embeddings, the attention
 181 mechanism, the multi-head attention mechanism, feedforward neural networks and
 182 positional encodings.

183 <**Fig. 4.** Working principle of large language models (LLMs) and training methods. >

184 **3.1.1 Word embedding**

185 Input representation refers to transforming original data (such as text, images or
 186 audio) into a numerical form that can be processed by the model, usually represented
 187 in the form of vectors. The core goal of this process is to map complex and diverse
 188 input data to a unified mathematical space to facilitate the model's understanding and
 189 learning [23]. The formula is represented as:

$$190 \quad E = Embedding(X)$$

191 Where X denotes the input sequence, and E refers to the embedded representation
 192 of the input sequence.

193 **3.1.2 Attention mechanism**

194 The attention mechanism, introduced in 2014 [33], originated in machine
 195 translation as a method inspired by human selective focus. It enables models to
 196 concentrate on the most relevant parts of input data, improving translation accuracy. By
 197 allowing the decoder to dynamically attend to different parts of the source sequence, it
 198 alleviates the need for the encoder to compress all information into a fixed-length vector,
 199 enabling more effective information flow and selective retrieval during decoding.

200 The attention mechanism calculates relevance scores using three components:
 201 query, key, and value. Analogous to an information retrieval process, the query
 202 represents the search intent, the key serves as the index, and the value holds the content.
 203 By comparing the query with the key, the model assigns weights to each input element,
 204 reflecting their relevance. These weights are then used to generate context-aware
 205 outputs. The calculation is as follows [34]:

$$206 \quad \text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

207 where Q represents query, K represents key, V denotes value, T is sequence length,
 208 and d denotes feature dimension.

209 3.1.3 Multi-head attention

210 The multi-head attention mechanism extends traditional attention by computing
 211 multiple attention functions in parallel, each as an independent “head.” This design
 212 significantly improves computational efficiency, particularly in scenarios involving
 213 large-scale data and deep neural networks, by facilitating parallelized operations.
 214 Moreover, it allows the model to capture information from different representational
 215 subspaces across positions, thereby avoiding the over-smoothing effect of single-head
 216 attention [35,36]. In practice, we compute attention over a set of queries in matrix Q ,
 217 along with the corresponding keys and values, to derive the output matrix as follows:

$$218 \quad \text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)W^0$$

$$219 \quad \text{where } \text{head} = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

220 where H denotes the number of heads, which refers to the number of parallel
 221 attention mechanisms, W^0 represents a trainable matrix used for the linear
 222 transformation of the output, W_i^Q , W_i^K , and W_i^V represents the linear transformation
 223 matrix for the i -th head, QW_i^Q , KW_i^K , and VW_i^V are the results after mapping the query,
 224 key, and value to different subspaces.

225 3.1.4 Feed-forward network (FFN)

226 Each encoder and decoder layer, aside from the attention sub-layer, contains a fully
 227 connected FFN applied independently and identically at each position. This FFN
 228 comprises two linear transformations separated by a ReLU activation [37]. Despite its
 229 simplicity, the feedforward layer is essential for the transformer's strong performance
 230 [38]. Dong et al. [39] highlighted that stacking self-attention modules alone can cause
 231 rank collapse, leading to token uniformity bias. The feedforward layer plays a key role
 232 in mitigating this issue.

233 3.1.5 Position encoding

234 The primary function of a positional encoder is to add a vector representing the
 235 position of each input element in a sequence. Since the transformer model lacks the
 236 inherent ability to process sequence order, position encoding is added to the input
 237 embedding in both the encoder and decoder stacks. This approach addresses the
 238 transformer's lack of positional information and enhances its efficiency and flexibility
 239 when processing sequence data [23]. By incorporating positional information, the
 240 model can better capture sequential relationships. Three common methods for position
 241 encoding are as follows:

242 Sinusoidal Positional Encoding, proposed by Bahdanau et al. [23], is a fixed
 243 method that doesn't require additional learning. It uses sine and cosine functions to
 244 generate positional encodings, enabling the model to capture relative positional
 245 relationships. The periodicity of these functions allows the model to handle sequences
 246 of varying lengths [40]. The calculation method of position encoding is as follows:

247 $PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{model}})$

248 $PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}})$

249 where pos represents the index of the current element (for instance the position of
 250 the word in the sentence), d represents the dimension of positional encoding, and i refers
 251 to the dimension index of the positional encoding.

252 Learned Positional Encoding treats position encodings as trainable parameters,
 253 allowing the model to learn the best representation for each position. Using an
 254 embedding layer, the model optimizes encodings through training data, updating them
 255 during training, similar to WE [23,41]. Shaw et al. [42] proposed adding learnable
 256 relative position embeddings to the attention mechanism, enabling adaptive
 257 adjustments of position encodings to enhance task-specific performance.

258 Relative Positional Encoding offers a more flexible alternative to fixed and learned
 259 position encodings by focusing on the relative distances between elements rather than
 260 their absolute positions [42–44]. This approach has demonstrated superior performance
 261 in specific tasks, particularly those involving long sequences. Notably, Alibi [45] and
 262 Rotary Position Embedding (RoPE) [46] are two widely adopted relative position
 263 encoding methods in LLMs. The technique has been integrated into several prominent
 264 architectures, including Transformer-XL [47], DeBERTa [48], T5 [49], and Roformer
 265 [46]. As large-scale models continue to evolve, relative positional encoding is expected
 266 to play an increasingly crucial role in managing complex tasks and extensive datasets.

267 3.2. Training methods

268 The training of a complete LLM typically involves three stages, i.e., PT from
 269 scratch, FT of an existing LLM, and alignment with specific application scenarios via
 270 prompt-based methods [50]. In this section, we elaborate on these three stages, PT, FT,
 271 and RAG (Fig. 4B).

272 3.2.1 PT

273 PT is the first stage in LLMs training, establishing a foundation for model
 274 capabilities [51]. During PT, LLMs are trained on large-scale text corpora using

275 unsupervised or self-supervised learning (SSL) methods to align modalities and acquire
 276 multimodal world knowledge [52]. This process enables LLMs to develop rich
 277 language understanding and generation skills, including vocabulary relationships,
 278 grammar, and context dependencies [53–55].

279 Common PT strategies include Masked Language Model (MLM), Autoregressive
 280 Language Model, and Next Sentence Prediction (NSP). These methods enable the
 281 model to learn language patterns from large text corpora and improve performance in
 282 various NLP tasks. The MLM captures context by masking part of the input and
 283 predicting missing words [56]. The Autoregressive Language Model learns language
 284 sequences from left to right, generating text gradually [57]. The NSP analyzes the
 285 logical relationship between sentences, capturing long-term dependencies and
 286 generating text autoregressively [58].

287 In addition to these common methods, large-scale PT models also employ
 288 advanced techniques such as contrastive learning and multi-task learning, enhancing
 289 the model’s expressiveness and transferability [59]. Furthermore, emerging methods
 290 like Generative Adversarial Networks and meta-learning provide new avenues for
 291 improving model learning [60,61]. To explore the practical effects of these PT strategies,
 292 we provide a detailed overview of current large language model PT methods, including
 293 their underlying structures, model names, parameters, and PT data scales (Table 1)
 294 [27,48,49,52, 62-83].

295 < **Table 1.** Comparison of the model structure and Pre-Training (PT) parameters of
 296 general large language models (LLMs). >

297 3.2.2 FT

298 After PT, LLMs have obtained general capabilities to solve various tasks.
 299 However, research shows that the performance of large models can be further optimized
 300 according to specific goals. The goal of FT is to utilize a pre-trained model that has
 301 already learned rich language knowledge and make it better adapt to specific tasks or
 302 domains by further training on specific datasets [84]. In this process, the PT language

303 model is initialized with the learned parameters first, and then trained on a specific
 304 dataset. In this way, the parameters of the model will be updated according to the data
 305 of specific tasks, enabling it to better adapt to the target task [19]. However, the cost of
 306 deploying LLM in practical work is very high. Therefore, how to reduce operating costs
 307 while maintaining performance has become a new research field. In this section, we
 308 summarize several common methods to improve the efficiency of LLM.

309 Supervised Fine-Tuning (SFT), also known as Instruction Fine-Tuning (IFT) [85],
 310 refers to the process of enhancing a pre-trained model's performance through labeled
 311 data. Initially, the model learns general knowledge via large-scale unsupervised
 312 learning. It is then fine-tuned on annotated data from specific domains or tasks, enabling
 313 it to generate more accurate and contextually appropriate outputs for novel inputs [86].
 314 This approach has been widely adopted in leading LLMs such as ChatGPT, FLAN [75],
 315 and OPT-IML [87].

316 Prompt Tuning [56]: Unlike traditional supervised learning, it utilizes a LLM that
 317 has been trained on a large-scale text corpus. By defining a new prompt function, the
 318 model can achieve better performance on specific tasks.

319 Prefix Tuning [88]: A lightweight FT method for natural language generation tasks.
 320 In prefix tuning, a set of prefix vectors trained for specific tasks are appended to the
 321 frozen transformer layers. The prefix vectors are virtual tokens and are attended to by
 322 the context tokens on the right. In addition, this method can keep the language model
 323 parameters unchanged, thus achieving the purpose of adapting to different tasks.

324 Adapter Tuning: Adapters are small neural modules inserted between or within
 325 transformer layers, enabling efficient task-specific FT without altering the original
 326 model parameters [89]. Comprising a dimensionality reduction layer, a nonlinear layer,
 327 and a dimensionality expansion layer, adapters introduce minimal trainable parameters
 328 while effectively adapting LLMs to downstream tasks. For instance, T5 model employs
 329 adapters for FT after PT [90].

330 Low-Rank Adaptation (LoRA) [91]: LoRA model by inserting trainable low-rank

331 matrices into key layers, avoiding changes to the original model architecture. Most pre-
 332 trained parameters are frozen, with only the low-rank components updated. This
 333 approach enables efficient adaptation with minimal parameter updates while preserving
 334 the original model's knowledge.

335 Robust Adaptation (RoSA): Nikdan et al. [92] proposed a parameter-efficient FT
 336 method inspired by robust principal component analysis. RoSA jointly trains low-rank
 337 and high-sparsity components to enhance model performance. Experimental results
 338 show that RoSA outperforms traditional low-rank and sparse FT approaches under
 339 limited computational and memory resources.

340 3.2.3 Distributed training

341 Due to the extremely large scale of LLMs, training a high-performance LLM poses
 342 tremendous challenges. To effectively learn the network parameters of LLMs,
 343 distributed training algorithms are usually required, combined with multiple parallel
 344 strategies to improve training efficiency. At present, several optimization frameworks
 345 for distributed training have been released to promote the implementation and
 346 deployment of parallel algorithms, such as DeepSpeed and Megatron-LM.

347 DeepSpeed [93]: A library for scalable distributed training and inference of DL
 348 models not only optimizes memory management but also significantly improves
 349 training efficiency.

350 Megatron-LM [94–96]: An NVIDIA-developed DL library offers optimized
 351 support for distributed training through model parallelism, data parallelism, and mixed-
 352 precision techniques, significantly enhancing training efficiency across GPUs.

353 3.3. Retrieval-augmented generation (RAG)

354 RAG is a technique that enhances LLM performance by integrating external
 355 knowledge into the generation process [56]. RAG consists of three key components:
 356 retrieval, augmentation, and generation. Specifically, it leverages retrievers to provide
 357 relevant context for LLMs, which then utilize their reasoning capabilities to decompose
 358 tasks, select appropriate tools, and generate responses [97–99]. Research indicates that

359 RAG substantially mitigates catastrophic forgetting caused by model weight updates,
 360 making it well-suited for domains requiring low tolerance for errors and rapidly
 361 evolving information. Compared to traditional FT, RAG allows timely incorporation of
 362 new medical knowledge without compromising previously learned information,
 363 thereby maintaining output accuracy in dynamic medical settings [50]. Notable
 364 implementations of RAG include QA-RAG [100], Almanac [101], Oncology-GPT-4
 365 [102], Impression GPT [103], and Retrieval-Augmented Lay Language Generation
 366 [104].

367 Finally, the three methods described above are all approaches to training LLMs.
 368 Compared to PT, FT significantly reduces computational and time costs. However, FT
 369 still requires additional model training with high-quality datasets, incurring
 370 considerable computational resources and manual effort. In contrast, RAG does not
 371 involve updating model parameters, making it a more efficient and convenient approach.
 372 Therefore, the choice of training method should be based on the specific requirements
 373 of different TCM tasks.

374 **4. TCM QA systems based on LLMs**

375 Nowadays, numerous LLMs have emerged successively, such as Qihuangwendao,
 376 Shuzhibencao, Huangdi, and Zhongjing. The development of these TCM-LLMs usually
 377 includes the following steps. First, the model is pre-trained through a large-scale general
 378 corpus to learn basic language knowledge and patterns. Then, on the basis of PT,
 379 professional data in the field of TCM is further used for FT. These professional data
 380 cover classical literature of TCM, case data, prescription compatibility, disease
 381 diagnosis and treatment, etc., aiming to endow the model with domain knowledge and
 382 context understanding ability of TCM. Finally, the fine-tuned LLM can be widely
 383 applied to TCM-related tasks, such as TCM diagnosis support, personalized treatment
 384 suggestions, drug recommendations, and medical QA [9,10,105]. Through this
 385 development process, the model can deeply master the terminology, knowledge
 386 systems, and diagnosis and treatment rules of TCM, thereby significantly improving its

387 performance in TCM-related tasks. This section will focus on introducing the currently
 388 developed TCM-LLMs and classifying them according to different task requirements
 389 and application scenarios. It will also focus on elaborating on specific aspects such as
 390 technical architecture, training dataset size, FT methods, and performance indicators.

391 **4.1. Huatuo**

392 HuaTuo is a TCM-LLM developed based on tuning LLaMA-7B architecture. It
 393 integrates both structured and unstructured knowledge from the Chinese Medical
 394 Knowledge Graph (CMeKG). Rather than a simple aggregation of resources, HuaTuo
 395 incorporates targeted architectural refinements informed by a comprehensive
 396 understanding of medical task requirements. As a result, it performs effectively in
 397 complex tasks such as medical question answering and the interpretation of domain-
 398 specific terminology. To facilitate SFT, over 8,000 high-quality, domain-specific
 399 instruction samples were curated by extracting knowledge instances from CMeKG and
 400 enhancing them using the OpenAI API. Although relatively limited in size, this dataset
 401 is highly specialized and contextually relevant, enabling the model to acquire
 402 professional knowledge efficiently, mitigate interference from general-purpose data,
 403 and improve its domain-specific reasoning and application capabilities.

404 During training, general-purpose instructions were excluded due to the specificity
 405 of the medical domain, and only the input components were retained. This design
 406 encourages the model to concentrate on learning knowledge and generating accurate,
 407 professional responses, thereby improving its practical utility. In the evaluation phase,
 408 HuaTuo introduced an evaluation metric called Safety, Usability, and Smoothness
 409 (SUS), with scores ranging from 1 (unacceptable) to 3 (good). Compared with models
 410 such as LLaMA, Alpaca, and ChatGLM, HuaTuo achieved superior performance,
 411 obtaining a safety score of 2.88, a knowledge usability score of 2.12, and a smoothness
 412 score of 2.47. These results indicate that HuaTuo performs well across multiple
 413 dimensions and holds significant potential for intelligent QA systems [106].

414 **4.2. Huangdi**

415 Huangdi is developed based on Ziya-LLaMA-13B-V1. Building on this
 416 foundation, it integrates corpora from TCM textbooks, various TCM websites, and
 417 other related data sources to construct a pre-trained language model with domain-
 418 specific knowledge in TCM. Utilizing extensive dialogue and instruction datasets, the
 419 model undergoes SFT, enabling it to effectively respond to questions related to classical
 420 TCM texts.

421 During the PT stage, data from "13th Five-Year Plan" TCM textbooks and folk
 422 medicine websites were used to establish a foundational understanding of TCM theory
 423 and clinical knowledge. The FT process consists of two parts: general IFT and
 424 instruction-based dialogue FT using classical TCM texts. The former uses 52k Chinese
 425 data from Alpaca-GPT4 to improve the generality of the model. The latter, based on
 426 TCM ancient books, constructs a professional dataset containing more than 500,000
 427 dialogue data, covering fields such as basic TCM theory, disease diagnosis, and the
 428 application of prescriptions and herbs.

429 In terms of data scale, the PT dataset is approximately 0.5 GB, and the processed
 430 ancient book dataset reaches 338 MB. These datasets together form diverse dialogue
 431 types, supporting the model in comprehensively mastering knowledge of TCM. FT
 432 adopts SFT and Direct Preference Optimization (DPO) optimization. The learning rate
 433 is set to 3×10^{-4} , lora rank is set to 16, and the number of training rounds is 6, enhancing
 434 the model's generation ability and user-friendliness. At the level of model performance
 435 evaluation, the train loss drops from 1.456 in the PT stage to 0.00276 in the DPO stage,
 436 and the eval loss is 1.258, indicating good generalization ability. Compared with Qwen,
 437 the Huangdi performs better in aspects such as TCM theory, diagnosis, and the
 438 application of prescriptions, demonstrating its application value and promotion
 439 potential in the TCM field.

440 4.3. Zhongjing

441 Zhongjing is developed by a research group at Zhengzhou University based on the
 442 Baichuan2-13B-Chat and Qwen1.5-1.8B-Chat models, with the goal of enhancing the

443 application capabilities of LLMs in the field of TCM. Its technical framework includes
 444 continuous PT, SFT, and reinforcement learning from human feedback (RLHF). During
 445 the continuous PT phase, Zhongjing was trained on diverse real-world medical text data
 446 from multiple sources, including electronic health records, medical consultation
 447 transcripts, textbooks, and other medical literature. In the next SFT stage, the model
 448 was trained using four types of instruction datasets. Among them, the Chinese multi-
 449 turn medical dialogue dataset, which includes 70,000 QA pairs across 14 medical
 450 departments and more than 10 real-world scenarios, significantly improved the model's
 451 conversational ability. During the RLHF stage, the researchers established a
 452 comprehensive set of annotation guidelines and enlisted six medical experts to rank
 453 20,000 model-generated sentences. A reward model was then trained using the
 454 Proximal Policy Optimization algorithm to better align the model's outputs with expert
 455 preferences.

456 To further optimize performance, the group avoided exclusive reliance on distilled
 457 data and carefully balanced the proportion of single-turn and multi-turn medical
 458 dialogues during the FT process. In terms of evaluation metrics, Zhongjing
 459 demonstrates strong performance across both single-turn and multi-turn interactions in
 460 three key dimensions: safety, professionalism, and fluency. In most scenarios,
 461 Zhongjing outperforms baseline models, and in multi-turn dialogue specifically, it
 462 surpasses all compared models except ChatGPT [9,13,107,108].

463 4.4. BianQue

464 Research has found that LLMs such as ChatGLM, ChatGPT, DoctorGLM, and
 465 ChatDoctor, perform well in generating general and widely applicable health
 466 suggestions in single-turn conversations. However, they exhibit a significant limitation:
 467 the lack of Chain-of-Thought (CoT) capability. Unlike real TCM experts, these models
 468 struggle to gather comprehensive patient information through continuous and
 469 interactive questioning. To alleviate this issue, BianQue has been developed. It has been
 470 developed based on the open-source ChatGLM-6B architecture, which features strong

471 Chinese language understanding and generation abilities. To support the training of the
 472 model, the group developed the BianQueCorpus, a health-related big dataset with tens
 473 of millions of samples. It integrated some multi-turn conversation datasets such as
 474 MedDialog-CN, IMCS-V2, and MedDG, and collected multi-turn health conversations
 475 (There are 243,7190 samples, of which 46.2% are the doctors' answers, and the rest are
 476 suggestions) in the real world through data outsourcing services. To ensure the quality
 477 of the dataset, researchers have developed an automated data cleaning mechanism based
 478 on regular expressions. In addition, the model uses ChatGPT to polish the doctors'
 479 suggestions for multi-turn conversations. This is because doctors often provide very
 480 brief responses through internet platforms, lacking detailed analysis and suggestions.

481 In terms of performance evaluation, the BianQue model performs outstandingly in
 482 multiple Chinese multi-turn medical dialogue datasets. It is evaluated using metrics
 483 such as Bilingual Evaluation Understudy (BLEU), Recall-Oriented Understudy for
 484 Gisting Evaluation (ROUGE), and the self-defined Proactive Questioning Ability
 485 (PQA). The results show that it outperforms baseline models like ChatGLM-6B,
 486 ChatGPT, and DoctorGLM in all indicators, demonstrating excellent generation quality
 487 and interactive questioning capabilities. Especially on the MedDG dataset, the BianQue
 488 has a remarkably high PQA value of 0.81, which proves that it has a strong proactive
 489 questioning ability and excellent medical dialogue generation performance [106,109].

490 4.5. TCMLLM-PR

491 In recent years, some open-source LLMs such as ChatGLM, ChatGPT, and
 492 LLaMA have all been trained on general-domain data, acquired good general-task
 493 processing capabilities, but their proficiency in specific medical fields is limited. To
 494 address this, researchers have successively developed specialized TCM-LLM.
 495 Although these LLMs have achieved certain progress, they still have not fully addressed
 496 the core issues in the TCM field, especially in the area of TCM prescription
 497 recommendation. To address the above challenges, researches proposed the TCMLLM-
 498 PR model, which is a LLM tailored to TCM prescription recommendation tasks using

499 the ChatGLM-6B architecture. The dataset for training TCMLLM-PR integrates multi-
 500 source heterogeneous information from eight channels, covering four TCM textbooks,
 501 *Pharmacopoeia of the People's Republic of China* (2020 Edition) [110], Chinese
 502 Medicine Clinical Cases, splenic-stomach disease, and hospital clinical records
 503 covering lung disease, diabetes, liver disease, stroke. TCMLLM-PR was then trained
 504 using the ChatGLM-6B architecture with P-Tuning v2 technology. On this basis, an
 505 instruction-tuning dataset containing 68,654 samples was constructed, with a total scale
 506 of approximately 10 M tokens. Ultimately, the model is able to accurately focus on the
 507 scenario of recommending TCM prescriptions.

508 In terms of evaluation metric, the model mainly uses the following indicators
 509 Precision@K, Recall@K, and F1 score@K. The experiment results demonstrate that
 510 TCMLLM-PR significantly outperforms baseline models on pharmacopoeia datasets
 511 and TCM textbooks, achieving F1@10 improvements of 59.48% and 31.80%,
 512 respectively. In the cross-dataset transfer task, it performed best when transferring from
 513 textbook data to the liver disease dataset, with F1@10 reaching 0.1551. The analysis of
 514 real-world cases further confirmed that this model performs outstandingly in the
 515 prescription recommendation task. The output results are highly consistent with the real
 516 doctors' prescriptions, demonstrating great potential for clinical applications [111].

517 4.6. TCMChat/TCMGPT

518 While LLMs have shown excellent performance in medical tasks like QA and
 519 diagnosis, TCM-LLMs (such as Ben Cao, Bian Que, HuaTuoGPT) still struggle with
 520 limited corpora, data inaccuracies, subjective evaluations, and inadequate tools. These
 521 problems have restricted their popularization and application. In this context, TCMChat
 522 emerged. The development of TCMChat begins with the utilization of the Baichuan2-
 523 7B chat base model. Moreover, it also adopts the typical transformer decoder
 524 architecture and has carried out a number of optimizations in the model structure. For
 525 example, it uses root mean square normalization (RMSNorm) instead of the
 526 normalization layer to achieve a more stable normalization process. It introduces rotary

527 positional encoding to enhance the sequence modeling ability, and replaces the
528 activation function with SwiGLU which has a stronger expressive ability. It is
529 committed to building a high-performance conversational large-scale model for TCM,
530 enhancing the application efficiency of artificial intelligence (AI) in the modernization
531 process of TCM, and providing technical support for the standardization and
532 popularization of TCM knowledge. To support model training, the research team has
533 constructed a large-scale and high-quality PT and SFT dataset. The PT data covers
534 multiple sources such as books (20 M), web crawlers (30 M), open-source data (352
535 M), and literature (715 M), and a total of approximately 1G of unsupervised corpus has
536 been constructed. The SFT dataset covers seven task scenarios, including TCM
537 knowledge QA, multiple-choice questions, reading comprehension, entity extraction,
538 medical case diagnosis, Chinese herbal medicines, formula recommendation, and
539 ADMET prediction, with a total of approximately 600,000 QA pairs.

540 TCMChat has parameters for two stages, PT and SFT. For the PT process, the
541 learning rate is 2×10^{-4} , the batch size is 32 per GPU, and the maximum context length
542 is 1024 tokens. For the SFT process, full-parameter FT is used. The learning rate has
543 been adjusted to 2×10^{-5} , the batch size per GPU is 16, and the maximum context length
544 is limited to 1,024 tokens. In addition, the model also uses the AdamW optimizer and
545 sets the weight decay to 1×10^{-4} to prevent overfitting. In terms of performance
546 evaluation, TCMChat demonstrates strong results across a range of tasks. For multiple-
547 choice questions related to herbs and formulas, the model achieves accuracy rates of
548 71.6% and 76.8%, respectively. In the reading comprehension task, it reaches a BLEU
549 score of 0.584 and a BertScore as high as 0.886. The entity extraction task yields an
550 impressive F1 score of 0.907. In medical case diagnosis, the model achieves an
551 accuracy of 0.847. For herb and formula recommendation, it records a Mean Reciprocal
552 Rank of 0.536 and a Normalized Discounted Cumulative Gain of 0.439. Lastly, in the
553 ADMET prediction task, TCMChat attains a classification accuracy of 0.818 and a
554 ROC-AUC of 0.830. Overall, TCMChat demonstrates powerful understanding and

555 generation capabilities in multiple sub-tasks, providing a strong support for AI systems
 556 in TCM [16].

557 **4.7. Qibo**

558 Qibo is a model based on LLaMA. During training, it first acquired the basic
 559 knowledge and the theoretical framework of TCM through continuous PT, developing
 560 capabilities in comprehension, dialectical analysis, and entity recognition of TCM.
 561 Subsequently, SFT with diverse datasets was employed to improve its dialogue and
 562 instruction following abilities. To enhance TCM inquiry and syndrome differentiation,
 563 Qibo incorporates a retrieval-enhanced approach using an external knowledge base.
 564 Furthermore, a CoT mechanism is implemented to simulate the real-world TCM
 565 consultation process, enabling multi-turn information integration and informed
 566 decision-making in prescription generation.

567 The PT data includes modern medical textbooks, TCM reading comprehension
 568 materials, TCM textbooks, TCM prescriptions, and so on. With a total size of
 569 approximately 2 GB, this dataset provides the model with a rich and comprehensive
 570 knowledge base of TCM. The SFT dataset contains seven types of tasks such as TCM
 571 QA, reading comprehension, and prescription recommendation, totaling approximately
 572 600,000 QA pairs. During the training, four types of data are used and converted into
 573 the Alpaca format, including single-round conversations, multi-turn conversations in
 574 TCM departments, NLP instruction tasks, and general medical conversations,
 575 comprehensively improving the generalization and robustness of the model in the fields
 576 of TCM.

577 In terms of model evaluation, Qibo demonstrates significant advantages in
 578 multiple dimensions. In subjective evaluation, the model performs outstandingly in
 579 terms of professionalism, safety, and fluency compared to the baseline models. Qibo-
 580 7B has an average subjective win rate of 63% on 150 TCM questions, and it particularly
 581 has a clear advantage in the safety dimension. The objective evaluation takes the form
 582 of multiple-choice questions. Based on the test of 3,175 questions related to TCM

583 professional practice examinations, Qibo's accuracy has increased by 23%-58%
 584 compared to the baseline models. In addition, in TCM NLP tasks, such as Entity
 585 Recognition Ability, TCM Reading Comprehension, and TCM Syndrome
 586 Differentiation. Qibo has achieved ROUGE-Longest (ROUGE-L) scores of 0.72, 0.61,
 587 and 0.55 respectively. Although it still does not reach the optimal performance of task-
 588 specific models. overall, it outperforms existing medical large language models, fully
 589 demonstrating its potential in TCM language understanding and application scenarios
 590 [112].

591 4.8. Lingdan

592 Based on the Baichuan2-13B-Base model, continuous PT is carried out to obtain
 593 the Lingdan model. Furthermore, the Lingdan Traditional Chinese Patent Medicine
 594 Chat (Lingdan-TCPM-Chat) model and the Lingdan-prescription recommendation
 595 (Lingdan-PR) model are developed. Experimental results show that these two models
 596 perform excellently in the tasks of TCM clinical knowledge answering and herbal
 597 prescription recommendation. Among them, the Lingdan-PR model has improved by
 598 18.39% in the Top@20 F1 metric compared with the best baseline model, providing
 599 strong support for promoting the integrated development of TCM and AI. To support
 600 the construction of the above models, the research team has created a large-scale TCM
 601 Pre-trained dataset that covers multi-source content such as ancient TCM books and
 602 textbooks. They have also used the Baichuan2-13B-Base model to translate ancient
 603 TCM texts and have linguistically processed the information about TCM herbs.

604 During the training stage, the Baichuan2-13B-Base model was used as the base
 605 model, and training was conducted on 6 NVIDIA A100-80G GPUs using Quantized-
 606 LoRA (QLoRA) and Zero Redundancy Optimizer. Training balance was achieved by
 607 configuring LoRA parameters and applying diverse data sampling strategies. Regarding
 608 the Lingdan-TCPM-Chat model, 200,000 single-turn dialogue data have been generated
 609 through knowledge QA transformation, and 1,599 multi-turn consultation data have
 610 been constructed based on the TCM Interactive Diagnostic Dialogue Framework. As

611 for the Lingdan-PR, a spleen and stomach herbal prescription recommendation dataset
 612 has been constructed based on the data from the Department of Spleen and Stomach
 613 Diseases in Guang'anmen Hospital. After data augmentation, FT has been carried out
 614 on two pre-trained models respectively, and finally, it has comprehensively
 615 outperformed existing baseline models in multiple Top@K indicators, verifying the
 616 effectiveness and robustness of the method.

617 The Lingdan outperforms the baseline models in multiple evaluation indicators. In
 618 the F1@5 indicator, it has increased by up to 5.82% compared to the highest score of
 619 the baseline models. In the F1@10 indicator, it has increased by 11.89%; and in the
 620 F1@20 indicator, it has increased by 18.39%. This shows that the Lingdan-PR model
 621 has better prediction accuracy, recall rate, and comprehensive performance in the task
 622 of TCM prescription recommendation [113].

623 4.9. Biancang

624 Biancang is a TCM-specific LLM. It uses a two-stage training process. First, it
 625 injected domain-specific knowledge, and then aligns it through targeted stimulation.
 626 During the FT phase, four major methods are employed. The first method is SFT, which
 627 responds to TCM task instructions through structured QA for training models. The
 628 second method is RLHF, which optimizes the model output based on preference scoring.
 629 The third method is knowledge enhancement, which integrating the structural
 630 information of TCM knowledge graph into the process of contextual understanding.
 631 The fourth method is multi-task learning, which integrates tasks such as QA,
 632 summarization, translation, and dialogue to enhance the model's generalization ability.
 633 In multiple evaluation tasks, Biancang demonstrates outstanding performance. In the
 634 Chinese medical exam, its accuracy rate reaches 94.97%, outperforming GPT-4
 635 (82.69%) and Qwen2-7B-Instruct (83.35%). Biancang-Qwen2.5-7B-Instruction has
 636 achieved further improvements, reaching an accuracy rate of 82.10% on the TCM
 637 syndrome differentiation test set. At the same time, the manual evaluation by TCM
 638 experts has also verified Biancang's leading performance in terms of professionalism,

639 fluency, and security, demonstrating its strong potential in the application of TCM-
 640 LLM [17].

641 4.10. Haiheqibo

642 Haiheqibo is a knowledge graph model developed for the TCM domain. The
 643 construction process of it can be divided into three main stages. Firstly, in the data
 644 processing stage, data is collected from medical books, open-source datasets, and
 645 crawled data (Wikipedia and Baidu encyclopedia). The data undergoes through
 646 processes such as unified formatting (format preprocessing), noise removal (data
 647 cleaning), sample deduplication, and quality assessment to ensure the accuracy and
 648 diversity of the training data. At each stage, sampled data is manually evaluated to
 649 enhance overall data quality. Subsequently, in the PT stage, the above high-quality data
 650 is used to construct a PT dataset. Continuous PT is then carried out based on the LLaMA
 651 model to generate the TCM-Base Model, which has basic knowledge of TCM but lacks
 652 dialogue ability. Then, in the FT stage, data including multi-turn dialogues, single-turn
 653 QA, and various NLP tasks are converted into instruction-style formats to construct an
 654 instruction dataset. Afterward, SFT is performed on the TCM-Base Model using this
 655 dataset to enhance its ability to follow task-specific instructions. Finally, Haiheqibo
 656 with both TCM knowledge and dialogue ability is obtained.

657 4.11. Congbaosuwen

658 On November 2, 2024, the 5.0 version of the Congbaosuwen was officially
 659 released. First, it can accurately understand users' needs. Whether it is professional
 660 Chinese medicine consultation or the public's inquiry about health preservation
 661 knowledge, it can achieve efficient interaction and enhance the user experience. Second,
 662 it features high flexibility. It supports specific vertical-field scenarios and multi-modal
 663 data. It can be optimized for scenarios such as clinical practice and teaching. For
 664 example, in tongue diagnosis, it can combine with images for assisted diagnosis. Third,
 665 its security has been significantly improved. The self-reflection mechanism ensures the
 666 accuracy and compliance of the content. Finally, it is compatible API service by

667 OpenAI, which simplifies development and deployment and accelerates the launch of
668 TCM products.

669 4.12. Pangu

670 In 2024, Zhejiang Jiuwei Health Technology Co., Ltd. and Huawei Cloud
671 Computing Technology Co., Ltd. jointly launched the PanGu. This model is a large-
672 scale pre-trained model based on DL technology, specifically designed and optimized
673 for the field of TCM. This model is trained using massive amounts of TCM data,
674 enabling it to deeply understand the language and culture of TCM, providing strong
675 support for the research, development, and application of TCM. At the level of data
676 quality, the PanGu integrates various types of data such as classic TCM literature, TCM
677 prescriptions, medicinal material information, and clinical cases, forming a vast and
678 comprehensive TCM knowledge base. These data not only cover all aspects of TCM
679 but also have been carefully cleaned and annotated to ensure data quality and accuracy.
680 In terms of technology, the PanGu adopts the transformer architecture in DL, which is
681 a neural network structure with powerful feature extraction and context-understanding
682 capabilities. Through large-scale PT, the model can automatically learn the complex
683 knowledge and patterns in the field of TCM, providing a solid foundation for
684 subsequent applications.

685 In application, it shows broad prospects and potential. First, in the
686 recommendation of TCM prescriptions, the model can intelligently recommend
687 personalized TCM prescriptions based on the patient's symptoms and constitution,
688 improving the accuracy and effectiveness of TCM treatment. Second, in the quality
689 control of medicinal materials, the model can assist in identifying the authenticity and
690 quality of medicinal materials by analyzing information such as the characteristics,
691 origin, and harvesting time of the medicinal materials, ensuring the quality and safety
692 of the medicinal materials. In addition, this model can also play an important role in
693 auxiliary disease diagnosis, new drug research and development, and health
694 management.

695 4.13. Tianhelingshu

696 Tianhelingshu is a professional LLM designed for the field of TCM acupuncture
697 and moxibustion. It is built on professional data, including classic Chinese medicine
698 works, acupuncture and moxibustion clinical practice evidence-based database, and
699 TCM evidence-based knowledge map. This model has systematically studied hundreds
700 of classic TCM works and been trained on tens of thousands of pieces of evidence-
701 based data. It has profound knowledge of TCM theory and can serve as an intelligent
702 assistant for TCM to provide users with accurate and professional answers. Whether it
703 is an in-depth discussion of TCM theory or a detailed analysis of health problems, the
704 model can quickly give detailed responses. When users seek advice on acupuncture and
705 moxibustion treatment, it can rapidly analyze users' conditions and put forward
706 personalized suggestions, including various acupuncture and moxibustion treatment
707 methods such as acupuncture, moxibustion, and acupressure.

708 4.14. Hengqin

709 Hengqin aggregates a vast amount of TCM data, including 10 B characters of TCM
710 knowledge texts and digital cases from TCM hospitals. Relying on a highly reliable
711 TCM diagnosis and treatment knowledge base, it assists doctors in accurate diagnosis
712 and treatment and provides personalized treatment plans. The Intelligent and
713 Automated Integrated Innovation Platform for New TCM Drugs, through engineering
714 development, realizes a one-stop solution for the entire experimental process of TCM
715 ingredient acquisition, structural characterization, and bioactivity determination based
716 on robotics and automation technologies.

717 Recently, Hengqin Rheumatoid Arthritis v1 (Hengqin-RA-v1) has been developed.
718 It is specifically designed for the diagnosis and treatment of RA. This model adopts a
719 progressive training workflow, optimizing RA-specific datasets while preserving
720 existing knowledge. Hengqin-RA-v1 integrates domain-specific knowledge through
721 segmented structured data and enhances model performance using instance-oriented
722 and entity-relationship-oriented retrieval enhancements. A sliding window strategy is

723 also employed during training to refine contextual logic and improve the model's
 724 understanding of the complex diagnostic context in TCM. Hengqin-RA-v1 outperforms
 725 other LLMs in the medical domain, achieving an accuracy rate of 54% in TCM
 726 examinations. This result significant outperformed better than both Chinese and non-
 727 Chinese models. It excels in generating diagnostic recommendations and treatment
 728 plans for rheumatoid arthritis, surpassing traditional approaches in certain scenarios and
 729 even outperforming human expert evaluations in some diagnostic cases [114].

730 4.15. Shennong

731 Shennong is jointly completed by the Intelligent Knowledge Management and
 732 Service Team of the School of Computer Science and Technology at East China Normal
 733 University. It aims to promote the development and implementation of large models in
 734 the field of TCM and enhance the knowledge of large models in TCM and their ability
 735 to answer medical consultations. Shennong is obtained by using FT with LoRA
 736 (rank=16). It is based on an open-source knowledge graph of TCM, with LLaMA
 737 serving as the base model. Through an entity-centered self-instruction method, more
 738 than 110,000 TCM instruction data is obtained by calling ChatGPT, promoting the
 739 inheritance of TCM empowered by large models. However, there are also some
 740 shortcomings at the same time. For example, the data relies on an open-source
 741 knowledge graph of TCM.

742 Compared with the Chinese LLaMA-7B model, the Shennong demonstrates
 743 superior overall performance in TCM QA tasks. First, Shennong exhibits more natural
 744 and human-centered language expression; its responses not only address the patient's
 745 condition but also convey empathy and attention to emotional needs, thereby enhancing
 746 the user experience. Second, owing to large-scale training on TCM-specific
 747 instructional data, the model possesses a strong foundation in domain knowledge. It can
 748 deliver detailed, actionable treatment suggestions tailored to specific symptoms,
 749 including common herbal formulas, methods of administration, and relevant
 750 precautions. In contrast, the Chinese LLaMA-7B model, lacking domain-specific

751 optimization, tends to produce relatively brief and generic responses, with limited
752 professionalism, which makes it insufficient for practical applications in TCM
753 diagnostic and therapeutic contexts. Consequently, Shennong is better suited for
754 intelligent diagnosis and decision-making tasks in TCM.

755 4.16. Shuzhibencao

756 Shuzhibencao was jointly developed by Huawei Cloud and Tasly Pharmaceutical
757 Group Co., Ltd. This model integrates an extensive database, including over 1,000
758 ancient texts and their translations, more than 90,000 traditional prescriptions, upwards
759 of 40,000 Chinese patent medicines, over 40 M literature abstracts, more than 3 M
760 natural products, data on over 20,000 target gene pathways, more than 100,000 clinical
761 treatment protocols, over 160,000 Chinese medicine patents, and a wide range of
762 pharmacopoeia and policy guidelines. This model has 38 B parameters and has been
763 pre-trained on a vast corpus of TCM texts. By integrating and combining the
764 reinforcement of vector library retrieval and FT in multiple scenarios of Chinese
765 medicine research and development, it can better assist researchers in mining and
766 summarizing the theoretical evidence of TCM.

767 Shuzhibencao was pre-trained based on billions of molecular structures and further
768 fine-tuned using 3.5 M unique natural product molecules. This enables it to more
769 accurately perform computational tasks, such as characterizing natural product
770 structures and improving the prediction of their downstream properties. By integrating
771 with appropriate algorithms, the model can also accelerate the screening and
772 optimization of medicinal materials and compound prescriptions [115].

773 4.17. Bencaozhiku

774 BenCaozhiku was released on 2024, during the second "Thousand Herbal
775 Genomes Project" symposium. This model integrates core foundational data for TCM
776 research, including 15 M genomic records of source species for medicinal herbs, over
777 30 M records on interactions between TCM compounds and their targets, and more than
778 4 M compounds. It forms a knowledge graph comprising over 20 M entities and more

779 than 2 B relationships, covering the entire TCM industry chain. Powered by the
 780 Wenxinyiyan LLM with hundreds of billions of parameters, and enhanced through
 781 instruction tuning and RAG techniques, the model supports three key functions:
 782 extraction and generation of TCM knowledge, delivery of domain-specific TCM
 783 solutions, and comprehensive digital services for the TCM industry. It achieves
 784 seamless integration of foundational research data with critical stages across the TCM
 785 value chain.

786 **4.18. Qihuangwendao**

787 In July 2023, Baidu Health and GuShengTang Incorporated have released the
 788 Qihuangwendao. Based on the training of TCM knowledge, this model takes more than
 789 1,000 ancient Chinese medical books and TCM documents such as *Huangdi Neijing*
 790 and *Treatise on Cold Damage and Miscellaneous Diseases* as its core data foundation.
 791 It covers 11 M pieces of data in the knowledge graph of TCM, 2 M pieces of real clinical
 792 diagnosis and treatment data of TCM, 100 thousand pieces of real medical case data of
 793 TCM experts, and 100 thousand pieces of data on pulse conditions, tongue
 794 manifestations, meridians, and acupoints. It has efficient operation capabilities and
 795 accurate syndrome differentiation capabilities, and has achieved a high level of
 796 specialization.

797 Qihuangwendao consists of two application functionalities: the medical LLM and
 798 the health preservation LLM. The medical LLM is further divided into two sub-models:
 799 providing prescriptions for confirmed diagnoses and conducting diagnostic assessments
 800 based on symptoms. Through the construction technology of the knowledge graph of
 801 the experience of famous veteran TCM doctors, their experiences are sorted out to
 802 provide knowledge support for the model. The natural language recognition technology
 803 of TCM is applied to improve the ability to interpret text related to diseases and
 804 symptoms. With the help of the technology for constructing standardized symptoms
 805 and signs in TCM, expressions are standardized for accurate analysis. The big data
 806 mining technology of TCM diagnosis and treatment is adopted to extract key

807 information from massive data and conduct training. Ultimately, the model can
 808 accurately diagnose complex diseases, conduct syndrome differentiation of symptoms,
 809 and make professional judgments. It can customize personalized treatment plans for
 810 different users, such as recommending TCM prescriptions, guiding meridian massage,
 811 and providing dietary therapy suggestions. The health preservation LLM, creating
 812 personalized multidimensional wellness plans to help maintain health and prevent
 813 diseases [9,13].

814 4.19. Shuzhiquihuang 2.0

815 On November 25 in 2024, at the fourth "Big Data and AI in TCM" Shanghai
 816 Forum, Northeast Normal University, in collaboration with multiple institutions,
 817 released the Shuzhiquihuang 2.0 multi-modal large model in the field of TCM. This
 818 model contains 32 B parameters and covers two main modules of TCM and western
 819 medicine: it has more than 200,000 and 100,000 instruction data respectively. In
 820 addition, the model covers over 80,000 TCM prescriptions, more than 40,000 TCM
 821 ingredients, encompasses 9,000+ kinds of TCM materials, incorporates 2,000+ TCM
 822 syndromes, and contains 1,000+ ancient books; including 18,000+ targets, 2,000+
 823 diseases, 2.4 M compounds (of which about 410,000 are natural products), and more
 824 than 2 M documents. With its advanced multi-modal capabilities and extensive
 825 knowledge base coverage, Shuzhiquihuang 2.0 shows a new breakthrough in the field of
 826 intelligent diagnosis and treatment.

827 4.20. MedChatZH

828 MedChatZH is a dialogue model specifically designed and optimized for TCM
 829 consultations, demonstrating significant advantages in aspects such as technical
 830 architecture, training data, and FT methods. This model is constructed based on the
 831 Baichuan-7B (similar to LLaMa). RMSNorm is adopted to normalize the input of each
 832 sublayer, which improves the stability of the training process and the effect of layer
 833 normalization. By incorporating the Rotational Position Embedding mechanism, it
 834 integrates relative and absolute positional information, thus significantly enhancing the

835 model's generalization and understanding capabilities for different text lengths and
 836 structures. At the data level, the PT dataset of MedChatZH encompasses more than
 837 1,000 TCM ancient books and modern books, including *Treatise on Febrile Diseases*,
 838 *the Yellow Emperor's Canon of Internal Medicine*, and *A Barefoot Doctor's Manual*. It
 839 covers a wide range of TCM theories and clinical experiences. In addition, the medical
 840 IFT dataset (med-mix-2M) is introduced, which contains 763,629 medical instructions
 841 and 1,305,194 general instructions. The data sources integrate multiple projects such as
 842 belle-3.5M, medical, medical-dialogue, and so on, laying a solid foundation for the
 843 application of the model in professional and general dialogue scenarios.

844 In terms of the FT, MedChatZH addresses the differences between the language
 845 style of TCM books and the requirements of modern conversations. First, use Baidu
 846 Wenxinyiyan to convert ancient Chinese texts into modern Chinese, and then leverage
 847 the ChatGPT API to optimize the translation quality, constructing a high-quality PT
 848 corpus. The model undergoes complex data processing steps, including heuristic and
 849 model-based filtering to remove irrelevant or sensitive content. Medical instruction data
 850 further undergoes a three-fold processing: First, personal privacy information is
 851 removed through regular expressions. Second, the trained Ziya-LLaMA-7B-Reward
 852 model is used to score the data, and low-quality samples with a score lower than 0.5 are
 853 eliminated. Third, the format of numerical symbols is unified to enhance data
 854 standardization. In the FT stage, reinforcement learning (RL) techniques is adopted to
 855 convert the instruction data into a QA template in the "Human-Assistant" format. Only
 856 the answer part generated by the model is used to calculate the loss and update. In the
 857 webMedQA dataset test, MedChatZH performs outstandingly on automatic evaluation
 858 metrics such as BLEU, GLEU, and ROUGE. Under the BLEU-1 metric, MedChatZH
 859 scores 56.14, while ChatGLM-Med scores 32.18 and BenTsao scores 32.02. In terms
 860 of the ROUGE-L metric, MedChatZH reaches 35.99, ChatGLM-Med is 26.14, and
 861 BenTsao is only 17.72. This indicates that the answers generated by MedChatZH are
 862 superior to models such as ChatGLM-Med and BenTsao in terms of similarity to

863 reference sentences, fluency, and quality evaluation based on word matching [3].

864 4.21. Chinese patent medicine instructions-ChatGLM (CPMI-ChatGLM)

865 CPMI-ChatGLM is based on the ChatGLM-6B model and adopts a prefix decoder-
 866 only transformer framework. It integrates the bidirectional and unidirectional attention
 867 mechanisms, which are used to process input and output information respectively. The
 868 model ensures the stability of the training process through the gradient scaling
 869 embedding layer and the post-LN layer normalization method. At the same time, it
 870 introduces RoPE to replace the traditional absolute position encoding, and adopts the
 871 GeLU activation function in the FFNs to enhance the expressive ability and
 872 generalization performance. In terms of training data, the core dataset of CPMI-
 873 ChatGLM is derived from the Entity Recognition of TCM, *Standard Therapeutic*
 874 *Guidelines for National Essential Drugs*, and instructions. After data cleaning,
 875 denoising, and the removal of drugs with unknown attributes, data augmentation is
 876 achieved in combination with ChatGLM. Eventually, a dataset containing 3,906 high-
 877 quality data records is constructed.

878 In terms of FT methods, CPMI-ChatGLM applies the parameter-efficient FT
 879 technology, focusing on comparing two methods, LoRA and P-Tuning v2. LoRA
 880 optimizes model parameters by introducing low-rank matrices to assist in the update
 881 process, while P-Tuning v2 optimizes performance by introducing continuous prompts
 882 at each layer of the model. Experiments show that FT based on P-Tuning v2 performs
 883 better in multiple evaluation metrics. For example, ROUGE-1 F1 score is 33.81%
 884 higher than that of LoRA. In addition, the model is also fine-tuned with instruction data
 885 containing TCM knowledge. For domain-specific tasks, a small set of instruction data
 886 is often sufficient to guide generation, in contrast to general-domain models. This
 887 strategy improves the performance of the model in tasks such as the recommendation.
 888 For performance evaluation, the model comprehensively assesses the similarity and
 889 quality of the generated text using automatic evaluation metrics such as BLEU,
 890 ROUGE, and BARTScore metrics. It achieves a score of 0.7641 on the BLEU-4 metric,

891 demonstrating excellent performance. Regarding manual evaluation, the SUS standard
 892 is introduced to evaluate the reliability and practicality of the content generated by the
 893 model from three aspects: safety, usability, and fluency. The results show that CPMI-
 894 ChatGLM scores highly in all three indicators, indicating that it has good application
 895 prospects in the field of Chinese herbal medicine instruction manuals [116].

896 **4.22. MING-mixture-of-expert (MING-MOE)**

897 MING-MOE, a medical LLM based on MOE, is designed to manage diverse and
 898 complex medical tasks. It does not require task-specific annotations, thus improving its
 899 usability across extensive datasets. MING-MOE adopts a Mixture of Low-Rank
 900 Adaptation technique. This technique allows for efficient parameter utilization by
 901 keeping the basic model parameters static while enabling adaptation through a minimal
 902 set of trainable parameters. Literature demonstrate that MING-MOE achieves state-of-
 903 the-art performance in more than 20 medical tasks, indicating a significant
 904 improvement over existing models. This approach not only expands the capabilities of
 905 medical language models but also enhances the inference efficiency. MING-MOE is a
 906 bilingual (Chinese and English) medical LLM built upon the foundation of the
 907 Qwen1.5-Chat. The FT process mainly focuses on enhancing the model's stability in
 908 handling specific tasks in the medical field. At the same time, a certain proportion of
 909 medical QA and interaction data are retained in the FT dataset reserve for the model's
 910 general capabilities. Based on the size of the base model, four different sizes were
 911 proposed, including MING-MOE (1.8 B), MING-MOE (4 B), MING-MOE (7 B), and
 912 MING-MOE (14 B).

913 The results of the automatic evaluation metrics showed that two studies reported
 914 accuracy in medical tasks. For TCM diagnosis, TCM-GPT achieved an accuracy of
 915 0.264, while MING-MOE achieved 41.58 (1.8 B), 50.31 (4 B), 57.03 (7 B), and 63.2
 916 (14 B). In TCM examinations, the reported accuracy was 0.29 for TCM-GPT and 33.96
 917 (1.8 B), 45.00 (4 B), 49.58 (7 B), and 59.79 (14 B) for MING-MOE. In summary,
 918 MING-MOE provides accurate answers and reasonable interpretations for a variety of

919 tasks, demonstrating strong capabilities in knowledge application and problem-solving,
 920 and has high application value in the field of TCM [117].

921 4.23. IvyGPT

922 IvyGPT is built on the LLM architecture and adopts a two-stage training strategy
 923 to form its technical system. In the first stage SFT process, it is optimized based on
 924 LLaMA-33B. By introducing the LoRA method, the weights of the original model are
 925 frozen, and trainable low-rank matrices are injected into each layer of the transformer.
 926 This significantly reduces the number of trainable parameters required in downstream
 927 tasks. On this basis, the QLoRA technology is further adopted to perform 4-bit
 928 quantization on the base model. This enables the efficient FT of large-scale models even
 929 in low memory resource environments, significantly improving the training efficiency
 930 and adaptability. In the second stage, RLHF is introduced to optimize the model. During
 931 the training process, the model generates responses based on real instructions and multi-
 932 turn conversations. Then, human evaluators score the responses from four dimensions:
 933 considering informativeness, coherence, adherence to human preferences, and factual
 934 accuracy, which is used to guide the learning of the reward model. At the same time, to
 935 prevent the model trained through RL from deviating too far from the model in the SFT
 936 stage. Constraints are introduced during the training process to ensure that the model
 937 achieves better results without deviating from the track. In addition, the top k responses
 938 generated by the model are scored, and by defining a reward function, the model is
 939 guided to generate answers that are more in line with human preferences, thus achieving
 940 the unity of performance and stability.

941 Compared with models such as HuaTuo, Shennong, ChatMed, MedicalGPT, and
 942 ChatGPT, IvyGPT achieved the highest semantic similarity score (93.58) in 100 QA
 943 tasks, demonstrating stronger semantic understanding and expression capabilities.
 944 Moreover, the average word count of responses generated by IvyGPT is 271.05, which
 945 is higher than that of other current models, indicating that its answers are more
 946 informative. Overall, IvyGPT performs outstandingly in terms of training efficiency,

947 output quality, and semantic consistency, showing great potential for medical
 948 applications [118].

949 Furthermore, we have also conducted a systematic compare of TCM-LLMs (Table
 950 S1). In practical applications, the selection of TCM-LLMs should be approached from
 951 multiple dimensions. First, it is essential to define the specific application scenario. For
 952 diagnostic reasoning tasks, models such as Zhongjing and Qibo demonstrate strong
 953 performance. For prescription recommendation, TCMLLM-PR and Lingdan are more
 954 suitable. In QA and consultation contexts, Huatuo, MedChatZH, and Biancang are
 955 recommended. Open-source models like GLM-130B, Huatuo, and TCMChat should be
 956 prioritized, as they facilitate local deployment and personalized FT, particularly
 957 important in handling sensitive data. Model adaptability should also be considered in
 958 relation to training methods. For instance, Huatuo and TCMLLM-PR employ IFT,
 959 while IvyGPT and Biancang incorporate RLHF. Regarding evaluation, Biancang
 960 achieves an accuracy of 94.97% in chronic disease management, and the BLEU score
 961 of MedChatZH also serves as a useful performance indicator. For tasks involving
 962 multimodal input and knowledge graph integration, models such as CPMI-ChatGLM
 963 and Shuzhiqihuang 2.0 are optimal. Ensuring the stable operation of deployed models
 964 through continuous optimization is also recommended.

965 < **Table S1.** Summary of existing intelligence traditional Chinese medicine (TCM)
 966 intelligent question-answering (QA) systems based on large language models (LLMs)
 967 in terms of open-source availability, model development, parameter size, training
 968 methods, features, and data sources. >

969 **5. Prospects for the application of TCM-LLMs**

970 According to the information, as of November 30, 2022, more than 100 M users
 971 have retrieved LLMs and medicine in PubMed. The number of included articles
 972 increased from 182 in 2020 to 867 in 2023. Thus, it can be seen that large models have
 973 become a new direction to help the development of TCM [9]. As a medical auxiliary
 974 tool, large models in the field of TCM have shown great application potential in medical

975 education, new drug research and development, and medical clinical practice. The
 976 application prospects in the field of TCM are very broad. In addition, we have drawn a
 977 figure illustrating the application prospects of LLMs in the field of TCM (Fig. 5).

978 <**Fig. 5.** Prospects for the application of traditional Chinese medicine (TCM)-large
 979 language models (LLMs).>

980 5.1 Medical education

981 5.1.1 Learning ancient books

982 Ancient TCM books not only record the theoretical innovations, clinical
 983 experiences and medical techniques of ancient practitioners, but also provide valuable
 984 references and inspirations for modern TCM diagnosis and treatment of various
 985 diseases. However, due to the complexity of content, systematically absorbing the vast
 986 information in these ancient books has always been a challenge. For this reason, many
 987 researchers are focusing on using advanced technologies to conduct in-depth studies on
 988 ancient texts. LLMs can establish a knowledge base of TCM through learning and
 989 analyzing TCM literature, providing researchers with more comprehensive and
 990 accurate knowledge of TCM. Since TCM research requires a large number literature as
 991 support, and these literature materials are often scattered in various fields and
 992 disciplines, LLMs can integrate these scattered literature materials to establish a
 993 comprehensive and accurate knowledge base of TCM [13, 119, 120].

994 LLMs achieve comprehensive learning and in-depth exploration of TCM
 995 knowledge by integrating a large amount of TCM literature, ancient books, medical
 996 records, and modern TCM research results. At the same time, they cover core TCM
 997 theories such as the theory of yin and yang, the five elements theory, syndrome
 998 differentiation and treatment, and meridian theory, as well as treatment plans for
 999 different diseases. The Huangdi constructed by Zhang and co-workers [121] and others
 1000 provides users with knowledge services in multiple aspects such as answering questions
 1001 about ancient books, TCM consultations, treatment suggestions, and preventive health
 1002 preservation through natural language conversations. This model not only deeply

1003 excavates the knowledge value contained in ancient books but also explores a new
1004 paradigm for the research and application of TCM ancient books.

1005 5.1.2 Intelligent teaching

1006 Compared with traditional teaching methods, LLMs can automatically generate
1007 the required teaching materials and even provide virtual patient cases for teachers and
1008 students to conduct discussion, analysis, and simulated diagnosis. At the same time,
1009 teachers can use large models to formulate targeted teaching plans and realize the
1010 personalization of TCM teaching. By inputting individual information such as students'
1011 strengths, weaknesses, learning goals, and preferences, the model can generate
1012 personalized learning plans that meet the specific needs of students, thus providing a
1013 basis for precise tutoring and significantly improving teaching efficiency [122]. In
1014 language learning, LLMs can simulate conversations in multiple languages, correct
1015 grammar, increase vocabulary, and help improve pronunciation to meet needs [123].

1016 5.2. Drug development

1017 New drug research and development is crucial for promoting the development of
1018 China's healthcare industry. It not only fills the gaps of existing drugs but also meets
1019 unmet treatment needs. Drug research and development can be roughly divided into
1020 four stages: (1) Target Selection and Validation. Drug research and development
1021 usually starts from identifying targets related to specific diseases. At this stage, multiple
1022 technical means such as cell and gene target assessment, genome and proteome analysis,
1023 and bioinformatics prediction need to be combined. (2) Compound Screening and Lead
1024 Optimization. By identifying active compounds for screening, methods such as
1025 combinatorial chemistry, high-throughput screening, and virtual screening are usually
1026 employed to find candidate compounds from public databases (such as PubChem,
1027 ChemBL Database, etc.) (3) Preclinical Research. Through structure-activity
1028 relationship research and computer simulation technology, combined with cell function
1029 testing, repeated iterative optimization is carried out to improve the functional
1030 performance of newly synthesized candidate drugs. (4) Clinical Trials. Before entering

1031 the preclinical research stage, candidate compounds need to undergo in vivo testing in
1032 animal models, including pharmacokinetic studies and toxicity tests, to ensure their
1033 basic safety and effectiveness [124–127].

1034 According to relevant literature, new drug research and development has the
1035 largest share in the global medical AI market, reaching 35% [128]. However, new drugs
1036 also face many difficulties in the research and development stage, including long
1037 research and development cycles, high costs, and low success rates [129]. Therefore,
1038 integrating LLMs into the field of drug discovery and development marks a major
1039 paradigm shift and provides new methods for understanding disease mechanisms,
1040 promoting drug discovery, and optimizing the clinical trial process [4]. Traditional drug
1041 screening methods include molecular docking, pharmacophore matching, and similarity
1042 search [128]. Different from the traditional drug research and development process, AI
1043 has been widely used in all aspects of new drug research and development through key
1044 AI technologies such as NLP, ML, DL, and knowledge graphs. including drug target
1045 identification, active compound screening, compound property prediction, molecular
1046 generation, and protein structure and protein-ligand interaction prediction [130,131].
1047 The TCM-LLMs can quickly discover the connection between drugs and diseases, and
1048 between diseases and genes through DL technology, thereby shortening the drug
1049 research and development cycle, reducing drug research and development costs, and
1050 reducing the risk of new drug research and development [132,133]. At the same time,
1051 LLMs can rapidly screen out required TCM articles through automated literature
1052 retrieval and analysis. They can extract key information and potential drug candidates,
1053 synthesize innovative new drug prescriptions. Additionally, they can extract available
1054 data from experimental reports and generate initial report texts to assist in designing
1055 experimental schemes and enhance research efficiency [11].

1056 5.3. Clinical diagnosis and treatment

1057 5.3.1 Clinical data processing

1058 These big models possess highly advanced computational capabilities. They can
 1059 rapidly and accurately search through a large amount of medical data, collect the latest
 1060 information on diseases, and provide reliable scientific bases for researchers. They can
 1061 conduct precise case analyses and thus more quickly raise the level of medical research.
 1062 LLMs can automatically extract and deeply understand massive amounts of medical
 1063 data from medical literature, electronic medical records, clinical trials, and more. This
 1064 helps construct a more comprehensive and accurate medical knowledge graph and
 1065 enables researchers to quickly sift out useful information from a large amount of
 1066 complex data [134]. In addition, these LLMs also have strong information integration
 1067 capabilities and can achieve cross-modal fusion between different data sources. Patel
 1068 et al. [135] found through research that big models can quickly generate standardized
 1069 discharge summaries, reducing the work burden on doctors and saving more time for
 1070 doctors to take care of patients.

1071 5.3.2 Clinical decision support

1072 LLM-based intelligent QA systems can automatically extract patients' symptom
 1073 information, analyze and diagnose it according to TCM theory, and provide reference
 1074 opinions for doctors. In terms of case analysis, through the analysis of text data such as
 1075 medical records and symptom descriptions, large models can help doctors quickly
 1076 understand patients' conditions and provide preliminary diagnostic suggestions. In
 1077 terms of drug recommendation, LLMs can automatically match suitable drugs
 1078 according to patients' conditions and symptoms and put forward reasonable medication
 1079 suggestions. In terms of disease prediction, through the analysis of a large number of
 1080 case data, LLMs can predict the occurrence of certain diseases [136].

1081 5.3.3 Providing diagnostic and treatment recommendations

1082 In the field of TCM, LLMs have been widely applied in fields such as diagnosis
 1083 and treatment, medicine, and health preservation due to their excellent performance
 1084 [137]. These models have unique advantages. For example, they can be perfectly
 1085 combined with "the four diagnostic methods", achieve a perfect match between TCM

1086 natural language and SSL, and be adaptable to the characteristics of TCM compound
1087 prescriptions. In addition, these models can help TCM professionals in diagnosis and
1088 treatment research and assist in TCM diagnosis and treatment [9]. Treating diseases
1089 according to their characteristics and providing precise treatment requires doctors to
1090 provide personalized diagnosis and treatment plans according to the conditions of
1091 different patients. LLMs can analyze patients' personal health information and other
1092 relevant data, striving for "one prescription for one person", "one prescription for one
1093 symptom", and "one strategy for one person", thereby improving the patient experience.

1094 We selected a question (Fig. 6), and the differences between intelligence TCM QA
1095 systems based on LLMs and general LLMs were compared. Six TCM-LLMs exhibit
1096 distinct characteristics. Shuzhiqihuang lacks detailed guidance on drug dosage.
1097 TCMChat adopts syndrome differentiation based on classical theories, but
1098 demonstrates limited applicability to contemporary diseases. Shuzhibencao offers
1099 flexibility and diversity in treatment approaches. Congbaosuwen primarily emphasizes
1100 dietary and lifestyle regulation, but lacks complete treatment protocols. Huatuo
1101 provides an integrated treatment plan, highlighting external therapies and requiring
1102 active patient participation. CMLM-Zhongjing provides a user-friendly and efficient
1103 solution without the necessity of local deployment. In the horizontal comparison with
1104 general LLMs, TCM-LLMs excels in providing a comprehensive and traditional
1105 diagnosis by emphasizing patterns of disharmony in the body and offering well-
1106 established herbal remedies rooted in ancient wisdom. Its response is more in-depth and
1107 patient-centric, deeply tied to the TCM approach. The general LLM focuses more on
1108 immediate symptom relief and general medical advice, which can be very useful for
1109 users seeking quick, practical solutions. However, it lacks the deeper, more nuanced
1110 diagnosis that a TCM-specific model would offer. It also offers more conventional
1111 treatments that align with modern medical practices. Therefore, based on a
1112 comprehensive analysis of the characteristics and applicability of these TCM-LLMs, it
1113 is recommended that patients flexibly apply these models during the treatment process,

1114 combining their personal symptoms and physical constitutions in a comprehensive and
 1115 complementary manner.

1116 < Fig. 6. Comparison of intelligence traditional Chinese medicine (TCM)
 1117 question-answering (QA) systems with general large language models (LLMs). >

1118 5.3.4 Equitable distribution of medical resources

1119 At present, China faces severe challenges in uneven distribution of medical
 1120 resources, shortage of primary care doctors, and prevention and treatment of chronic
 1121 diseases. Especially in remote areas, the problems of scarce medical resources and
 1122 limited diagnosis and treatment levels are more prominent, directly affecting the health
 1123 security and quality of life of the people [138]. The development of LLMs in the field
 1124 of TCM provides new ideas and technical support for solving these problems. By
 1125 integrating abundant knowledge of TCM and modern medical research results, LLMs
 1126 can realize the sharing of medical resources nationwide, provide real-time diagnostic
 1127 assistance, treatment suggestions, and chronic disease management plans for primary
 1128 care doctors. The development of LLMs in the field of TCM can promote the sharing
 1129 of medical resources nationwide, improve the medical diagnosis level in areas with
 1130 limited medical resources, and provide innovative solutions for China's primary health
 1131 care services [139].

1132 **6. Discussion**

1133 Currently, TCM intelligent QA systems based on LLMs have been widely applied.
 1134 However, these systems continue to face numerous challenges in handling upstream
 1135 tasks. Firstly, in actual diagnostic consultations, the massive clinical data and patient
 1136 health information recorded by hospitals are sensitive and must be collected, stored,
 1137 and used in strict compliance with regulations. To prevent patient information leaks,
 1138 advanced technical measures must be adopted alongside strict adherence to privacy
 1139 regulations. Federated Learning, an emerging privacy-preserving framework where
 1140 "data remains local and only model parameters are uploaded", has garnered widespread
 1141 attention [140]. This mechanism allows medical institutions to participate in

1142 collaborative model training without moving their local data from its original storage
 1143 environment, effectively alleviating conflicts between data sharing limitations and
 1144 privacy risks. Drawing on the principles of Federated Learning, its application to the
 1145 training and optimization of TCM-LLMs may offer a viable solution to the current
 1146 challenges of high data sensitivity and difficult cross-institutional collaboration.

1147 Simultaneously, due to the large scale of current model parameters, it is quite
 1148 difficult for patients to deploy the model on terminal devices. To enhance the
 1149 accessibility and practicality of the TCM-LLMs in practical applications, model
 1150 lightweighting has become a crucial direction. Through technologies such as model
 1151 compression, knowledge distillation, and parameter pruning, the computational
 1152 resources required during the inference phase can be significantly reduced, making it
 1153 feasible to deploy the model on mobile devices and facilitating the provision of
 1154 convenient medical services. In addition, MoE architecture achieves an optimized
 1155 balance between performance and computational efficiency by selectively activating
 1156 some expert sub-modules, providing a promising research direction for the efficient
 1157 deployment and application of TCM-LLM in the future [141].

1158 TCM diagnosis involves inspection, auscultation and olfaction, inquiry, and
 1159 palpation. However, most of the currently constructed TCM-LLMs still mainly rely on
 1160 text data and lack the ability to deeply integrate unstructured and multimodal data such
 1161 as images (such as tongue images, facial complexion), voice, odor descriptions, and
 1162 pulse conditions [142,143]. This holistic approach to data integration can lead to more
 1163 comprehensive and accurate medical insights, enabling TCM-LLMs to provide a more
 1164 complete picture of a patient's health status. Comprehensive data integration not only
 1165 helps to enhance the model's comprehensive perception of patients' health conditions,
 1166 but also significantly improves the application performance of TCM-LLMs in clinical
 1167 auxiliary diagnosis and intelligent QA systems. This enables them to respond more
 1168 accurately to complex medical queries and provide personalized diagnosis and
 1169 treatment recommendations. However, there are still many technical challenges at

1170 present. For example, cross-modal learning involves how to transfer information and
1171 knowledge between different modalities, especially how to design effective multi-
1172 modal alignment methods. Data from different modalities often contain features at
1173 different levels, how to integrate this information through algorithms while maintaining
1174 its independence and diversity is a major difficulty in the development of multi-modal
1175 TCM-LLMs.

1176 Additionally, in the process of drug development, compound-related data exhibit
1177 highly complex and heterogeneous characteristics [144]. Common data types include
1178 SMILES representations of molecules, two-dimensional molecular structure diagrams,
1179 three-dimensional conformational information, medicinal material source literature,
1180 and formula context [145]. These data originate from diverse sources, ranging from
1181 modern laboratory measurements to historical documents or pharmacopoeia extractions,
1182 leading to inconsistent formats, semantic discrepancies, and uneven data quality. For
1183 example, while SMILES data are concise and easy to process, they lack spatial
1184 positional information and struggle to accurately reflect stereochemical properties;
1185 image or structural data, though information-rich, still lack standardized expression and
1186 annotation methods. Effectively integrating multimodal data remains a key technical
1187 bottleneck in constructing unified and high-quality molecular representation models.
1188 Additionally, TCM drug development emphasizes the holistic synergistic effects of
1189 compound formulas, with their mechanisms of action generally following a systemic
1190 regulatory pattern of multi-component-multi-target-multi-pathway. Traditional
1191 modeling paradigms centered on single-drug-single-target-single-action-pathway are
1192 ill-suited to this context, failing to effectively capture synergistic or antagonistic
1193 interactions between components [146]. Therefore, there is an urgent need to develop
1194 large models capable of integrating structured chemical data with TCM compound
1195 formula knowledge to support key processes such as efficacy prediction, mechanism-
1196 of-action modeling, and candidate drug screening.

1197 In future research, QA must exhibit traceability and interpretability in clinical
1198 applications to gain the trust of both healthcare professionals and patients. Future
1199 research should incorporate techniques such as causal reasoning, symbolic logic, and
1200 visual path analysis to make the model's diagnostic and recommendation processes
1201 transparent and controllable, reducing the risk of "black-box" phenomena.

1202 **7. Conclusion**

1203 This work provides a comprehensive review of the development of LLMs and
1204 offers an in-depth analysis of their key components, including model architecture,
1205 training methods, and core development technologies. Building on this foundation, we
1206 systematically analysis existing TCM intelligent QA systems based on LLM,
1207 highlighting their promising application prospects in areas such as medical education,
1208 drug development, and clinical diagnosis and treatment. By reviewing these systems,
1209 the study aims to offer valuable references for future research and development of
1210 LLMs in the TCM domain.

1211

References:

- [1] Y. Yao, J. Duan, K. Xu, et al., A survey on large language model (llm) security and privacy: The good, the bad, and the ugly, *High-Confid. Comput.* 4 (2024), 100211.
- [2] J. Yang, H. Jin, R. Tang, et al., Harnessing the power of llms in practice: A survey on chatgpt and beyond, *ACM. Trans. Knowl. Discov. Data.* 18 (2024) 1–32.
- [3] Y. Tan, Z. Zhang, M. Li, et al., MedChatZH: A tuning LLM for traditional Chinese medicine consultations, *Comput. Biol. Med.* 172 (2024), 108290.
- [4] Y. Zheng, H.Y. Koh, M. Yang, et al., Large language models in drug discovery and development: From disease mechanisms to clinical trials, *arXiv.* 2024. <https://arXiv.org/abs/2409.04481>.
- [5] J. Cosentino, A. Belyaeva, X. Liu, et al., Towards a Personal Health Large Language Model, *arXiv.* 2024. <https://arXiv.org/abs/2406.06474>.
- [6] C. Wu, Z. Lin, W.L. Fang, et al., A medical diagnostic assistant based on LLM, In *China Health Information Processing Conference*, Springer, Berlin, 2024, pp. 135–147.
- [7] X. Meng, X. Yan, K. Zhang, et al., The application of large language models in medicine: A scoping review, *Iscience* 27 (2024), 109713.
- [8] K. Liu, H. Zhang, H. Liu, et al., Research and Practice on The Establishment of Intelligent TCM Dialectical Treatment Platform, *Chinese Journal of Health Informatics and Management* 20 (2023) 333–338.
- [9] Z. Chen, W. Peng, D. Zhang, et al., Application, Challenges, and Prospects of Large Language Model in the Field of Traditional Chinese Medicine, *Med. J. PUMCH.* 16 (2025) 83–89.
- [10] W. Xiao, C. Song, S. Chen, et al., Key technologies and construction strategies of large language models for traditional Chinese medicine, *CHM.* 55 (2024) 5747–5756.
- [11] L. Yang, Z. Wang, K. Yao, et al., Prospective Reflections on Application of Large Language Models in Field of Traditional Chinese Medicine, *Chin. Arch. Tradit. Chin. Med.* 43 (2025) 16–24.

- [12] X. Wang, T. Yang, K. Hu, Research on Personalized Prescription Recommendation of Traditional Chinese Medicine Based on Large Language Pre-Training Model, *Chin. Arch. Tradit. Chin. Med.* 42 (2024) 15–18+264.
- [13] C. Bai, J. Wang, The Application of the Artificial Intelligence Large Language Model in the Field of Traditional Chinese Medicine, *Journal of Xichang University* 38 (2024) 62–69.
- [14] Y. Yin, L. Zhang, Y. Wang, et al., Question answering system based on knowledge graph in traditional Chinese medicine diagnosis and treatment of viral hepatitis B, *Biomed. Res. Int.* (2022), 7139904.
- [15] M.U. Hadi, A.-T. Qasem, R. Qureshi, et al., Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects, *Authorea Preprints* 1 (2023) 1–26.
- [16] Y. Dai, X. Shao, J. Zhang, et al., TCMChat: A generative large language model for traditional Chinese medicine, *Pharmacol. Res.* 210 (2024), 107530.
- [17] S. Wei, X. Peng, Y.-F. Wang, et al., BianCang: A Traditional Chinese Medicine Large Language Model, *arXiv*. 2024. <https://arXiv.org/abs/2411.11027>.
- [18] Q. Zhang, T. Gui, R. Zheng, et al., Large-scale language models: from theory to practice, 2023, pp. 5–6.
- [19] W.X. Zhao, K. Zhou, J. Li, et al., A survey of large language models, *arXiv*. 2023. <https://arXiv.org/abs/2303.18223>.
- [20] B. Wang, A. Wang, F.X. Chen, et al., Evaluating word embedding models: Methods and experimental results, *APSIPA trans. signal.* 8 (2019), e19.
- [21] T. Mikolov, K. Chen, G. Corrado, et al., Efficient estimation of word representations in vector space, *arXiv*. 2013. <https://arXiv.org/abs/1301.3781>.
- [22] J. Pennington, R. Socher, C.D. Manning, Glove: Global vectors for word representation, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543.
- [23] A. Vaswani, N. Shazeer, N. Parmar, et al. Attention is all you need, *arXiv*. 2017,

<https://arXiv.org/abs/1706.03762>.

- [24] Y. Wang, Q. Li, Z. Dai, et al., Current status and trends in large language modeling research, Chin. J. Eng. 46 (2024) 1411–1425.
- [25] J. Devlin, M.-W. Chang, K. Lee, et al., Bert: Pre-training of deep bidirectional transformers for language understanding, Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, vol 1 (long and short papers), 2019, pp. 4171–4186.
- [26] J. Howard, S. Ruder, Universal language model fine-tuning for text classification, arXiv. 2018. <https://arXiv.org/abs/1801.06146>.
- [27] A. Radford, J. Wu, R. Child, et al., Language models are unsupervised multitask learners, OpenAI blog 1 (2019), 9.
- [28] T.B. Brown, B. Mann, N. Ryder, et al., Language models are few-shot learners, arXiv. 2020. <https://arXiv.org/abs/2005.14165>.
- [29] Y.Q. Shen, L. Heacock, J. Elias, et al., ChatGPT and other large language models are double-edged swords, Radiology. 307 (2023), e230163.
- [30] J.T. Halloran, M.S. Gulati, P.F. Roysdon, Mamba state-space models can be strong downstream learners, arXiv. 2024. <https://arXiv.org/abs/2406.00209>.
- [31] J.W. Zuo, M. Velikanov, D.E. Rhaiem, et al., Falcon mamba: The first competitive attention-free 7b language model, arXiv. 2024. <https://arXiv.org/abs/2410.05355>.
- [32] B. Peng, E. Alcaide, Q. Anthony, et al., Rwkv: Reinventing rnns for the transformer era, arXiv. 2023. <https://arXiv.org/abs/2305.13048>.
- [33] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, arXiv. 2014. <https://arXiv.org/abs/1409.0473>.
- [34] J. Liu, K.H. Lim, R.K.-W. Lee, et al., Towards Objective and Unbiased Decision Assessments with LLM-Enhanced Hierarchical Attention Networks, arXiv. 2024. <https://arXiv.org/abs/2411.08504>.
- [35] L. Yi, X. Zhou, W. He, et al., LongHeads: Multi-Head Attention is Secretly a Long Context Processor, arXiv. 2024. <https://arXiv.org/abs/2402.10685>.

- [36] G. Xiao., J. Tang, J. Zuo, et al., DuoAttention: Efficient Long-Context LLM Inference with Retrieval and Streaming Heads, arXiv. 2024. <https://arXiv.org/abs/2410.10819>.
- [37] Z. Liu, Q. Song, Q.C. Xiao, et al., FFSplit: Split Feed-Forward Network For Optimizing Accuracy-Efficiency Trade-off in Language Model Inference, arXiv. 2024. <https://arXiv.org/abs/2401.04044>.
- [38] T. Lin, Y. Wang, X. Liu, et al., A survey of transformers, AI open 3 (2022) 111–132.
- [39] Y. Dong, J.-B. Cordonnier, A. Loukas, Attention is not all you need: Pure attention loses rank doubly exponentially with depth, International Conference on Machine Learning, PMLR. (2021) 2793–2803.
- [40] A. Haviv, O. Ram, O. Press, et al., Transformer language models without positional encodings still learn positional information, arXiv. 2022. <https://arXiv.org/abs/2203.16634>.
- [41] S. Sukhbaatar, A. Szlam, J. Weston, et al., End-to-end memory networks, arXiv. 2015. <https://arXiv.org/abs/1503.08895>.
- [42] P. Shaw, J. Uszkoreit, A. Vaswani, Self-attention with relative position representations, arXiv. 2018. <https://arXiv.org/abs/1803.02155>.
- [43] N.-Q. Pham, T.-L. Ha, T.-N. Nguyen, et al., Relative Positional Encoding for Speech Recognition and Direct Translation, arXiv. 2020. <https://arXiv.org/abs/1803.02155>.
- [44] K. Wu, H. Peng, M. Chen, et al., Rethinking and Improving Relative Position Encoding for Vision Transformer, Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 10033–10041.
- [45] O. Press, N.A. Smith, M. Lewis, et al., Train Short, Test Long: Attention with linear biases enables input length extrapolation, arXiv. 2021. <https://arXiv.org/abs/2108.12409>.

- [46] J. Su, Y. Lu, S. Pan, et al., Roformer: Enhanced transformer with rotary position embedding, *Neurocomputing* 568 (2024), 127063.
- [47] Z. Dai, Z. Yang, Y. Yang, et al., Transformer-xl: Attentive language models beyond a fixed-length context, *arXiv*. 2019. <https://arXiv.org/abs/1901.02860>.
- [48] P. He, X. Liu, J. Gao, et al., Deberta: Decoding-enhanced bert with disentangled attention, *arXiv*. 2020. <https://arXiv.org/abs/2006.03654>.
- [49] C. Raffel, N. Shazeer, A. Roberts, et al., Exploring the limits of transfer learning with a unified text-to-text transformer, *J. Mach. Learn. Res.* 21 (2020) 1–67.
- [50] H. Zhou, F. Liu, B. Gu, et al., A survey of large language models in medicine: Progress, application, and challenge, *arXiv*. 2023. <https://arXiv.org/abs/2311.05112>.
- [51] T.B. Brown, B. Mann, N. Ryder, et al., Language Models are Few-Shot Learners, *NeurIPS*. 33 (2020) 1877–1901.
- [52] A. Chowdhery, S. Narang, J. Devlin, et al., Palm: Scaling language modeling with pathways, *J. Mach. Learn. Res.* 24 (2023) 1–113.
- [53] C. Zhang, S. Bengio, M. Hardt, et al., Understanding deep learning (still) requires rethinking generalization, *Commun. Acm.* 64 (2021) 107–115.
- [54] N. Mehrabi, F. Morstatter, N. Saxena, et al., A survey on bias and fairness in machine learning, *Acm. Comput. Surv.* 54 (2021) 1–35.
- [55] T. Li, S. Shetty, A. Kamath, et al., CancerGPT for few shot drug pair synergy prediction using large pretrained language models, *NPJ. Digit. Med.* 7 (2024), 40.
- [56] A. Wettig, T. Gao, Z. Zhong, et al., Should you mask 15% in masked language modeling? *arXiv*. 2022. <https://arXiv.org/abs/2202.08005>.
- [57] Z. Du, Y. Qian, X. Liu, et al., Glm: General language model pretraining with autoregressive blank infilling, *arXiv*. 2021. <https://arXiv.org/abs/2103.10360>.
- [58] H. An, Y. Chen, Z. Sun, et al., SentenceVAE: Enable next-sentence prediction for large language models with faster speed, higher accuracy and longer context, *arXiv*. 2024. <https://arXiv.org/abs/2408.00655>.

- [59] A. Bindal, S. Ramanujam, D. Golland, et al., Improved Content Understanding With Effective Use of Multi-task Contrastive Learning, arXiv. 2024. <https://arXiv.org/abs/2405.11344>.
- [60] Y. Ling, X. Jiang, Y. Kim, MALLM-GAN: Multi-Agent Large Language Model as Generative Adversarial Network for Synthesizing Tabular Data, arXiv. 2024. <https://arXiv.org/abs/2406.10521>.
- [61] S.A. Sahoo, Meta-Learning for Large Language Models: Teaching LLMs to Learn New Tasks with Minimal Data, Available at SSRN 4977093. 2024.
- [62] H. Naveed, A.U. Khan, S. Qiu, et al., A comprehensive overview of large language models, arXiv. 2023. <https://arXiv.org/abs/2307.06435>.
- [63] J. Devlin, M.-W. Chang, K. Lee, et al., Bert: Pre-training of deep bidirectional transformers for language understanding, arXiv. 2019. <https://doi.org/10.18653/v1/N19-1423>.
- [64] Y. Liu, M. Ott, N. Goyal, et al., Roberta: A robustly optimized bert pretraining approach, arXiv. 2019. <https://arXiv.org/abs/1907.11692>.
- [65] K. Clark, M.-T. Luong, Q.V. Le, et al., ELECTRA: Pre training text encoders as discriminators rather than generators, in Proc. 8th ICLR. Addis Ababa, Ethiopia, 2020, pp. 118.
- [66] H. Wang, J. Li, H. Wu, et al., Pre-trained language models and their applications, Engineering 25 (2023), 51–65.
- [67] A. Liu, B. Feng, B. Xue, et al., Deepseek-v3 technical report, arXiv. 2024. <https://arXiv.org/abs/2412.19437>.
- [68] Vicuna, An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. <https://lmsys.org/blog/2023-03-30-vicuna>. (Accessed 20 December 2024).
- [69] R. Taori, I. Gulrajani, T. Zhang, et al., Stanford alpaca: An instruction-following llama model, 2023. https://github.com/tatsu-lab/stanford_alpaca.
- [70] A.Q. Jiang, A. Sablayrolles, A. Mensch, et al., Mistral 7B, arXiv. 2023. <https://arXiv.org/abs/2310.06825>.

- [71] H. Touvron, T. Lavril, G. Izacard, et al., Llama: Open and efficient foundation language models, arXiv.2023. <https://arXiv.org/abs/2302.13971>.
- [72] H. Touvron, L. Martin, K. Stone, et al., Llama 2: Open foundation and fine-tuned chat models, arXiv. 2023. <https://arXiv.org/abs/2307.09288>.
- [73] B. Cohen-Wang, H. Shah, K. Georgiev, et al., Contextcite: Attributing model generation to context, NeurIPS. 37 (2024) 95764–95807.
- [74] J. Bai, S. Bai, Y. Chu, et al., Qwen technical report, arXiv. 2023. <https://arXiv.org/abs/2309.16609>.
- [75] H.W. Chung, L. Hou, S. Longpre, et al., Scaling instruction-finetuned language models, J. Mach. Leaen. Res. 25 (2024) 1–53.
- [76] A. Mihalache, J. Grad, N.S. Patil, et al., Google Gemini and Bard artificial intelligence chatbot performance in ophthalmology knowledge assessment, Eye. 38 (2024) 2530–2535.
- [77] L. Ouyang, J. Wu, X. Jiang, et al., Training language models to follow instructions with human feedback, NeurIPS. 35 (2022) 27730–27744.
- [78] J. Achiam, S. Adler, S. Agarwal, et al., Gpt-4 technical report, arXiv. 2023. <https://arXiv.org/abs/2303.08774>.
- [79] R. Kurokawa, Y. Ohizumi, J. Kanzawa, et al., Diagnostic performances of Claude 3 Opus and Claude 3.5 Sonnet from patient history and key images in Radiology’s “Diagnosis Please” cases, Jpn. J. Radiol. 42 (2024) 1399–1402.
- [80] M. Lewis, Y. Liu, N. Goyal, et al., Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, arXiv. 2019. <https://arXiv.org/abs/1910.13461>.
- [81] A. Zeng, X. Liu, Z. Du, et al., Glm-130b: An open bilingual pre-trained model, arXiv. 2022. <https://arXiv.org/abs/2210.02414>.
- [82] P.M. Abhinav, R.M. SujayKumar, O. Christopher, Machine Translation with Large Language Models: Decoder Only vs. Encoder-Decoder. arXiv. 2024. <https://arXiv.org/abs/2409.14747>.

- [83] Y. Tay, M. Dehghani, V.Q. Tran, et al., UI2: Unifying language learning paradigms, arXiv. 2022. <https://arXiv.org/abs/2205.05131>.
- [84] L. Wu, Z. Zheng, Z. Qiu, et al., A survey on large language models for recommendation, *World Wide Web.* 27 (2024), 60.
- [85] J. Wei, M. Bosma, V.Y. Zhao, et al., Finetuned language models are zero-shot learners, arXiv. 2021. <https://arXiv.org/abs/2109.01652>.
- [86] V. Sanh, A. Webson, C. Raffel, et al., Multitask prompted training enables zero-shot task generalization, arXiv. 2021. <https://arXiv.org/abs/2110.08207>.
- [87] S. Iyer, X.V. Lin, R. Pasunuru, et al., Opt-iml: Scaling language model instruction meta learning through the lens of generalization, arXiv. 2022. <https://arXiv.org/abs/2212.12017>.
- [88] X.L. Li, P. Liang, Prefix-tuning: Optimizing continuous prompts for generation, arXiv. 2021. <https://arXiv.org/abs/2101.00190>.
- [89] D. Yin, L. Hu, B. Li, et al., Adapter is all you need for tuning visual tasks. arXiv. 2023. <https://arXiv.org/abs/2311.15010>.
- [90] N. Houlsby, A. Giurgiu, S. Jastrzebski, et al., Parameter-efficient transfer learning for NLP, *Proceedings of the 36 th International Conference on Machine Learning*, Long Beach, California, 2019, pp. 2790–2799.
- [91] E. Hu, Y. Shen, P. Wallis, et al., Lora: Low-rank adaptation of large language models, arXiv. 2021. <https://arXiv.org/abs/2106.09685>.
- [92] M. Nikdan, S. Tabesh, E. Crnčević, et al., Rosa: Accurate parameter-efficient fine-tuning via robust adaptation, arXiv. 2024. <https://arXiv.org/abs/2401.04679>.
- [93] J. Rasley, S. Rajbhandari, O. Ruwase, et al., DeepSpeed: System optimizations enable training deep learning models with over 100 billion parameters, *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, 2020, pp. 3505–3506.

- [94] M. Shoeybi, M. Patwary, R. Puri, et al., Megatron-lm: Training multi-billion parameter language models using model parallelism, arXiv. 2019. <https://arXiv.org/abs/1909.08053>.
- [95] D. Narayanan, M. Shoeybi, J. Casper, et al., Efficient large-scale language model training on gpu clusters using megatron-lm, arXiv. 2021. <https://arXiv.org/abs/zenodo.5181820>.
- [96] V. Korthikanti, J. Casper, S. Lym, et al., Reducing activation recomputation in large transformer models, Proceedings of Machine Learning and Systems, 5 (2023) 341–353.
- [97] D. Gao, L. Ji, L. Zhou, et al., Assistgpt: A general multi-modal assistant that can plan, execute, inspect, and learn, arXiv. 2023. <https://arXiv.org/abs/2306.08640>.
- [98] P. Lu, B. Peng, H. Cheng, et al., Chameleon: Plug-and-play compositional reasoning with large language models, arXiv. 2023. <https://arXiv.org/abs/2304.09842>.
- [99] B. Paranjape, S. Lundberg, S. Singh, et al., Art: Automatic multi-step reasoning and tool-use for large language models, arXiv. 2023. <https://arXiv.org/abs/2303.09014>.
- [100] J. Kim, M. Min, From rag to qa-rag: Integrating generative ai for pharmaceutical regulatory compliance process, arXiv. 2024. <https://arXiv.org/abs/2402.01717>.
- [101] C. Zakka, A. Chaurasia, R. Shad, et al., Almanac-retrieval-augmented language models for clinical medicine, NEJM. AI. 1 (2024), AIoa2300068.
- [102] D. Ferber, I.C. Wiest, G. Wölflein, et al., GPT-4 for Information Retrieval and Comparison of Medical Oncology Guidelines, NEJM. AI. 1 (2024), AIcs2300235.
- [103] C. Ma, Z. Wu, J. Wang, et al., An Iterative Optimizing Framework for Radiology Report Summarization with ChatGPT, IEEE Trans. Artif. Intell. 8 (2024) 4163–4175.
- [104] Y. Guo, W. Qiu, G. Leroy, et al., Retrieval augmentation of large language models for lay language generation, J. Biomed. Inform. 149 (2024), 104580.
- [105] C. Wang, M. Li, J. He, et al., A survey for large language models in biomedicine, arXiv. 2024. <https://arXiv.org/abs/2409.00133>.

- [106] H. Wang, C. Liu, N. Xi, et al., Huatuo: Tuning llama model with chinese medical knowledge, arXiv. 2023. <https://arXiv.org/abs/2304.06975>.
- [107] S. Yang, H. Zhao, S. Zhu, et al., Zhongjing: Enhancing the chinese medical capabilities of large language model through expert feedback and real-world multi-turn dialogue, Proceedings of the AAAI conference on artificial intelligence 38 (2024) 19368–19376.
- [108] Y. Kang, Y. Chang, J. Fu, et al., CMLM-ZhongJing: Large language model is good story listener, GitHub Repository 2023. <https://github.com/pariskang/CMLM-ZhongJing>. (Accessed 20 May 2025).
- [109] Y. Chen, Z. Wang, H. Zheng, et al., Bianque: Balancing the questioning and suggestion ability of health llms with multi-turn health conversations polished by chatgpt, arXiv. 2023. <https://arXiv.org/abs/2310.15896>.
- [110] Chinese Pharmacopoeia Commission. People's Republic of China, 2020 edition. Beijing: The Medicine Science and Tech nology Press of China, 2020.
- [111] H. Tian, K. Yang, X. Dong, et al. TCMLLM-PR: evaluation of large language models for prescription recommendation in traditional Chinese medicine, Digital Chinese Medicine 7 (2024) 343–355
- [112] H. Zhang, X. Wang, Z. Meng, et al., Qibo: A Large Language Model for Traditional Chinese Medicine, arXiv. 2024. <https://arXiv.org/abs/2403.16056>.
- [113] R. Hua, X. Dong, Y. Wei, et al., Lingdan: enhancing encoding of traditional Chinese medicine knowledge for clinical reasoning tasks with large language models[J]. J. Am. Med. Inform. Assn. 2024, 31(9) 2019–2029.
- [114] Y. Liu, S. Luo, Z. Zhong, et al., Hengqin-RA-v1: Advanced Large Language Model for Diagnosis and Treatment of Rheumatoid Arthritis with Dataset based Traditional Chinese Medicine, arXiv. 2025. <https://arXiv.org/abs/2501.02471>.
- [115] W. Zhu, W. Yue, X. Wang, ShenNong-TCM: A Traditional Chinese Medicine Large Language Model, GitHub 2023. <https://github.com/michael-wzhu/ShenNong-TCM-LLM>. (Accessed 20 May 2025).

- [116] C. Liu, K. Sun, Q. Zhou, et al., CPMI-ChatGLM: parameter-efficient fine-tuning ChatGLM with Chinese patent medicine instructions, *Sci. rep.* 14 (2024), 6403.
- [117] Y. Liao, S. Jiang, Y. Wang, et al. MING-MOE: Enhancing medical multi-task learning in large language models with sparse mixture of low-rank adapter experts, *arXiv*. 2024. <https://arXiv.org/abs/2404.09027>.
- [118] R. Wang, Y. Duan, C. Lam, et al. Ivygpt: Interactive chinese pathway language model in medical domain, *CAAI International Conference on Artificial Intelligence*. Singapore: Springer Nature, Singapore, 2023, pp. 378–382.
- [119] L. Gao, C.-H. Jia, W. Wang, Recent advances in the study of ancient books on traditional Chinese medicine, *World J. Tradit. Chin. Med.* 6 (2020) 61–66.
- [120] I. Abdelaziz, A. Fokoue, O. Hassanzadeh, et al., Large-scale structural and textual similarity-based mining of knowledge graph to predict drug–drug interactions, *J. Web. Semant.* 44 (2017) 104–117.
- [121] J. Zhang, S. Yang, J. Liu, et al., AIGC Empowering the Revitalization of Ancient Books on Traditional Chinese Medicine: Building the Huang-Di Large Language Model, *Library Tribune* 44 (2024) 103–112.
- [122] W. Dai, J. Lin, F. Jin, et al., Can large language models provide feedback to students? A case study on ChatGPT, *ICALT*. IEEE. (2023) 323–325.
- [123] J.C. Young, M. Shishido, Investigating OpenAI’s ChatGPT potentials in generating Chatbot’s dialogue for English as a foreign language learning, *Int. J. Adv. Comput. Sci. Appl.* 14 (2023) 65–72.
- [124] S. Kim, P.A. Thiessen, E.E. Bolton, et al., PubChem substance and compound databases, *Nucleic. Acids. Res.* 44 (2016) D1202–D1213.
- [125] S.R. Lynch, T. Bothwell, L. Campbell, A comparison of physical properties, screening procedures and a human efficacy trial for predicting the bioavailability of commercial elemental iron powders used for food fortification, *Int. J. Vitam. Nutr. Res.* 77 (2007) 107–124.

- [126] T. Andrysek, Impact of physical properties of formulations on bioavailability of active substance: current and novel drugs with cyclosporine, *Mol. Immunol.* 39 (2003) 1061–1065.
- [127] H.C.S. Chan, H. Shan, T. Dahoun, et al., Advancing drug discovery via artificial intelligence, *Trends. Pharmacol. Sci.* 40 (2019) 592–604.
- [128] F. Huang, H. Yang, X. Zhu, Progress in the Application of Artificial Intelligence in New Drug Discovery, *Progress in Pharmaceutical Sciences* 45 (2021) 502-511.
- [129] Z. Liu, R.A. Roberts, M. Lal-Nag, et al., AI-based language models powering drug discovery and development, *Drug Discov. Today* 26 (2021) 2593–2607.
- [130] L. Liang, C. Deng, Y. Zhang, et al., Application and Challenges of Artificial Intelligence in Drug Discovery, *Progress in Pharmaceutical Sciences* 44 (2020) 18–27.
- [131] L. Patel, T. Shukla, X. Huang, et al., Machine learning methods in drug discovery, *Molecules*. 25 (2020), 5277.
- [132] J.M. Stokes, K. Yang, K. Swanson, et al. A deep learning approach to antibiotic discovery, *Cell.* 180 (2020) 688–702.
- [133] T. Wu, R. Lin, P. Cui, et al. Deep learning-based drug screening for the discovery of potential therapeutic agents for Alzheimer's disease, *J. Pharm. Anal.* 14 (2024), 101022.
- [134] X. Yang, A.K. Chen, N. PourNejatian, et al., A large language model for electronic health records, *Npj. Digit. Med.* 5 (2022), 194.
- [135] S.B. Patel, K. Lam, ChatGPT: the future of discharge summaries? *Lancet Digit. Health* 5 (2023) e107–e108.
- [136] B. Fatani, ChatGPT for future medical and dental research, *Cureus*. 15 (2023), e37285.
- [137] W. Ma, M. Meng, H. Dai, et al., A Comprehensive Review of the Applications of Large Language Models in Clinical Medicine with ChatGPT as a Representative, *Journal of Medical Intelligence* 44 (2023) 9–17.

- [138] R. Khera, A.J. Butte, M. Berkwits, et al., AI in medicine—JAMA’s focus on clinical outcomes, patient-centered care, quality, and equity, *Jama-J. Am Med. Assoc.* 330 (2023) 818–820.
- [139] W. Yan, J. Hu, H. Ceng, et al., The Application of Large Language Models in Primary Healthcare Services and the Challenges, *Chinese General Practice* 28 (2025) 1–6.
- [140] L. Li, Y. Fan, M. Tse, et al. A review of applications in federated learning, *Compu. Ind. Eng.* 149 (2020), 106854.
- [141] X. Du, T. Gunter, X. Kong, et al. Revisiting MoE and Dense Speed-Accuracy Comparisons for LLM Training, *arXiv*. 2024. <https://arXiv.org/abs/2405.15052>.
- [142] Y. Bao, H. Ding, Z. Zhang, et al., Intelligent Acupuncture: Data-driven Revolution of Traditional Chinese Medicine, *Acupuncture and Herbal Medicine* 3 (2023) 271–284.
- [143] X. Liu, T. Gong, Artificial Intelligence and Evidence-Based Research Will Promote the Development of Traditional Medicine, *Acupuncture and Herbal Medicine* 4 (2024) 134–135.
- [144] C. Chakraborty, M. Bhattacharya, S.-S. Lee, Artificial intelligence enabled ChatGPT and large language models in drug target discovery, drug discovery, and development, *Mol. Ther. Nucleic Acids* 33 (2023) 866–868.
- [145] S. Liu, H. Wang, W. Liu, et al. Pre-training molecular graph representation with 3d geometry, *arXiv*. 2021. <https://arXiv.org/abs/2110.07728>
- [146] Y. Li, X. Liu, J. Zhou, et al. Artificial intelligence in traditional Chinese medicine: advances in multi-metabolite multi-target interaction modeling, *Front. Pharmacol.* 16 (2025), 1541509.

Figures captions:

Fig. 1. Summary of the data in the field of traditional Chinese medicine (TCM). TCM data has rich types, mainly covering the following four aspects: experience data, pharmacological data, chemical data, and clinical data. The sources of these data are very extensive, including ancient books, experts, experiments, public databases, and literature. Their types show a high degree of diversity, and the formats also have multimodal characteristics, involving various forms such as text, image, and table.

Fig. 2. Overview analysis diagram of intelligent question-answering (QA) systems based on large language models (LLMs). (A) Co-occurrence of keywords. (B) Emergent map of key words in TCM-LLMs.

Fig. 3. Background of large language models (LLMs) development. (A) Open source LLMs. (B) Closed source LLMs. WE: word embedding; PT: pre-training; FT: fine-tuning; RAG: retrieval-augmented generation; NLP: natural language processing.

Fig. 4. Working principle of large language models (LLMs) and training methods. (A) Diagram of Transformer architecture. (B) Three training methods of LLMs. RAG: retrieval-augmented generation; Norm: normalization.

Fig. 5. Prospects for the application of traditional Chinese medicine (TCM)-large language models (LLMs).

Fig. 6. Comparison of intelligence traditional Chinese medicine (TCM) question-answering (QA) systems with general large language models (LLMs). (A) Intelligence TCM QA systems based on LLMs. (B) General LLMs.

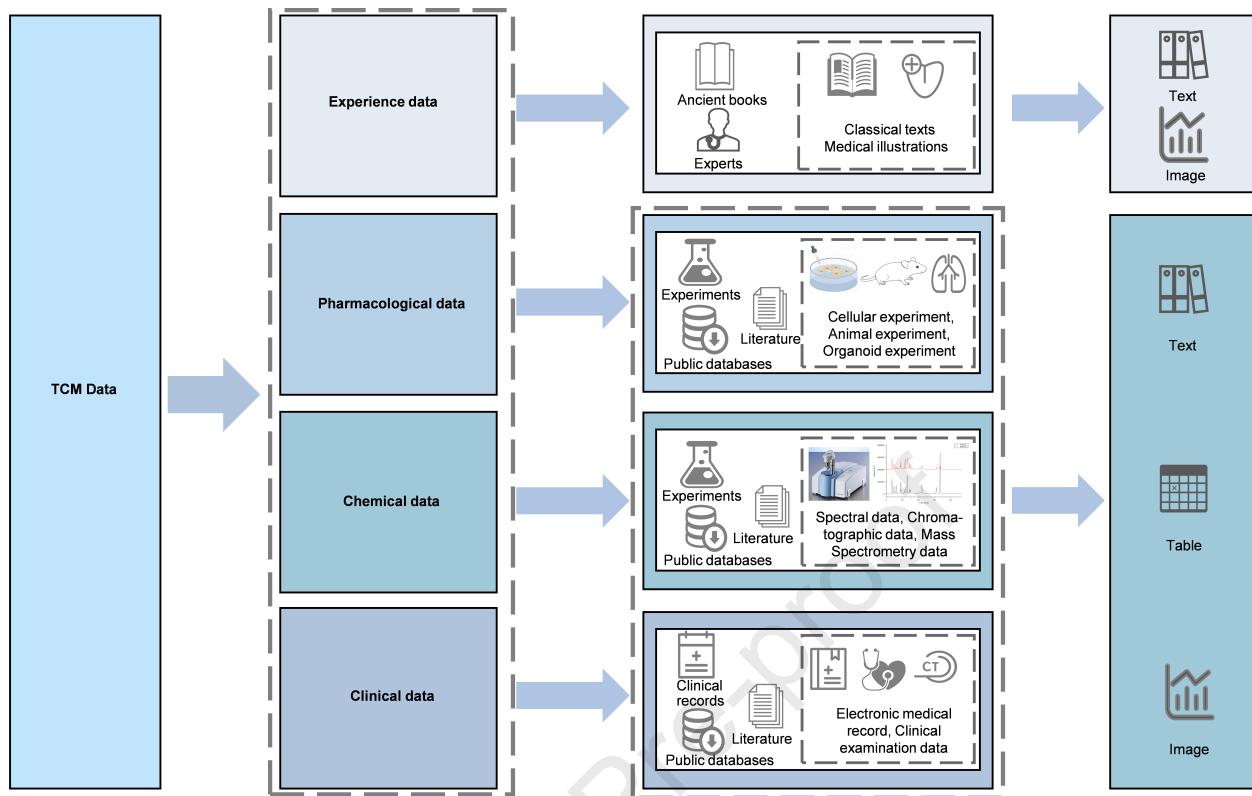
Acknowledgments

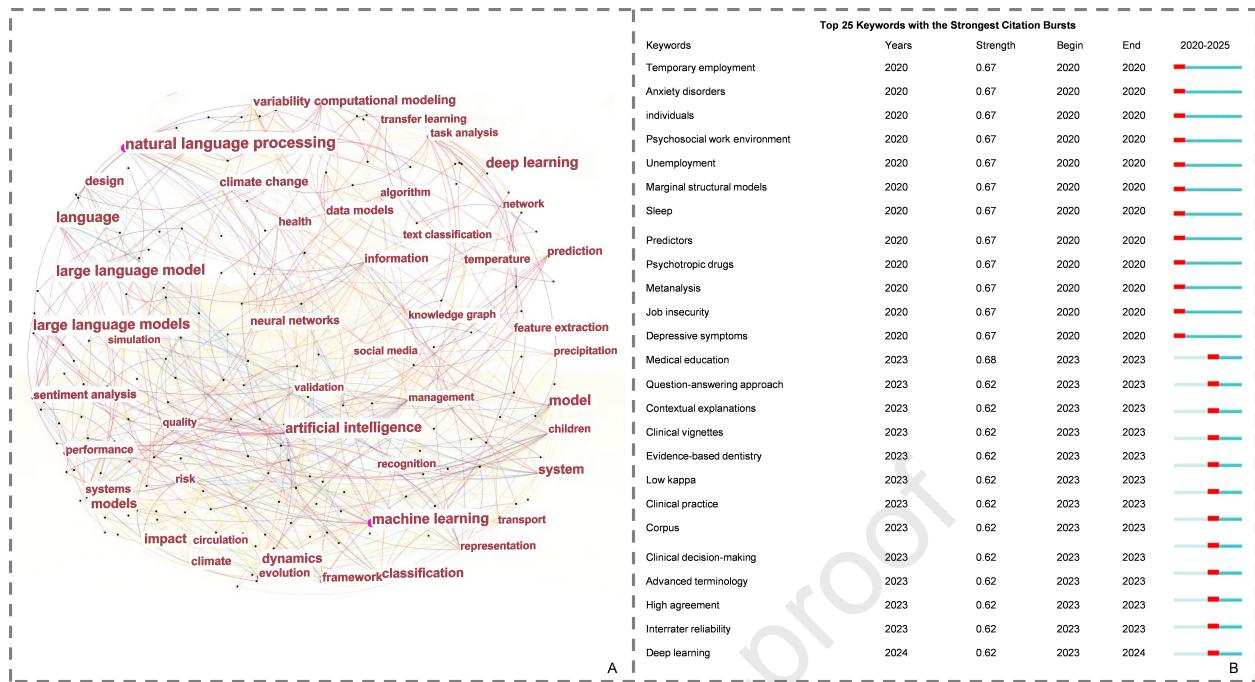
This work was supported by the Special Project for Technological Innovation in New Productive Forces of Modern Chinese Medicines (Grant No.: 24ZXZKSY00010), Science and Technology Program of Tianjin (Grant No.: 24ZXZSSS00460).

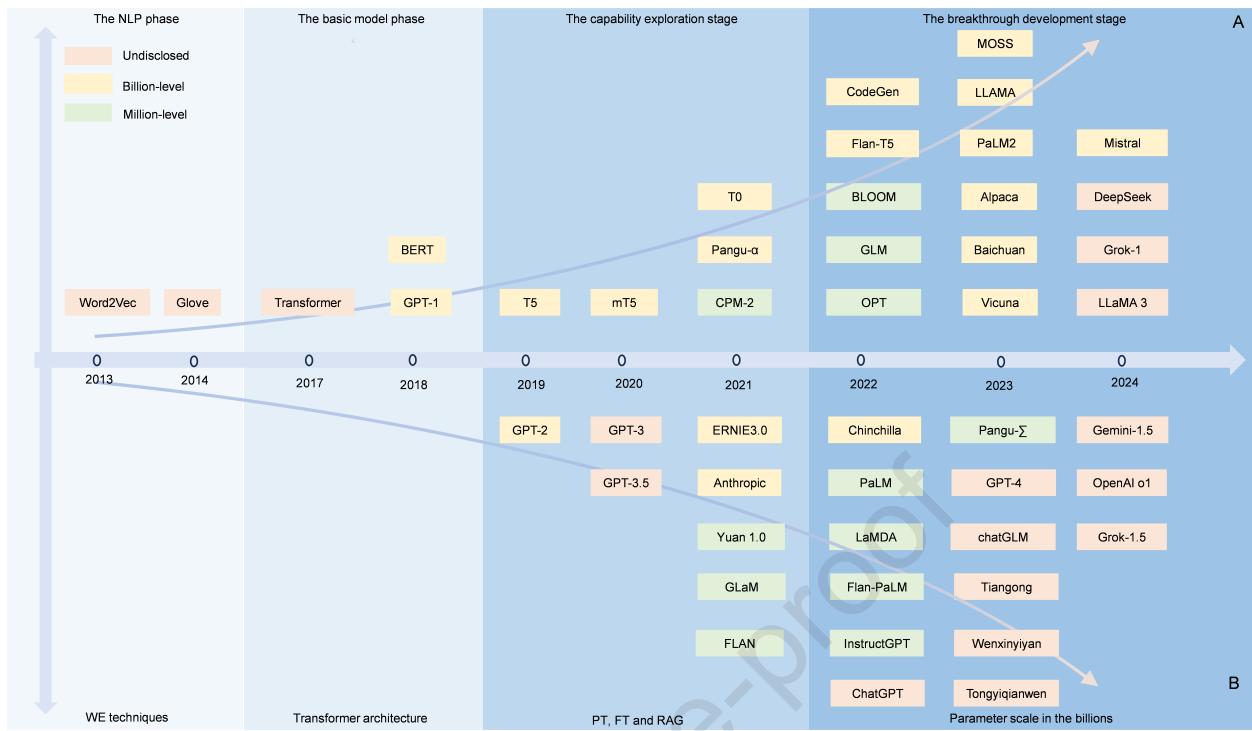
Table 1. Comparison of the model structure and Pre-Training (PT) parameters of general large language models (LLMs).

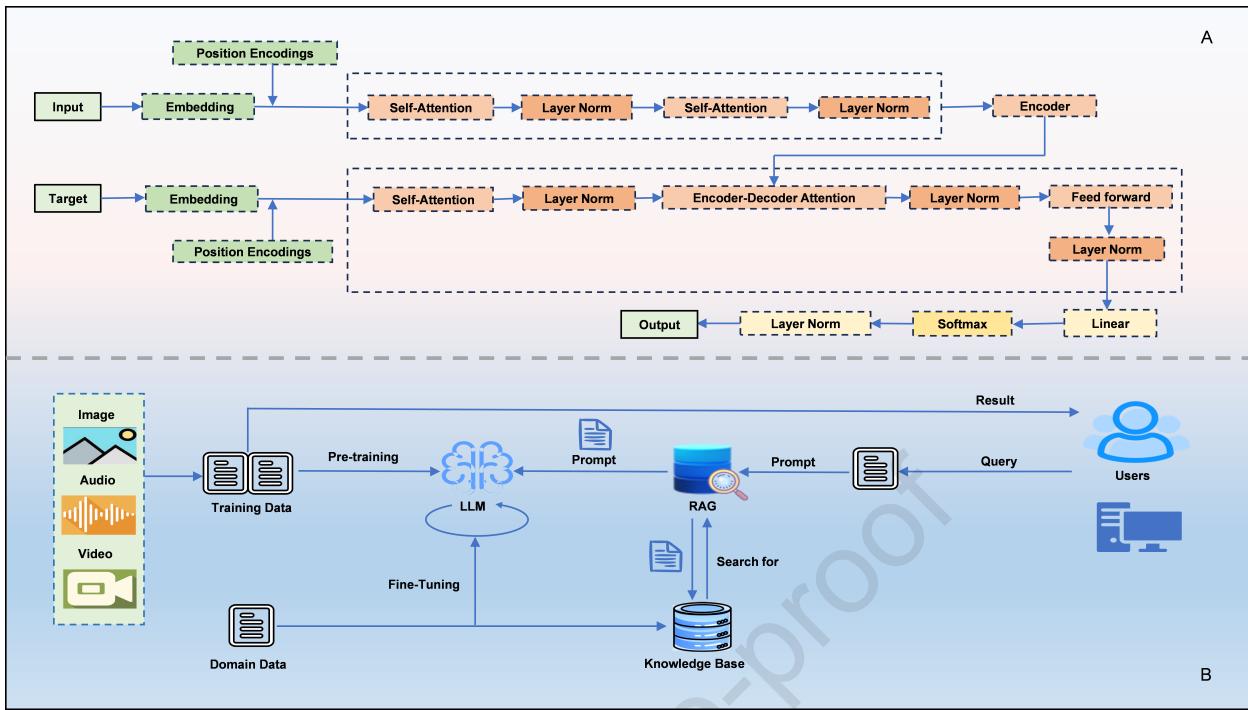
Num	Model name	Model structure	Parameters	Refs.
1	GPT-2	Decoder only	1.5B	[27]
2	DeBERTa	Encoder only	1.5B	[48]
3	T5	Encoder and decoder	11B	[49]
4	PaLM	Decoder only	8B/62B/540B	[52]
5	GPT-3	Decoder only	6.7B/13B/175B	[62]
6	BERT	Encoder only	110M/340M	[63]
7	RoBERTa	Encoder only	355M	[64]
8	ELECTRA	Encoder only	335M	[65]
9	XLNet	Encoder only	360M	[65]
10	CPM	Encoder only	2.6B	[66]
11	CPM-2	Encoder only	11B	[66]
12	GLaM	Encoder only	1.2 trillion	[66]
13	Gopher	Decoder only	280B	[66]
14	DeepSeek-V3	Decoder only	671B	[67]
15	Vicuna	Decoder only	7B/13B	[68]
16	Alpaca	Decoder only	7B/13B	[69]
17	Mistral	Decoder only	7B	[70]
18	LLaMA	Decoder only	7B/13B/33B/65B	[71]
19	LLaMA-2	Decoder only	7B/13B/34B/70B	[72]
20	LLaMA-3	Decoder only	8B/70B	[73]
21	Qwen	Decoder only	1.8B/7B/14B/72B	[74]
22	FLAN-PaLM	Decoder only	540B	[75]
23	Gemini	Decoder only	–	[76]
24	GPT-3.5	Decoder only	–	[77]
25	GPT-4	Decoder only	–	[78]
26	Claude-3	Decoder only	–	[79]
27	BART	Encoder and decoder	140M/400M	[80]
28	GLM	Encoder and decoder	130B	[81]
29	mT5	Encoder and decoder	300M	[82]
30	UL2	Encoder and decoder	19.5B	[83]

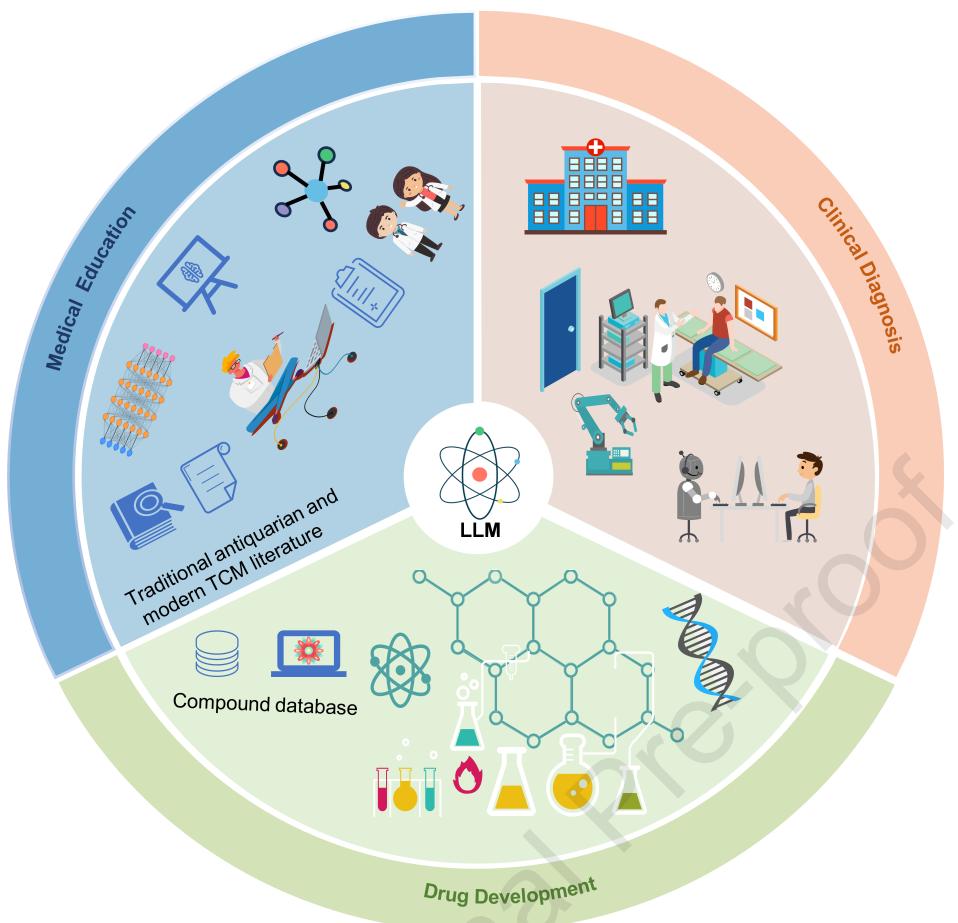
–: no data; M: Million; B: Billion.

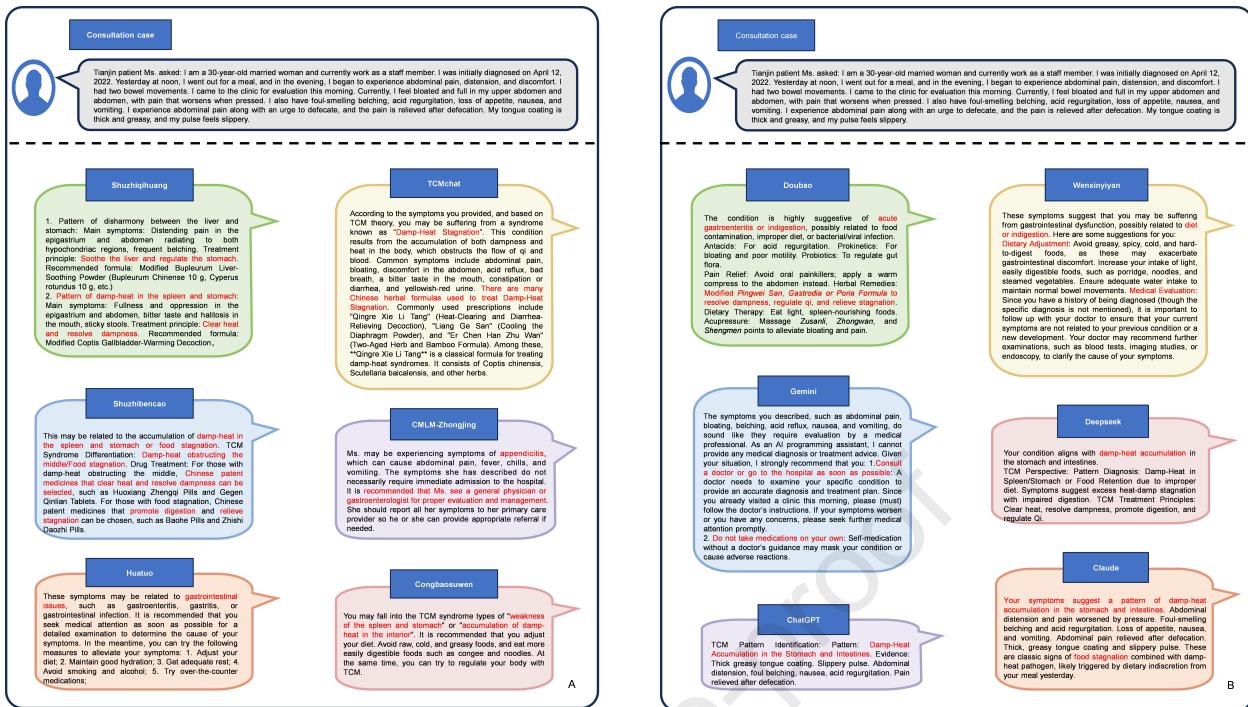












Highlights

- This article is the first systematic review of TCM-LLMs.
- Discussed the characteristics of LLMs in four different stages of development.
- Summarized and compared the working principles and key technologies of LLMs.
- Evaluated the advantages and limitations of open-source and closed-source TCM-LLMs.
- Discussed the prospects and potential impact of TCM-LLM applications.

Declaration of competing interest

The authors declare that there are no conflicts of interest.

CRediT authorship contribution statement

Qilan Xu: Writing – original draft, Methodology, Investigation, Formal analysis.
Tong Wu: Writing – original draft, Methodology, Investigation, Formal analysis.
Yiwen Wang: Methodology, Investigation. **Xingyu Li:** Formal analysis. **Heshui Yu:** Conceptualization. **Shixin Cen:** Writing – review & editing, Visualization, Supervision.
Zheng Li: Conceptualization, Funding acquisition, Writing – review & editing.