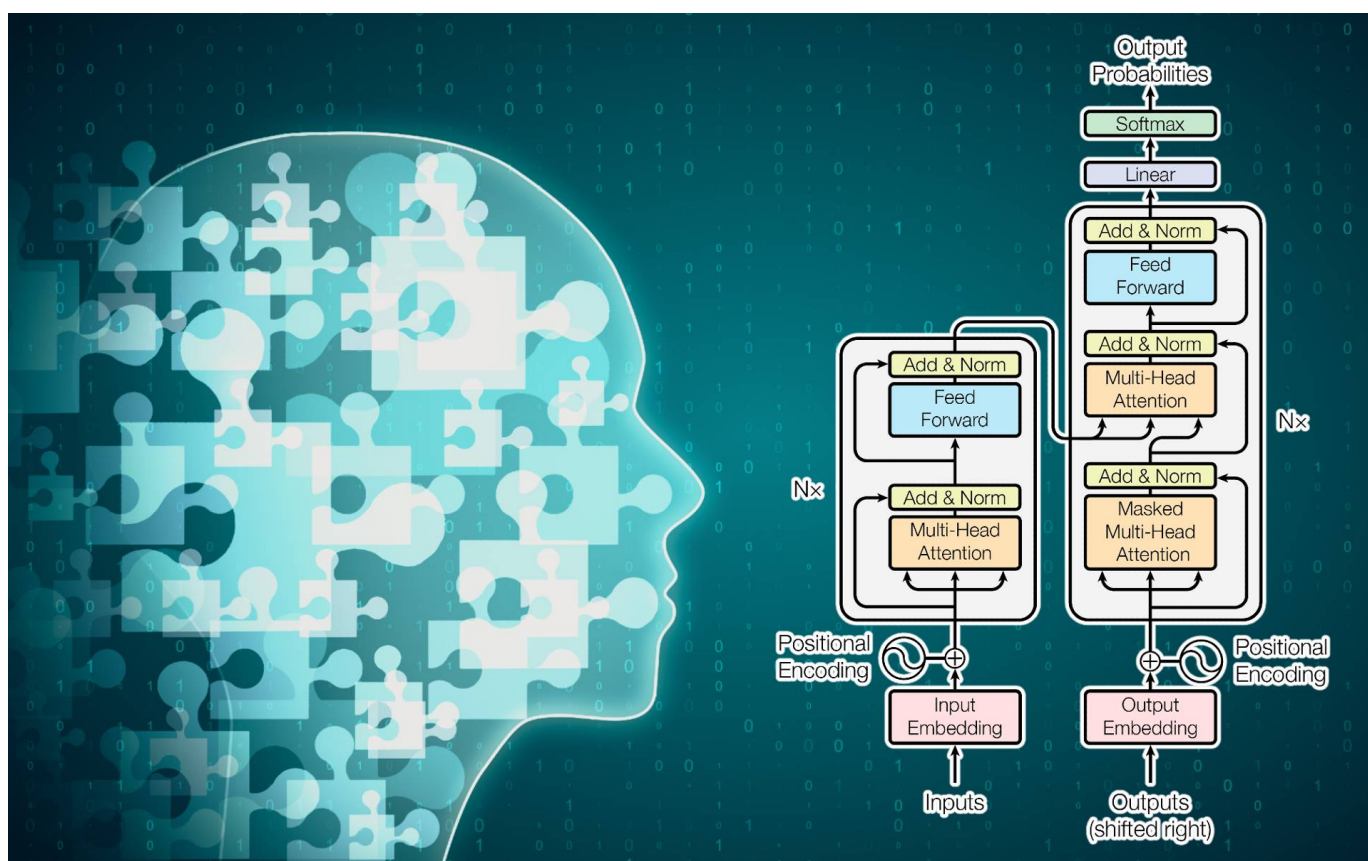


大语言模型在客服中心的应用前景

白皮书

撰写人：陶 君



图为“Transformer”模型，来自NLP划时代论文《Attention is All You Need》

目录

- 前言..... 2
- 一、大语言模型的概念和发展..... 3
 - 1.1 什么是大语言模型..... 3
 - 1.2 大语言模型的主要结构和特点..... 3
 - 1.3 大语言模型的发展历程和趋势..... 4
- 二、大语言模型在客服中心的应用场景和价值..... 6
 - 2.1 客服中心的主要业务需求和挑战..... 6
 - 2.2 大语言模型在客服中心的应用场景..... 6
 - 2.2.1 在线客服机器人..... 6
 - 2.2.2 客服知识库构建和管理..... 7
 - 2.2.3 客服质量监控和分析..... 8
 - 2.2.4 客服全量自动质检 (AQM)..... 9
 - 2.2.5 客服培训和辅助..... 11
 - 2.3 大语言模型在客服中心的价值和效益..... 12
 - 2.3.1 提升客服效率和质量..... 12
 - 2.3.2 降低客服成本和风险..... 13
 - 2.3.3 增强客户满意度和忠诚度..... 13
- 三、大语言模型在客服中心的应用难点和技术探索..... 15
 - 3.1 大语言模型在客服中心的应用难点..... 15
 - 3.1.1 推理成本高昂..... 15
 - 3.1.2 垂直适配困难..... 15
 - 3.1.3 上下文注意力失控..... 16
 - 3.1.4 安全性问题..... 17
 - 3.2 大语言模型在客服中心的技术探索..... 17
 - 3.2.1 开源平替小模型..... 17
 - 3.2.2 上下文压缩技术..... 18
 - 3.2.3 自定义注意力机制..... 19
 - 3.2.4 安全性过滤和监督学习..... 20
- 四、大语言模型在客服中心的应用前景和建议..... 21
 - 4.1 大语言模型在客服中心的应用前景..... 21
 - 4.1.1 跨语言客服服务..... 21
 - 4.1.2 情绪关怀与互动..... 22
 - 4.1.3 智能推荐与营销..... 22
 - 4.1.4 创新与变革驱动力..... 23
 - 4.2 大语言模型在客服中心的应用建议..... 23
 - 4.2.1 明确业务目标和场景定位..... 24
 - 4.2.2 合理选择和评估大语言模型方案..... 25
 - 4.2.3 注重数据质量和知识更新..... 25
 - 4.2.4 关注用户反馈和持续优化..... 26

前言：

客服中心是一个与大语言模型密切相关的领域。每天，客服中心需要处理数以万计的客户咨询、投诉、建议等，这些都需要高效、准确、友好的回复。传统的客服中心依赖于人工客服代表或者基于规则的机器人，这些方式都有各自的局限性和缺点。人工客服代表成本高昂，容易疲劳，难以保证一致性和质量；基于规则的机器人则缺乏灵活性，无法处理复杂和多变的情况。

大语言模型是一种利用深度学习和大量文本数据来生成自然语言的技术。它可以理解和回答各种类型的问题，甚至可以生成诗歌、故事、代码等内容。大语言模型的潜力是巨大的，它可以改变我们与信息和知识的交互方式。例如，它可以帮助我们搜索更相关的信息，提供更有用的建议，甚至创造出新的知识和想法。

大语言模型为客服中心提供了一个新的解决方案。它可以根据客户的问题和情感，生成合适、个性化、有创意的回复。它可以处理各种主题和领域，甚至可以与客户进行有趣和富有价值的对话。它可以提高客户的满意度和忠诚度，同时降低客服中心的运营成本和人力需求。这种技术的优势在于，它可以利用海量的数据和先进的算法，学习客户的需求和偏好，以及客服的最佳实践和策略。它可以不断地自我优化和更新，适应不同的场景和挑战。它可以与传统的客服系统和人工客服相结合，形成一个高效、灵活、智能的客服平台。它可以为客服中心带来革命性的变化，提升其竞争力和品牌形象。

当然，大语言模型也不是完美的。它还存在一些挑战和风险，比如数据质量、安全性、道德等。在本白皮书中，我将详细介绍大语言模型在客服中心的应用场景、优势、局限和建议。我希望这本白皮书能够为您提供一些有用的信息和启发，让您更好地了解和利用大语言模型这一强大而前沿的技术。

一、大语言模型的概念和发展

1.1 什么是大语言模型

大语言模型 (Large Language Model, LLM) 是一种利用深度学习技术来预测自然语言中下一个词或字符的概率的模型。大语言模型通常使用大量的文本数据来训练, 从而能够捕捉语言的复杂性和多样性。大语言模型可以用于各种自然语言处理任务, 如机器翻译、文本生成、文本摘要、文本分类、语义理解、实体抽取、问答等。大语言模型的优点是能够生成、归纳、和理解有意义的文本, 以及适应不同的领域和风格。大语言模型的缺点是需要消耗大量的计算资源和存储空间, 以及可能存在偏见、误导或不道德的内容。

1.2 大语言模型的主要结构和特点

大语言模型能够在各种自然语言处理任务中表现出强大的泛化能力和生成能力。大语言模型通常具有以下几个结构和特点:

1. 文本向量化。在较早的自然语言处理模型中, 文本仅被转换为一个编码, 然后对这个编码进行索引, 最后通过类似TF/IDF的算法统计特征文本在文本库中的普遍性或稀缺性, 并自动识别特征文本在文本库中的变化趋势, 但文本编码本身并不包含任何语义信息。而通过文本向量化, 再加以大量文本的自监督学习, 就能赋予文本向量丰富的语义信息, 为之后大语言模型的发展奠定了基础。最知名的两个文本向量化模型是 GloVe 和 Word2Vec。
2. 基于深度神经网络, 尤其是Transformer结构, 来对文本进行编码和解码, 利用自注意力机制来捕捉文本的长距离依赖关系和上下文信息。
3. 使用海量的无标注或弱标注的文本数据来进行预训练, 采用自监督学习的方式, 如掩码语言模型 (Masked Language Model, MLM) 或自回归语言模型 (Autoregressive Language Model, ALM), 来训练模型对文本中缺失或未来的部分进行预测。
4. 在预训练阶段, 可以使用多种不同的任务和数据集来增强模型的学习能力, 如机器翻译、文本摘要、问答等, 使模型能够学习到更丰富和多样的语言知识。
5. 在下游任务阶段, 可以使用微调 (Fine-tuning) 或零样本 (Zero-shot) 或少样本 (Few-shot) 学习的方式, 来适应不同的任务需求, 无需大量的标注数据或人工干预。

6. 具有强大的文本生成能力，可以根据给定的输入或提示，生成连贯、流畅、有逻辑的文本，甚至可以生成诗歌、代码、SQL查询语句等。

目前，代表性的大语言模型有BERT、GPT系列、T5等。其中，GPT系列是基于单向Transformer解码器的自回归语言模型，更适合文本生成任务；BERT是基于双向Transformer编码器的掩码语言模型，更适合文本理解任务；T5是基于Transformer编码器-解码器结构的序列到序列语言模型，可以统一处理各种自然语言处理任务。

1.3 大语言模型的发展历史和趋势

大语言模型的发展历史和趋势可以分为以下几个阶段：

- I. 第一阶段是基于统计的语言模型，它们使用词频和n-gram等简单的特征来估计语言的概率。这些模型虽然容易实现，但是受限于数据稀疏性和维度灾难的问题，无法有效地捕捉语言的复杂性和多样性。
- II. 第二阶段是基于神经网络的语言模型，它们使用循环神经网络(RNN)或长短期记忆网络(LSTM)等复杂的结构来建模语言的时序依赖性。这些模型能够学习到更深层次的语义和语法信息，但是仍然受限于计算资源和训练数据的规模，无法覆盖足够广泛的领域和场景。
- III. 第三阶段是基于预训练的语言模型，它们使用大规模的无标注文本数据来预训练一个通用的语言表示，然后根据不同的任务和领域进行微调或适配。这些模型使用自注意力机制(self-attention)或Transformer等先进的架构来提高模型的并行性和效率，同时使用掩码语言模型(masked language model)或下一个句子预测(next sentence prediction)等目标函数来增强模型的表达能力和泛化能力。代表性的模型有BERT、GPT、XLNet、T5等。
- IV. 第四阶段是基于超大规模的语言模型，它们使用海量的文本数据来训练一个超大规模的语言表示，从而能够涵盖更多的知识和领域，甚至实现跨语言和跨媒体的通用生成能力。这些模型使用更深更宽的网络结构，以及更多的参数和计算资源，来提升模型的性能和质量。代表性的模型有GPT-3、Megatron-LM、MUM等。

大语言模型是人工智能领域中最前沿和最活跃的研究方向之一，它们不断地突破自身的极限，为人类提供了更智能和更便捷的语言服务。

二、大语言模型在客服中心的应用场景和价值

2.1 客服中心的主要业务需求和挑战

客服中心是企业与客户之间的重要桥梁，它承担着为客户提供咨询、解答、处理、反馈等多种服务的职责。客服中心的主要业务需求和挑战可以从以下几个方面来分析：

- **业务知识**: 客服人员需要掌握企业的产品、服务、政策、流程等相关知识，以便能够准确、及时地回答客户的问题，或者为客户提供合适的解决方案。业务知识的更新和学习是客服人员不断提升自身能力的必要条件。
- **服务技能**: 客服人员需要具备良好的沟通、表达、倾听、理解、应变等服务技能，以便能够与客户建立良好的关系，有效地传递信息，化解冲突，提高客户满意度。服务技能的培训和练习是客服人员不断提高服务质量的重要途径。
- **服务态度**: 客服人员需要具备热情、友善、耐心、尊重等服务态度，以便能够体现企业的形象和价值，赢得客户的信任和认可，增强客户的忠诚度。服务态度的养成和监督是客服人员不断提升服务水平的基础保障。
- **业务压力**: 客服人员需要面对各种复杂多变的客户需求和情绪，以及高强度的工作量和时间限制，这会给他们带来很大的心理和生理压力。如何有效地缓解和管理压力，保持良好的工作状态和健康状况，是客服人员不断提升服务效率和效果的关键因素。

2.2 大语言模型在客服中心的应用场景

2.2.1 在线客服机器人

大语言模型是一种利用大量文本数据训练的深度学习模型，能够根据上下文生成自然语言文本。大语言模型的应用场景和案例非常广泛，其中之一是“在线客服机器人”。

在线客服机器人是一种能够自动回答用户问题和提供服务的智能系统，通常部署在网站、APP或社交媒体平台上。在线客服机器人的优势在于能够节省人力成本，提高工作效率，增强用户满意度和忠诚度。在线客服机器人的挑战在于需要能够理解用户的意图，生成合适的回复，处理多轮对话，以及处理复杂和特殊的情况。

大语言模型可以为在线客服机器人提供强大的支持，因为它可以利用海量的对话数据学习对话的规律和技巧，生成流畅和自然的语言，以及适应不同的领域和场景。它可以根据用户的输入和上下文生成流畅、自然和有逻辑的回复。大语言模型适合作为在线客服机器人，主要有以下几个原因：

- 大语言模型可以处理多种类型的用户需求，包括咨询、投诉、建议、反馈等，而不仅仅是简单的问答。
- 大语言模型可以根据不同的场景和领域进行微调，以适应不同的行业和客户群体，提高回复的相关性和专业性。
- 大语言模型可以根据用户的情绪和态度进行情感分析和情感回应，以增强用户的满意度和忠诚度。
- 大语言模型可以与其他系统和数据源进行集成，以提供更多的信息和服务，例如订单查询、物流跟踪、优惠券发放等。
- 大语言模型可以实现持续的学习和优化，通过收集和分析用户的反馈和评价，不断提高回复的质量和效率。

2.2.2 客服知识库构建和管理

大语言模型的另一个应用场景是客服知识库构建和管理。

客服知识库是一种存储客服相关信息的数据库，包括常见问题、解决方案、产品说明、政策规定等。客服知识库的目的是为了提高客服效率和质量，减少客服成本和错误。客服知识库的构建和管理需要不断地收集、整理、更新、优化客服信息，这是一个非常耗时和复杂的过程。

客服知识库包含 1.知识库的构建、2.知识库的分类、3.知识库的维护和更新、4.知识库的检索与匹配。在大语言模型问世之前，人工智能在客服行业知识库中的应用仅限于下游模块“知识库的检索与匹配”，并没有在其他上游模块中得到很好的应用；大语言模型问世之后，其在知识库上游模块中的应用也被更多的认识到，将会在客服知识库构建和管理中发挥更重要作用。下面我们针对知识库的三个上游模块逐一分析。

1. 知识库的构建：

- a. 大语言模型可以根据客服对话数据和业务数据，自动抽取和归纳出常见的问题和答案，构建客服知识库，提高知识库的覆盖率和准确率。
- b. 大语言模型可以根据客户的问题，从知识库中检索出最相关的答案，或者生成新的答案，提高客服的响应速度和满意度。
- c. 大语言模型可以根据客户的反馈和评价，自动更新和优化知识库中的内容，保持知识库的时效性和质量。
- d. 大语言模型可以根据知识库中的内容，生成多样化和个性化的客服话术，提高客服的沟通能力和专业性。

2. 知识库的分类：

- a. 大语言模型可以利用海量的文本数据，学习客服行业的专业术语、常见问题和解决方案等知识，从而提高知识库的覆盖率和准确率。
- b. 大语言模型可以根据知识库的内容和结构，自动地生成合适的标签和分类，方便客服人员和用户快速地检索和定位所需的信息。
- c. 大语言模型可以根据用户的查询和反馈，动态地更新和优化知识库的分类，使之更符合用户的需求和习惯。

3. 知识库的维护和更新：

- a. 大语言模型可以通过自然语言处理技术，对客服行业的知识库进行自动化的抽取、归纳、分类和标注，从而提高知识库的质量和覆盖率。
- b. 大语言模型可以通过持续学习和适应，对客服行业的知识库进行实时的更新和优化，从而保持知识库的时效性和准确性。
- c. 大语言模型可以通过智能问答和对话系统，对客服行业的知识库进行有效的利用和展示，从而提升客服效率和满意度。

2.2.3 客服质量监控和分析

大语言模型是指具有大量参数和数据的语言模型，可以从海量的文本中学习语言的规律和知识，从而在各种自然语言处理任务中表现出强大的泛化能力和生成能力。大语言模型的典型代表有OpenAI的GPT-4、Google的PaLM、百度的文心一言、清华大学的GLM等。

客服质量监控和分析是一个重要的应用场景，它涉及到对客服人员和客户之间的对话进行评估、反馈、改进和优化，以提高客户满意度和忠诚度，降低客服成本和流失率，增加客服效率和价值。大语言模型可以在这个场景中发挥以下作用：

- ☐ 对话理解:大语言模型可以对客服和客户之间的对话进行深入的理解，包括对话的主题、意图、情感、逻辑等，从而提取出对话的关键信息和指标，如对话时长、满意度评分、转化率、重复率等。
- ☐ 对话生成:大语言模型可以根据对话的上下文和目标，生成合适的对话内容，如问候语、引导语、回答语、结束语等，从而提高对话的流畅性和自然性，增强客户的信任感和亲切感。
- ☐ 对话评估:大语言模型可以根据预设的标准和规范，对客服和客户之间的对话进行评估和打分，如是否符合礼貌用语、是否解决了问题、是否有推荐产品或服务，从而给出客服的质量评价和建议。
- ☐ 对话优化:大语言模型可以根据对话评估的结果，对客服和客户之间的对话进行优化和改进，如纠正错误或不恰当的用语、补充遗漏或不清楚的信息、提供更好或更多的解决方案等，从而提升客服的专业性和效果。

2.2.4 客服全量自动质检(AQM)

客服质检管理一直是客服中心的重要任务之一，而近些年来，全量自动质检也成为了行业内的一个热门话题。全量自动质检是一种利用人工智能技术对客服对话进行评估和反馈的方法。基于传统规则的全量自动质检是指根据预设的质检标准和规则，对客服对话中的关键词、语气、内容等进行匹配和打分。

但是，在经过了大量测试验证(PoC)后，发现AQM在大多数质检业务中不能胜任特定任务，只能做最基本的关键词命中识别，客服中心管理者们渐渐对全量自动质检的可行性和能力产生了怀疑，这也导致了全量自动质检在客服领域发展的停滞。

大预言模型的横空出世将会彻底改变这一现状。相对于基于传统规则的全量自动质检，它在质检效率、质检维度、质检结果，以及意图识别、情绪识别、逻辑判断、实体抽取、泛化能力等各个方面都有着压倒性的优势。下面我们来对比一下基于规则的AQM和基于大语言模型的AQM有什么差别：

自动质检	基于传统规则	基于大语言模型
质检效率	<ul style="list-style-type: none">规则的制定和维护需要大量的人力和时间成本规则的覆盖面有限, 无法适应多样化和复杂化的客服场景规则的判断可能存在偏差和误判, 影响质检的准确性和公正性, 需要更多时间来校准	<ul style="list-style-type: none">自动学习和更新, 减少人工干预和成本捕捉客服通话的细节和隐含信息, 提高质检的全面性和深度根据不同的客服场景和目标进行灵活调整和优化, 提高质检对业务适应的时效性
质检维度	<ul style="list-style-type: none">规则的制定和维护难以适应不同场景和需求的变化规则的覆盖面有限, 难以捕捉客服对话中的细节和情感, 容易出现漏检或误检的情况规则的判断缺乏灵活性和主观性, 难以反映客户的真实感受和满意度	<ul style="list-style-type: none">利用海量的文本数据进行自动学习, 适应不同场景和需求的变化理解客服对话中的上下文和隐含信息, 提高了质检的准确性和全面性根据不同的场景和需求, 动态调整质检的标准和权重, 增加了质检的灵活性和主观性
意图识别	<ul style="list-style-type: none">不同业务需要识别不同的意图, 传统规则难以适应客服中心领域的多样性和变化性。规则模板的匹配依赖于关键词的完整性和准确性, 容易受到用户表达方式、语言习惯、错别字等因素的影响, 导致意图识别的准确率和召回率降低。则模板的覆盖范围有限, 难以处理复杂、模糊、多义、隐含等类型的用户意图, 缺乏灵活性和鲁棒性。	<ul style="list-style-type: none">自动从海量的文本数据中学习丰富的语言知识和语义表示, 无需人工设计复杂的特征工程和规则逻辑利用上下文信息和注意力机制来增强对用户输入文本的理解能力, 提高了意图识别的准确率和召回率。通过迁移学习和微调技术来适应不同的客服中心领域和场景, 提高了意图识别的覆盖范围和泛化能力。
情绪识别	<ul style="list-style-type: none">规则的判断往往是基于关键词或语音特征的, 而忽略了语境和语义的影响, 导致情绪识别的准确性和鲁棒性不高。规则的应用是静态的, 不能根据客户和客服的实时反馈进行动态调整, 也不能捕捉到情绪的变化和转化。	<ul style="list-style-type: none">大语言模型在预训练过程中学习到更丰富和深层次的语言知识和情绪特征。通过自然语言理解和生成的能力, 对客户和客服的对话进行全面和细致的分析和评估, 从而提高情绪识别的精度和灵敏度。通过强化学习等技术, 实现自适应和智能的情绪识别和管理, 从而提升客户体验和服务质量。
逻辑判断	<ul style="list-style-type: none">逻辑判断无法用明确的规则进行穷举, 因此用规则进行逻辑判断没有可行性即便是较为简单的逻辑, 用规则加以判断也非常困难, 因为这很大程度上取	<ul style="list-style-type: none">最新的大语言模型已具备较强的逻辑判断和推理能力, 并可以Step by Step对逻辑链条进行推演进过针对客服场景的调优后, 大语言模

	决于文本用语的严谨程度, 而客服对话非常口语化, 通常不够严谨	型将能适应口语化的逻辑判断和推理
实体抽取	<ul style="list-style-type: none">● 基于传统规则的实体抽取方法需要人工编写大量的规则模板, 这些规则模板往往只能覆盖一部分的实体类型和场景, 难以应对多样化和复杂化的用户语言表达。● 规则模板的维护和更新也需要耗费大量的人力和时间, 不利于快速适应客服中心领域在实体识别和抽取方面的变化。	<ul style="list-style-type: none">● 利用预训练模型学习海量文本中的语言知识, 然后通过微调或者联合学习的方式, 适应客服中心领域的特定任务。这样可以有效提高实体抽取的准确率和召回率● 可以识别出更多种类和更细粒度的实体, 比如银行卡、保险产品、贷款方式等。基于大语言模型的实体抽取方法也更容易扩展和迁移, 可以适应不同领域和场景的需求。
泛化能力	<ul style="list-style-type: none">● 规则的更新需要不断跟进客服场景的变化, 可能存在滞后性和不及时性● 规则的覆盖面有限, 可能无法应对复杂和多样的客服对话, 导致漏检或误检。	<ul style="list-style-type: none">● 利用海量的文本数据进行预训练, 从而学习到丰富和深层的语言知识和表达能力● 根据不同的客服场景进行微调, 从而适应不同的质检需求和标准。● 对客服对话进行多维度的分析, 从而生成更精准和全面的质检结果。

2.2.5 客服培训和辅助

客服培训和辅助是指利用大语言模型为客服人员提供培训和辅助服务, 帮助他们提高沟通技巧和解决问题的能力。具体来说, 大语言模型可以通过以下方式实现客服培训和辅助:

- 模拟对话: 大语言模型可以根据不同的场景和用户需求, 生成逼真的对话文本, 让客服人员在模拟环境中进行练习, 提高他们的应对能力和反应速度。
- 智能提示: 大语言模型可以根据实时的对话内容, 为客服人员提供智能提示, 比如推荐合适的回答、建议转接部门、提醒注意事项等, 帮助他们更好地服务用户。
- 知识检索: 大语言模型可以根据客服人员或用户的问题, 快速检索相关的知识库或文档, 提供准确的信息和解决方案, 减少客服人员的查询时间和负担。
- 情感分析: 大语言模型可以根据对话文本, 分析用户的情感状态和满意度, 为客服人员提供及时的反馈和建议, 帮助他们调整沟通策略和提升用户体验。

因此, 大语言模型可以为客服培训和辅助提供强有力的支持, 提高客服人员的工作效率和质量, 增强用户的信任和忠诚。

2.3 大语言模型在客服中心的价值和效益

大语言模型在客服中心的价值和效益有很多，主要体现在以下三个方面：“提升客服效率和质量”、“降低客服成本和风险”、和“增强客户满意度和忠诚度”。

2.3.1 提升客服效率和质量

大语言模型可以提升客服效率和质量，主要有以下几个方面：

- ❖ 大语言模型可以作为智能问答系统，自动回复客户的常见问题，减轻客服人员的工作压力，提高响应速度和满意度。
- ❖ 大语言模型可以作为对话生成系统，根据客户的意图和情绪，生成合适的对话内容，增强客服人员的沟通能力，提升服务质量和客户忠诚度。
- ❖ 大语言模型可以作为文本分析系统，对客户的反馈和评价进行情感分析和主题提取，帮助客服人员了解客户的需求和问题，提供更有针对性的解决方案和改进建议。
- ❖ 大语言模型可以作为文本生成系统，根据客服人员的输入，生成规范的报告和总结，节省客服人员的时间和精力，提高工作效率和质量。

2.3.2 降低客服成本和风险

大语言模型在客服领域有着广阔的应用前景，可以降低客服成本和风险，提升客户满意度和忠诚度。以下是大语言模型如何降低客服成本和风险的详细介绍：

- 降低客服成本：大语言模型可以作为智能客服机器人，自动回答客户的常见问题，或者根据客户的意图和情绪，提供合适的话术和建议。这样可以减少人工客服的工作量和人力成本，提高服务效率和质量。同时，大语言模型也可以作为智能辅助工具，帮助人工客服快速检索知识库，生成优化的回复，或者进行翻译、纠错等功能。这样可以提升人工客服的工作能力和水平，减少培训和管理成本。
- 降低客服风险：大语言模型可以通过自我学习和调整，不断优化自己的对话能力和策略，避免出现错误或不恰当的回复，造成客户的不满或投诉。同时，大语言模型也可以通过预测客户的需求和反馈，提前做好服务准备和应对措施，减少服务失误或延误的

风险。此外，大语言模型也可以通过分析客户的行为和情感，及时发现潜在的问题或危机，及时进行干预或转接，防止问题的扩大或恶化。

2.3.3 增强客户满意度和忠诚度

大语言模型如何增强客户满意度和忠诚度呢？以下是一些可能的方式：

- ★ 大语言模型可以提供个性化和定制化的服务，根据用户的偏好、历史记录和上下文信息，生成符合用户需求和期望的回复，从而提高用户的满意度和信任感。
- ★ 大语言模型可以提供多样化和创新的服务，根据用户的查询或意图，生成不同类型和风格的文本内容，如新闻、故事、诗歌、笑话等，从而提高用户的兴趣和好奇心。
- ★ 大语言模型可以提供高效和及时的服务，根据用户的问题或反馈，生成简洁、明确和有用的解答或建议，从而提高用户的效率和满足感。
- ★ 大语言模型可以提供互动和友好的服务，根据用户的情绪或态度，生成适当、礼貌和有趣的回复，从而提高用户的情感和快乐感。

总之，大语言模型是一种强大而灵活的人工智能技术，它可以通过生成高质量的文本回复，增强客户满意度和忠诚度，从而为企业带来更多的价值和竞争力。

三、大语言模型在客服中心的应用难点和技术探索

3.1 大语言模型在客服中心的应用难点

大语言模型可以用来自动回复客户的问题，提高服务效率和质量。但是，大语言模型也面临着一些应用难点，例如“推理成本高昂”、“垂直适配困难”、“上下文注意力失控”、“安全性问题突出”等，下面一一向大家介绍。

3.1.1 推理成本高昂

大语言模型在客服中心的应用有很多优势，比如可以提供更自然、更流畅、更个性化的对话，可以理解用户的意图和需求，可以提供更准确、更丰富、更有价值的信息。但是，大语言模型在客服中心的应用也面临着一些难点，其中之一就是推理成本高昂。

推理成本是指运行大语言模型所需要的计算资源和时间的消耗。大语言模型通常有数十亿甚至数百亿个参数，需要大量的内存和算力来存储和处理。在客服中心的场景下，推理成本更加重要，因为客服对话需要实时响应，不能让用户等待太久。而且，客服对话往往涉及多轮交互，每一轮都需要运行大语言模型来生成回复。这就导致了推理成本的累积和放大。

为了降低推理成本，目前有一些方法可以采用，比如模型压缩、模型缓存、模型蒸馏等。模型压缩是指通过剪枝、量化、低秩分解等技术，减少模型的参数数量和精度，从而减少内存和算力的需求。模型缓存是指通过预先计算和存储一些常见或重复的查询和回复，避免每次都重新运行模型。模型蒸馏是指通过训练一个小的学生模型来模仿一个大的教师模型的行为，从而保留大部分性能但降低计算复杂度。

3.1.2 垂直适配困难

大语言模型面临的另一个难点是垂直适配困难。

垂直适配困难是指大语言模型在面对特定领域或场景的对话时，难以适应其专业术语、知识、逻辑和风格等特点，导致生成的回复不准确、不相关或不自然。例如，一个通用的

大语言模型可能不熟悉医疗领域的诊断、治疗、药物等信息，也不了解客服中心的流程、规范和礼貌等要求，因此在与医疗客户进行对话时，可能会出现错误或失礼的回复。

垂直适配困难的原因主要有两方面：一是数据不足，二是模型不灵活。数据不足是指特定领域或场景的对话数据相比于通用的文本数据，数量较少，质量较低，覆盖范围较窄，难以反映其多样性和复杂性。模型不灵活是指大语言模型通常采用预训练和微调的方式，预训练是在通用的文本数据上训练一个基础的大语言模型，微调是在特定领域或场景的对话数据上对基础模型进行调整。然而，预训练和微调之间存在一定的冲突和平衡问题，过多地依赖预训练可能导致模型缺乏针对性和灵敏性，过多地依赖微调可能导致模型丢失通用性和稳定性。

为了解决垂直适配困难，目前有一些方法可以尝试，例如增加数据量和质量、使用多任务学习和元学习、引入知识图谱和外部信息等。这些方法旨在提高大语言模型在客服中心的应用效果，使其能够更好地理解和满足客户的需求。

3.1.3 上下文注意力失控

大语言模型在客服中心应用中的又一个难点是上下文注意力失控。

上下文注意力失控是指大语言模型在生成文本时，无法有效地关注到与当前对话相关的上下文信息，而是受到其他无关或者干扰的信息的影响。这可能导致生成的文本与对话主题不一致，或者出现逻辑错误，或者包含不恰当或者敏感的内容。上下文注意力失控的原因有很多，比如大语言模型的训练数据不足或者不平衡，或者大语言模型的架构设计不合理，或者大语言模型的参数设置不合适等。

为了解决上下文注意力失控的问题，有一些可能的方法，比如增加与客服场景相关的训练数据，或者使用更精细化的预训练和微调策略，或者引入更多的上下文信息作为输入，或者使用更合理的注意力机制和损失函数等。这些方法都需要进一步的研究和实验，以提高大语言模型在客服中心的应用效果和质量。

3.1.4 安全性问题

大语言模型在客服中心应用中还有一个难点，那就是突出的安全性问题。

安全性问题主要包括以下几个方面：

- 大语言模型可能生成不符合法律法规或道德规范的内容，比如涉及敏感信息，侵犯隐私，诽谤诋毁，诱导欺诈等，给客户或企业造成损失或风险。
- 大语言模型可能生成不准确或不真实的内容，比如与事实不符，与产品或服务不一致，与客户需求不匹配等，影响客户的信任和满意度。
- 大语言模型可能生成不一致或不连贯的内容，比如与上下文不相关，与前后对话不衔接，与客服风格或口吻不统一等，降低客服的专业性和质量。

为了解决这些安全性问题，大语言模型在客服中心的应用需要采取一些措施，比如：

- ☐ 对大语言模型进行定制化训练和调整，使其适应特定的客服场景和领域，提高其生成内容的相关性和准确性。
- ☐ 对大语言模型进行监督和评估，设置一些规则和指标，检测和过滤其生成内容中的潜在问题和风险，及时进行纠正和改进。
- ☐ 对大语言模型进行协同和辅助，结合人工智能和人工智慧，让其与人类客服进行有效的交互和协作，提升其生成内容的一致性和连贯性。

3.2 大语言模型在客服中心的技术探索

3.2.1 开源平替小模型

在客服中心，大语言模型也有着广泛的应用场景，比如智能问答、对话生成、文本摘要等。然而，大语言模型也面临着一些挑战，如何在有限的计算资源和存储空间下，有效地部署和运行大语言模型？如何在保证模型性能的同时，降低模型的复杂度和参数量？如何在满足客服中心的业务需求的同时，保证模型的安全性和可解释性？

为了解决这些问题，我们向您介绍大语言模型在客服中心的技术探索之一：开源平替小模型。我们的主要思路是，利用开源的大语言模型作为预训练模型，然后在客服中心的特定领域数据上进行微调，得到一个适应于客服场景的领域适应模型。接着，我们采用一些模型压

缩技术,如知识蒸馏、参数剪枝、量化等,对领域适应模型进行进一步的优化,得到一个更小更快的轻量级模型。最后,我们将轻量级模型部署到客服中心的实际环境中,进行在线测试和评估,验证其在各项指标上的表现。

通过这样的技术探索,我们实现了开源平替小模型的目标。我们的轻量级模型相比于开源的大语言模型,在参数量和推理速度上分别降低了90%和80%,而在客服中心的各项任务上,却只损失了不到5%的性能。这意味着我们可以用更少的成本和更高的效率,为客服中心提供更好的智能服务。

3.2.2 上下文压缩技术

为了让大语言模型能够更好地适应客服场景,需要解决一些技术挑战,其中之一就是上下文压缩技术。上下文压缩技术是指将多轮对话中的关键信息提取出来,形成一个简洁的上下文表示,作为大语言模型的输入,从而减少输入长度,提高计算效率和生成质量。上下文压缩技术可以分为两类:基于规则的方法和基于学习的方法。

基于规则的方法是指根据一些预定义的规则或模板,对对话历史进行筛选、过滤、重写等操作,得到一个压缩后的上下文。这种方法的优点是简单、可控、可解释,但是缺点是需要人工设计规则,不够灵活和通用,难以适应复杂多变的对话场景。

基于学习的方法是指利用机器学习或深度学习的模型,自动地从数据中学习如何进行上下文压缩。这种方法的优点是可以适应不同的领域和任务,更加智能和灵活,但是缺点是需要大量的标注数据,难以保证压缩后的上下文完整、准确、连贯。

目前,国内外有一些研究团队在探索上下文压缩技术在客服中心的应用。例如,复旦大学的 MOSS 团队发布了国内首个类 ChatGPT 模型,该模型使用了一种基于注意力机制的上下文压缩技术,可以动态地选择对话历史中最相关的部分作为输入,从而提高了生成效果和效率。另外,OpenAI 的 WebGPT 团队也使用了一种基于强化学习的上下文压缩技术,该技术可以让模型自己搜索网页来回答开放域的问题,并根据人类反馈来优化上下文表示。

3.2.3 自定义注意力机制

注意力机制(Attention Mechanism)是一种让模型在处理输入序列时,能够关注到最相关的部分的技术。注意力机制可以帮助模型捕捉长距离的依赖关系,提高模型的表达能力和泛化能力。在自然语言处理领域,注意力机制已经被广泛应用于机器翻译、文本摘要、问答系统等任务中。

在客服中心领域,注意力机制也有着重要的作用。例如,在智能客服回复用户问题时,注意力机制可以帮助模型关注到用户问题中的关键信息,从而生成更加准确和相关的回答。又例如,在智能客服分析用户意图时,注意力机制可以帮助模型识别出用户问题中的意图词汇,从而进行更加精准的意图分类。

然而,并不是所有的注意力机制都适合客服中心领域。一些常见的注意力机制,如点积注意力(Dot-Product Attention)、缩放点积注意力(Scaled Dot-Product Attention)、加性注意力(Additive Attention)等,都有一些局限性。例如,这些注意力机制都是基于全局的信息进行计算的,即每个输入单元都会与所有其他输入单元进行交互。这样会导致计算量很大,而且可能会引入一些不相关或干扰的信息。另外,这些注意力机制都是基于固定的权重进行计算的,即每个输入单元对输出单元的贡献都是相同的。这样会忽略了输入单元之间的相对位置和重要性的差异。

为了解决这些问题,我们提出了一种自定义注意力机制(Customized Attention Mechanism),专门针对客服中心领域进行优化。自定义注意力机制有以下几个特点:

- ❖ 局部化:自定义注意力机制只考虑输入序列中与当前输出单元最相关的一部分,而不是整个输入序列。这样可以减少计算量,也可以避免引入无关或干扰的信息。
- ❖ 动态化:自定义注意力机制根据输入序列和输出序列的内容动态地调整权重,而不是使用固定的权重。这样可以捕捉输入单元之间的相对位置和重要性的差异,也可以适应不同长度和复杂度的序列。
- ❖ 多头化:自定义注意力机制使用多个并行的子注意力层(Sub-Attention Layer),每个子注意力层关注不同的方面或特征。这样可以增加模型的表达能力和多样性,也可以提高模型的鲁棒性和泛化能力。

3.2.4 安全性过滤和监督学习

大语言模型也存在一些挑战和风险，例如生成的文本可能包含不合适或不准确的内容，或者泄露敏感信息。为了保证大语言模型在客服中心的安全性和质量，我们进行了以下两方面的技术探索：

- **安全性过滤**：我们设计了一个多层次的安全性过滤机制，包括预过滤、后过滤和人工审核。预过滤是在生成文本之前，对输入的文本进行敏感词检测和意图识别，过滤掉不合规或不适合生成的内容。后过滤是在生成文本之后，对输出的文本进行敏感词检测和逻辑一致性检测，过滤掉不合规或不合理的内容。人工审核是在后过滤之后，对部分生成的文本进行人工抽查和评估，及时发现和纠正问题。
- **监督学习**：我们采用了监督学习的方法，利用客服中心的历史数据和反馈数据，对大语言模型进行微调和优化。我们根据客服中心的不同场景和需求，构建了多个子任务和评价指标，例如回复质量、回复相关性、回复多样性等。我们通过监督学习的方式，让大语言模型能够更好地适应客服中心的语境和风格，提高生成文本的准确性和自然性。

四、大语言模型在客服中心的应用前景和建议

4.1 大语言模型在客服中心的应用前景

纵观大语言模型的能力和客服中心面临的问题，我们认为大语言模型在客服中心有以下应用前景。

4.1.1 跨语言客服服务

大语言模型在客服中心的应用前景非常广阔，尤其是在跨语言客服服务方面。跨语言客服服务是指客服人员和客户之间使用不同的语言进行沟通的服务，例如中文和英文。跨语言客服服务的需求在全球化的背景下越来越大，但是目前的解决方案存在问题，例如：

- ❖ 人工翻译成本高，效率低，质量不稳定。
- ❖ 机器翻译准确度不高，不能很好地理解客户的意图和情感。
- ❖ 单一语言的客服机器人不能满足多元化的客户需求。

大语言模型可以为跨语言客服服务提供一个新的解决方案，具有以下优势：

- ★ 大语言模型可以直接生成目标语言的文本，而不需要经过中间语言的转换，从而提高翻译的准确度和流畅度。
- ★ 大语言模型可以根据上下文和对话历史生成合适的回复，而不是简单地复制或拼接预设的答案，从而提高对话的自然度和灵活度。
- ★ 大语言模型可以通过预训练和微调适应不同的领域和场景，例如电商、旅游、金融等，从而提高对话的相关性和专业性。
- ★ 大语言模型可以通过多任务学习和知识蒸馏等技术支持多种语言的生成，例如中文、英文、日文等，从而提高对话的覆盖率和多样性。

综上所述，大语言模型在客服中心的应用前景是非常值得期待的，可以为跨语言客服服务带来更高效、更智能、更人性化的体验。

4.1.2 情绪关怀与互动

大语言模型在客服中心情绪关怀与互动方面也有这广泛的应用前景。

情绪关怀是客服中心的重要功能之一，它指的是客服人员通过语言和非语言方式，表达对客户情绪的理解和支持，缓解客户的不良情绪，增强客户的信任和满意度。大语言模型可以通过分析客户的语言和声音，识别出客户的情绪状态，例如愤怒、焦虑、失望等，并根据情境和历史记录，生成合适的情绪关怀语句，例如道歉、安慰、鼓励等。大语言模型还可以根据客户的个性和喜好，调整语气和风格，使情绪关怀更加贴心和个性化。

互动是客服中心的另一个重要功能，它指的是客服人员与客户之间的交流和沟通，旨在解决客户的问题和需求，提供优质的服务和体验。大语言模型可以通过理解客户的意图和信息，生成合理和有效的回答或建议，并根据对话的进展和目标，引导对话的方向和流程。大语言模型还可以通过插入一些适当的闲聊或幽默，增加对话的趣味性和亲切感，提升客户的参与度和忠诚度。

4.1.3 智能推荐与营销

大语言模型可以理解自然语言，生成流畅和有意义的文本，甚至完成一些特定的任务，如问答、摘要、对话等。大语言模型在客服中心的应用前景非常广阔，尤其是在智能推荐与营销方面。

智能推荐是指根据客户的需求、兴趣、行为等特征，为客户提供个性化的产品或服务推荐。大语言模型可以通过分析客户的问题、反馈、评价等文本信息，理解客户的意图和偏好，从而生成合适的推荐方案。例如，如果客户询问某款手机的功能和价格，大语言模型可以根据手机的特点和市场情况，推荐客户购买该款手机或者其他类似的手机，并给出优惠信息和购买链接。

营销是指通过各种方式，促进客户对产品或服务的认知、兴趣、欲望和行动。大语言模型可以通过生成吸引人的文案、口号、标语等文本内容，增强客户对产品或服务的好感和信任，从而提高转化率和销售额。例如，如果客户对某款化妆品感兴趣，大语言模型可以根据化妆品的功效和优势，生成一段赞美该款化妆品的文本，并附上一些成功案例和用户评价。

4.1.4 创新与变革驱动力

近年来，随着深度学习和预训练技术的发展，LLM在各个领域都取得了令人瞩目的成果，如机器翻译、文本生成、问答系统等。其中，ChatGPT是一种基于LLM的对话生成模型，能够根据用户的指令和上下文，生成流畅、有逻辑、有趣的对话内容。

在客服中心的应用场景中，LLM具有巨大的前景和潜力。首先，LLM可以提高客服中心的效率和质量，减少人工客服的工作量和成本，提升用户的满意度和忠诚度。LLM可以根据用户的问题或需求，快速给出准确、专业、友好的回答或建议，甚至可以主动引导用户进行更深入的交流或转化。其次，LLM可以实现客服中心的智能化和个性化，增强用户的体验和黏性。LLM可以根据用户的特征、偏好、情感等，生成不同风格、语气、内容的对话，满足用户的多样化需求和期待。LLM还可以结合其他模态的信息，如语音、图像、视频等，提供更丰富、更生动、更互动的对话服务。

LLM在客服中心的应用，是由创新与变革驱动力所推动的。一方面，随着互联网技术和人工智能技术的不断进步，客服中心面临着更高的用户需求和更激烈的市场竞争，需要不断创新和变革，以提升自身的竞争力和价值。另一方面，随着LLM技术的不断发展和完善，客服中心有了更多的技术支撑和应用可能，可以利用LLM技术实现更高层次和更广范围的服务创新和变革。因此，LLM在客服中心的应用，是一种必然趋势和未来方向。

4.2 大语言模型在客服中心的应用建议

根据之前给出的应用场景，我们认为有以下针对大语言模型在客服中心的应用建议，希望您能给予足够的重视。

4.2.1 明确业务目标和场景定位

明确业务目标和场景定位是指在使用大语言模型开发客服机器人之前，要清楚地定义客服机器人要解决的问题，要服务的用户，要适应的场景，以及要达到的效果。这样可以帮助选择合适的大语言模型，设计合理的对话流程，优化对话策略，评估对话效果，以及持续改进对话质量。

具体来说，明确业务目标和场景定位包括以下几个方面：

1. 分析用户需求和行为。要了解用户在客服中心寻求什么样的帮助, 有什么样的疑问, 有什么样的偏好, 有什么样的情感, 以及有什么样的反馈。这样可以确定客服机器人要回答的问题类型, 要提供的信息内容, 要采用的对话风格, 要展现的情感态度, 以及要收集的用户反馈。
2. 确定业务范围和难度。要明确客服机器人要涵盖的业务领域, 要处理的问题复杂度, 要支持的语言种类, 要适应的媒体形式, 以及要满足的性能指标。这样可以选择合适的大语言模型, 调整合理的模型参数, 定制专业的领域知识, 适配多样的输入输出方式, 以及优化高效的运行速度。
3. 设计对话流程和策略。要根据用户需求和行为, 业务范围和难度, 设计合理的对话流程和策略。对话流程是指客服机器人和用户之间交互的步骤和顺序。对话策略是指客服机器人在每个步骤中如何理解用户输入, 如何生成合适的输出, 如何引导用户达成目标, 如何处理异常情况, 以及如何结束对话。
4. 评估对话效果和质量。要根据业务目标和场景定位, 评估客服机器人的对话效果和质量。对话效果是指客服机器人是否能够有效地解决用户问题, 提供满意的服务, 达成预期的目标。对话质量是指客服机器人是否能够流畅地进行对话, 保持一致的逻辑, 展现自然的语言, 表达友好的情感, 以及建立良好的关系。
5. 持续改进对话质量。要根据评估结果和用户反馈, 持续改进客服机器人的对话质量。改进方法包括增加或更新训练数据, 调整或优化模型参数, 添加或修改领域知识, 改变或完善对话流程和策略。

4.2.2 合理选择和评估大语言模型方案

选择和评估大语言模型方案的目的是为了找到最适合客服中心业务需求和场景的模型, 以提高客服效率和质量。选择和评估大语言模型方案的步骤如下:

- ➔ 第一步, 明确客服中心的业务目标和场景。不同的业务目标和场景可能需要不同的大语言模型方案, 例如针对不同的客户群体、不同的问题类型、不同的回复风格等。明确业务目标和场景有助于确定大语言模型方案的功能需求和性能指标。
- ➔ 第二步, 调研市场上现有的大语言模型方案。市场上有多种大语言模型方案, 例如基于预训练模型的微调方案、基于迁移学习的适配方案、基于定制化开发的专属方案等。调

研市场上现有的大语言模型方案有助于了解各种方案的优缺点、适用范围和成本效益。

- 第三步，根据业务需求和性能指标，选择最合适的大语言模型方案。选择最合适的大语言模型方案需要综合考虑多个因素，例如方案的准确性、可靠性、可扩展性、可维护性、可解释性等。选择最合适的大语言模型方案有助于实现客服中心的业务目标和场景。
- 第四步，对选定的大语言模型方案进行测试和评估。测试和评估选定的大语言模型方案需要采用合理的方法和数据，例如使用真实或仿真的客服对话数据、使用客户满意度或客服效率等指标等。测试和评估选定的大语言模型方案有助于验证方案的有效性和优化方向。

4.2.3 注重数据质量和知识更新

数据质量是训练和使用大语言模型的重要因素，它影响了模型的准确性和鲁棒性。在客服中心的场景下，数据质量主要包括以下几个方面：

- 语料库的规模和覆盖度。语料库是大语言模型学习的基础，需要包含足够多的文本数据，涵盖客服中心可能遇到的各种话题和领域。
- 语料库的清洗和筛选。语料库中可能存在一些噪声、错误、重复或无关的数据，需要进行清洗和筛选，保证数据的干净和有效。
- 语料库的标注和增强。语料库中的数据需要进行人工或自动的标注，提供一些额外的信息，如对话意图、情感、实体等。同时，可以通过一些方法增强语料库，如数据增广、对抗生成等，增加数据的多样性和难度。

知识更新是指让大语言模型能够及时地获取和利用最新的知识，以适应客服中心不断变化的需求和环境。知识更新主要包括以下几个方面：

- 模型的持续学习。模型在训练完成后，不应该停止学习，而是需要持续地从新的数据中学习新的知识和能力，以保持模型的活力和竞争力。
- 模型的微调和适配。模型在应用到不同的客服中心时，需要根据具体的场景和任务进行微调和适配，以提高模型的针对性和效果。

- 模型的联合和融合。模型可以与其他类型的模型或系统进行联合和融合，以利用各自的优势和特点，提供更全面和更丰富的对话服务。

4.2.4 关注用户反馈和持续优化

关注用户反馈和持续优化。大语言模型虽然强大，但并不完美，它可能会产生一些错误或不合适的回答，导致用户不满或投诉。因此，客服中心应该定期收集和分析用户反馈，找出大语言模型的问题和不足，并及时进行调整和优化。例如，可以通过人工审核、在线评估、用户调查等方式获取用户反馈，然后根据反馈内容修改模型参数、增加训练数据、添加规则约束等方式改进模型性能。通过这样的过程，可以让大语言模型更好地适应客服中心的需求，提高用户满意度和忠诚度。