

# ChatGPT 调研报告 (仅供内部参考)

哈尔滨工业大学  
自然语言处理研究所 (HIT-NLP)  
2023 年 3 月 6 日

## 序言

2022 年 11 月 30 日, OpenAI 推出全新的对话式通用人工智能工具——ChatGPT。ChatGPT 表现出了非常惊艳的语言理解、生成、知识推理能力, 它可以很好地理解用户意图, 做到有效的多轮沟通, 并且回答内容完整、重点清晰、有概括、有逻辑、有条理。ChatGPT 上线后, 5 天活跃用户数高达 100 万, 2 个月活跃用户数已达 1 个亿, 成为历史上增长最快的消费者应用程序。除了被广大用户追捧外, ChatGPT 还受到了各国政府、企业界、学术界的广泛关注, 使人们看到了解决自然语言处理这一认知智能核心问题的一条可能的路径, 并被认为是向通用人工智能迈出了坚实的一步, 将对搜索引擎构成巨大的挑战, 甚至将取代很多人的工作, 更将颠覆很多领域和行业。

哈工大自然语言处理研究所组织多位老师和同学撰写了本调研报告, 从技术原理、应用场景、未来发展等方面对 ChatGPT 进行了尽量详尽的介绍及总结。

本报告仅供内部参考。

## 主要编撰人员

第一章由车万翔、杨沐昀、张伟男、赵妍妍、冯骁骋、孙承杰、李佳朋编写; 第二章由张伟男、隋典伯、高翠芸、朱庆福、李明达、王雪松编写; 第三章由刘铭、朱聪慧、汤步洲编写; 第四章由徐永东、高翠芸、朱庆福编写; 第五章由杨沐昀、张伟男、韩一、庄子彧编写; 第六章由隋典伯、高翠芸编写; 第七章由车万翔、刘铭编写。参与各章审校工作的还有: 崔一鸣、徐志明等。

报告整体由车万翔统稿。

# 目录

<b>第一章 ChatGPT 的背景与意义</b>	<b>6</b>
1.1 自然语言处理的发展历史	6
1.2 大规模预训练语言模型的技术发展历程	8
1.3 ChatGPT 技术发展历程	8
1.3.1 ChatGPT 的相关技术	10
1.3.2 ChatGPT 技术发展脉络的总结	11
1.3.3 ChatGPT 的未来技术发展方向	12
1.4 ChatGPT 的优势与劣势	13
1.4.1 ChatGPT 的优势	13
1.4.2 ChatGPT 的劣势	15
1.5 ChatGPT 的应用前景	16
1.5.1 在人工智能行业的应用前景及影响	17
1.5.2 在其他行业的应用前景及影响	17
1.6 ChatGPT 带来的风险与挑战	19
<b>第二章 ChatGPT 相关核心算法</b>	<b>24</b>
2.1 基于 Transformer 的预训练语言模型	24
2.1.1 编码预训练语言模型 (Encoder-only Pre-trained Models)	24
2.1.2 解码预训练语言模型 (Decoder-only Pre-trained Models)	25
2.1.3 基于编解码架构的预训练语言模型 (Encoder-decoder Pre-trained Models)	28
2.2 提示学习与指令精调	30
2.2.1 提示学习概述	30

## ChatGPT 调研报告

2.2.2	ChatGPT 中的指令学习 . . . . .	31
2.3	思维链 (Chain of Thought, COT) . . . . .	32
2.4	基于人类反馈的强化学习 (Reinforcement Learning with Human Feedback, RLHF) . . . . .	33
<b>第三章</b>	<b>大模型训练与部署</b>	<b>35</b>
3.1	大模型并行计算技术 . . . . .	35
3.2	并行计算框架 . . . . .	36
3.3	模型部署 . . . . .	40
3.3.1	预训练模型部署的困难 . . . . .	40
3.3.2	部署框架和部署工具 . . . . .	41
3.3.3	部署技术和优化方法 . . . . .	43
3.4	预训练模型的压缩 . . . . .	45
3.4.1	模型压缩方案概述 . . . . .	45
3.4.2	结构化模型压缩策略 . . . . .	45
3.4.3	非结构化模型压缩策略 . . . . .	46
3.4.4	模型压缩小结 . . . . .	46
<b>第四章</b>	<b>ChatGPT 相关数据集</b>	<b>48</b>
4.1	预训练数据集 . . . . .	48
4.1.1	文本预训练数据集 . . . . .	48
4.1.2	代码预训练数据集 . . . . .	50
4.2	人工标注数据规范及相关数据集 . . . . .	52
4.2.1	指令微调工作流程及数据集构建方法 . . . . .	53
4.2.2	常见的指令微调数据集 . . . . .	53
4.2.3	构建指令微调数据集的关键问题 . . . . .	54
<b>第五章</b>	<b>大模型评价方法</b>	<b>59</b>
5.1	模型评价方式 . . . . .	59
5.1.1	人工评价 . . . . .	59
5.1.2	自动评价 . . . . .	60
5.2	模型评价指标 . . . . .	62
5.2.1	准确性 . . . . .	62
5.2.2	不确定性 . . . . .	63
5.2.3	攻击性 . . . . .	63

5.2.4	毒害性 . . . . .	64
5.2.5	公平性与偏见性 . . . . .	65
5.2.6	鲁棒性 . . . . .	66
5.2.7	高效性 . . . . .	67
5.3	模型评价方法小结 . . . . .	68
<b>第六章</b>	<b>现有大模型及对话式通用人工智能系统</b>	<b>69</b>
6.1	现有大模型对比 . . . . .	69
6.2	对话式通用人工智能系统调研 . . . . .	72
6.2.1	对话式通用人工智能系统 . . . . .	72
6.2.2	不同系统之间的比较 . . . . .	75
<b>第七章</b>	<b>自然语言处理的未来发展方向</b>	<b>80</b>
7.1	提高 ChatGPT 的能力 . . . . .	80
7.2	加深对模型的认识 . . . . .	81
7.3	实际应用 . . . . .	82
7.4	从语言到 AGI 的探索之路 . . . . .	83

# 第一章 ChatGPT 的背景与意义

本章首先介绍自然语言处理、大规模预训练语言模型以及 ChatGPT 技术的发展历程，接着就 ChatGPT 的技术优点和不足进行分析，然后讨论 ChatGPT 可能的应用前景，最后展望 ChatGPT 普及后可能带来的风险与挑战。

## 1.1 自然语言处理的发展历史

人类语言（又称自然语言）具有无处不在的歧义性、高度的抽象性、近乎无穷的语义组合性和持续的进化性，理解语言往往需要具有一定的知识和推理等认知能力，这些都为计算机处理自然语言带来了巨大的挑战，使其成为机器难以逾越的鸿沟。因此，自然语言处理被认为是目前制约人工智能取得更大突破和更广泛应用的瓶颈之一，又被誉为“**人工智能皇冠上的明珠**”。国务院 2017 年印发的《新一代人工智能发展规划》将知识计算与服务、跨媒体分析推理和自然语言处理作为新一代人工智能关键共性技术体系的重要组成部分。

自然语言处理自诞生起，经历了五次研究范式的转变（如图 1.1 所示）：由最开始基于小规模专家知识的方法，逐步转向基于机器学习的方法。机器学习方法也由早期基于浅层机器学习的模型变为了基于深度学习的模型。为了解决深度学习模型需要大量标注数据的问题，2018 年开始又全面转向基于大规模预训练语言模型的方法，其突出特点是充分利用**大模型、大数据和大计算**以求更好效果。

近期，ChatGPT 表现出了非常惊艳的语言理解、生成、知识推理能力，它可以极好地理解用户意图，真正做到多轮沟通，并且回答内容完整、重点清晰、有概括、有逻辑、有条理。ChatGPT 的成功表现，使人们看到了解决自然语言处理这一认知智能核心问题的一条可能的路径，并被认为向通用人工智能迈出了坚实的一步，将对搜索引擎构成巨大的挑战，甚至将取代很

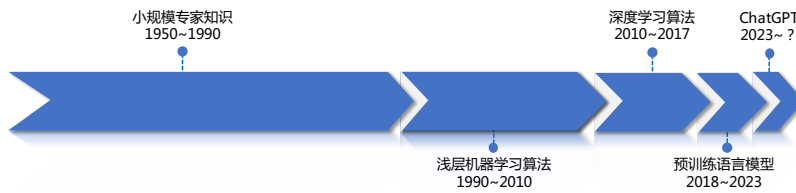


图 1.1: 自然语言处理研究范式的发展历程

多人的工作，更将颠覆很多领域和行业。

那么，ChatGPT 到底解决了什么本质科学问题，才能变得如此强大并受到广泛的关注呢？我们认为，**ChatGPT 是继数据库和搜索引擎之后的全新一代的“知识表示和调用方式”**。

知识在计算机内的表示是人工智能的核心问题。如表 1.1 所示，早期，知识以结构化的方式存储在数据库中，人类需要掌握机器语言（如 SQL），才能调用这些知识；后来，随着互联网的诞生，更多文本、图片、视频等非结构化知识存储在互联网中，人类通过关键词的方式调用搜索引擎获取知识；现在，知识以参数的形式存储在大模型中（从 2018 年开始），ChatGPT 主要解决了用自然语言直接调用这些知识的问题，这也是人类获取知识最自然的方式。

表 1.1: 知识表示和调用方式的演进

知识表示方式	表示方式的精确度	知识调用方式	调用方式的自然度	研究领域	代表应用	代表公司
关系型数据库	高	SQL	低	数据库	DBMS	Oracle、Microsoft
互联网	中	Keywords	中	信息检索	搜索引擎	Google、Microsoft
大模型	低	自然语言	高	自然语言处理	ChatGPT	OpenAI、Microsoft、Google

另外，从自然语言处理技术发展阶段的角度看（如图 1.1），可以发现一个有趣的现象，即每一个技术阶段的发展时间，大概是上一个阶段的一半。小规模专家知识发展了 40 年，浅层机器学习是 20 年，之后深度学习大概 10 年，预训练语言模型发展的时间是 5 年，那么以 ChatGPT 为代表的技

术能持续多久呢？如果大胆预测，可能是 2 到 3 年，也就是到 2025 年大概又要更新换代了。

## 1.2 大规模预训练语言模型的技术发展历程

大规模预训练语言模型（简称大模型）作为 ChatGPT 的知识表示及存储基础，对系统效果表现至关重要，接下来对大模型的技术发展历程加以简要介绍。

2018 年，OpenAI 提出了第一代 GPT（Generative Pretrained Transformer）模型<sup>[1]</sup>，将自然语言处理带入“预训练”时代。然而，GPT 模型并没有引起人们的关注，反倒是谷歌随即提出的 BERT（Bidirectional Encoder Representations from Transformers）模型<sup>[2]</sup>产生了更大的轰动。不过，OpenAI 继续沿着初代 GPT 的技术思路，陆续发布了 GPT-2<sup>[3]</sup> 和 GPT 模型 GPT-3<sup>[4]</sup>。

尤其是 GPT-3 模型，含有 1,750 亿超大规模参数，并且提出“提示语”（Prompt）的概念，只要提供具体任务的提示语，即便不对模型进行调整也可完成该任务，如：输入“我太喜欢 ChatGPT 了，这句话的情感是 \_\_\_”，那么 GPT-3 就能够直接输出结果“褒义”。如果在输入中再给一个或几个示例，那么任务完成的效果会更好，这也被称为语境学习（In-context Learning）。更详细的技术细节推荐阅读相关的综述文章<sup>[5-8]</sup>。

不过，通过对 GPT-3 模型能力的仔细评估发现，大模型并不能真正克服深度学习模型鲁棒性差、可解释性弱、推理能力缺失的问题，在深层次语义理解和生成上与人类认知水平还相去甚远。直到 ChatGPT 的问世，才彻底改变了人们对于大模型的认知。

## 1.3 ChatGPT 技术发展历程

2022 年 11 月 30 日，OpenAI 推出全新的对话式通用人工智能工具——ChatGPT。据报道，在其推出短短几天内，注册用户超过 100 万，2 个月活跃用户数已达 1 个亿，引爆全网热议，成为历史上增长最快的消费者应用程序，掀起了人工智能领域的技术巨浪。

ChatGPT 之所以有这么多个活跃用户，是因为它可以通过学习和理解人类语言，以对话的形式与人类进行交流，交互形式更为自然和精准，极大地改变了普通大众对于聊天机器人的认知，完成了从“人工智障”到“有趣”



的印象转变。除了聊天，ChatGPT 还能够根据用户提出的要求，进行机器翻译、文案撰写、代码撰写等工作。ChatGPT 拉响了大模型构建的红色警报，学界和企业界纷纷迅速跟进启动研制自己的大模型。

继 OpenAI 推出 ChatGPT 后，与之合作密切的微软迅速上线了基于 ChatGPT 类技术的 New Bing，并计划将 ChatGPT 集成到 Office 办公套件中。谷歌也迅速行动推出了类似的 Bard 与之抗衡。除此之外，苹果、亚马逊、Meta（原 Facebook）等企业也均表示要积极布局 ChatGPT 类技术。国内也有多家企业和机构明确表态正在进行类 ChatGPT 模型研发。百度表示正在基于文心大模型进行文心一言的开发，阿里巴巴表示其类 ChatGPT 产品正在研发之中，华为、腾讯表示其在大模型领域均已有相关的布局，网易表示其已经投入到类 ChatGPT 技术在教育场景的落地研发，京东表示将推出产业版 ChatGPT，科大讯飞表示将在数月后进行产品级发布，国内高校复旦大学则推出了类 ChatGPT 的 MOSS 模型。

除了国内外学界和企业界在迅速跟进以外，我国国家层面也对 ChatGPT 有所关注。2023 年 2 月 24 日，科技部部长王志刚表示：“ChatGPT 在自然语言理解、自然语言处理等方面有进步的地方，同时在算法、数据、算力上进行了有效结合。”科技部高新技术司司长陈家昌在回应 ChatGPT 相关提问时也表示，ChatGPT 最近形成了一种现象级的应用，表现出很高的人机交互水平，表现出自然语言的大模型已经具备了面向通用人工智能的一些特征，在众多行业领域有着广泛的应用潜力。<sup>1</sup>

ChatGPT 是现象级应用，标志着语言大模型已经具备了一些通用人工智能特征，在众多行业领域有着广泛的应用潜力。”这标志着在未来，ChatGPT 相关技术有可能会成为国家战略支持的重点。

从技术角度讲，ChatGPT 是一个聚焦于对话生成的大语言模型，其能够根据用户的文本描述，结合历史对话，产生相应的智能回复。其中 GPT 是英文 Generative Pretrained Transformer 的缩写。GPT 通过学习大量网络已有文本数据（如 Wikipedia, reddit 对话），获得了像人类一样流畅对话的能力。虽然 GPT 可以生成流畅的回复，但是有时候生成的回复并不符合人类的预期，OpenAI 认为符合人类预期的回复应该具有真实性、无害性和有用性。为了使生成的回复具有以上特征，OpenAI 在 2022 年初发表的工作“Training language models to follow instructions with human feedback”中提到引入人工反馈机制，并使用近端策略梯度算法（PPO）对大模型进行

---

<sup>1</sup>[https://www.sohu.com/a/645545405\\_120109837](https://www.sohu.com/a/645545405_120109837)

训练。这种基于人工反馈的训练模式能够很大程度上减小大模型生成回复与人类回复之间的偏差，也使得 ChatGPT 具有良好的表现。

### 1.3.1 ChatGPT 的相关技术

接下来将简要介绍 ChatGPT 相关技术的发展历程。ChatGPT 核心技术主要包括其具有良好的自然语言生成能力的大模型 GPT-3.5 以及训练这一模型的钥匙——基于人工反馈的强化学习（RLHF）。

GPT 家族是 OpenAI 公司推出的相关产品，这是一种生成式语言模型，可用于对话、问答、机器翻译、写代码等一系列自然语言任务。每一代 GPT 相较于上一代模型的参数量均呈现出爆炸式增长。OpenAI 在 2018 年 6 月发布的 GPT 包含 1.2 亿参数，在 2019 年 2 月发布的 GPT-2 包含 15 亿参数，在 2020 年 5 月发布的 GPT-3 包含 1750 亿参数。与相应参数量一同增长的还有公司逐年积淀下来的恐怖的数据量。可以说大规模的参数与海量的训练数据为 GPT 系列模型赋能，使其可以存储海量的知识、理解人类的自然语言并且有着良好的表达能力。

除了参数上的增长变化之外，GPT 模型家族的发展从 GPT-3 开始分成了两个技术路径并行发展<sup>2</sup>，一个路径是以 Codex 为代表的代码预训练技术，另一个路径是以 InstructGPT 为代表的文本指令（Instruction）预训练技术。但这两个技术路径不是始终并行发展的，而是到了一定阶段后（具体时间不详）进入了融合式预训练的过程，并通过指令学习（Instruction Tuning）、有监督精调（Supervised Fine-tuning）以及基于人类反馈的强化学习（Reinforcement Learning with Human Feedback, RLHF）等技术实现了以自然语言对话为接口的 ChatGPT 模型。

RLHF 这一概念最早是在 2008 年 TAMER: Training an Agent Manually via Evaluative Reinforcement<sup>[9]</sup>一文中被提及的。在传统的强化学习框架下代理（Agent）提供动作给环境，环境输出奖励和状态给代理，而在 TAMER 框架下，引入人类标注人员作为系统的额外奖励。该文章中指出引入人类进行评价的主要目的是加快模型收敛速度，降低训练成本，优化收敛方向。具体实现上，人类标注人员扮演用户和代理进行对话，产生对话样本并对回复进行排名打分，将更好的结果反馈给模型，让模型从两种反馈模式——人类评价奖励和环境奖励中学习策略，对模型进行持续迭代式微调。这一框架的提出成为后续基于 RLHF 相关工作的理论基础。

---

<sup>2</sup><https://openai.com/blog/>

欢迎访问：电子书学习和下载网站 (<https://www.shgis.com>)

文档名称：【哈尔滨工业大学】ChatGPT调研报告.pdf

请登录 <https://shgis.com/post/1663.html> 下载完整文档。

手机端请扫码查看：

