



美国人工智能 研究与发展战略规划

美国科学与技术委员会
美国网络和信息技术研究和发展小组委员会
2016 年 10 月

AI 世代 编译整理

目录

| | |
|-------------------------------------|----|
| 目录 | 2 |
| 概要 | 3 |
| 简介 | 5 |
| 美国人工智能研发战略计划的目的 | 5 |
| 期望的结果 | 6 |
| 推进人工智能国家优先级的目标 | 7 |
| 人工智能如何改善教育机会和社会健康？ | 8 |
| 人工智能的现状 | 10 |
| 研发战略 | 13 |
| 战略 1：长期投资 AI 研究 | 14 |
| 战略 2：开发有效的人类和 AI 协作方法 | 19 |
| 战略 3：了解和解决 AI 的伦理、法律和社会影响 | 22 |
| 战略 4：保证 AI 系统的安全 | 24 |
| 战略 5：开发可供 AI 培训和测试的公共共享数据集和环境 | 27 |
| 战略 6：制定标准、参照来衡量和评估 AI 技术 | 29 |
| 战略 7：更好地理解国家 AI 研发劳动力需求 | 32 |
| 建议 | 33 |
| 首字母缩略词 | 34 |

概要

作为一种转型技术，人工智能（AI）有可能带来巨大的社会经济效益。人工智能有可能变革我们生活、工作、学习、发现和沟通的方式。人工智能研究可以推进我们的国家优先级，包括促进经济繁荣，改善教育机会和生活质量，以及加强国家安全和国土安全。由于这些潜在效益，美国在过去多年中持续投资人工智能。然而，与美国联邦政府感兴趣的任何重量级技术类似，人工智能不仅带来了大量机会，在设定政府的人工智能研发方向时，也有许多问题需要考虑。

2016 年 5 月 3 日，美国政府宣布成立新的国家科学技术委员会（NSTC）机器学习和人工智能子委员会，协调联邦政府关于人工智能的活动。2016 年 6 月 15 日，这一子委员会要求网络和信息技术研发（NITRD）子委员会制定《美国人工智能研发战略计划》。随后，NITRD 成立了关于人工智能的工作组，定义美国联邦关于人工智能研发的优先级。这一优先级设定中尤其考虑了业界不太关注的领域。

《美国人工智能研发战略计划》设定了联邦政府投资人工智能研究的一系列目标，其中既包括政府主导的研究，也包括非政府组织，例如学术机构主导的研究。这些研究的最终目标是创造新的人工智能知识和技术，给社会带来全方位的效益，同时使负面影响最小化。为了实现这一目标，《人工智能研发战略计划》列出了以下优先战略：

战略 1：对人工智能研究进行长期投资。对于能驱动发现和洞察，帮助美国维持人工智能技术全球领先地位的下一代人工智能，需要优先开展投资。

战略 2：开发人类和人工智能协作的有效方法。人工智能不应取代人类，而是应当与人类合作，帮助人类取得更好的表现。研究工作需要创造人类和人工智能系统之间有效的互动。

战略 3：理解并处理人工智能的道德、法律和社会影响。我们预计，人工智能技术的行为方式应当符合人类的正式或非正式准则。研究工作需要理解并处理人工智能的道德、法律和社会影响，探索方法去设计人工智能系统，使其符合道德、法律和社会目标。

战略 4：确保人工智能系统的实体安全 and 信息安全。在人工智能大规模应用之前，我们需要确保系统运行的实体安全 and 信息安全。人工智能系统应当受控、具备良好定义，以及被正确地理解。进一步的研究需要开发出稳定、可靠、可信的人工智能系统。

战略 5 :开发用于人工智能训练和测试的共享公共数据集和环境。训练数据集和资源的深度、质量和准确性将极大地影响人工智能的表现。研究者需要开发高质量的数据集和环境,使人工智能可以以负责任的方式访问高质量数据集,以及测试和训练资源。

战略 6 :通过标准化方法去衡量人工智能技术。人工智能发展中必要的一部分是标准、评分、测试平台,以及社区参与,从而衡量人工智能的进步程度。对于开发普适的评估技术,额外研究将是必须的。

战略 7 :更好地理解美国全国人工智能研发工作者的需求。人工智能的进步需要强有力的人工智能研究者社区。更好地理解当前和未来人工智能研发工作者的需求是必要的,这将有助于保证我们有足够的人工智能专家,去面对本计划中列出的研发战略领域。

《人工智能研发战略规划》最后提出了两方面建议:

建议 1 :根据本计划中的战略 1 至战略 6 ,开发人工智能研发执行框架,识别科技机会,支持高效的人工智能研发投资合作。

建议 2 :根据本计划中的战略 7 ,研究美国的行业形势,培育并维持健康的人工智能研发工作者团队。

简介

美国人工智能研发战略计划的目的

1956 年，美国全国计算机科学的研究员聚集在新汉普歇尔的达特茅斯学院，讨论人工智能这一新的计算科学分支。在他们设想的世界中，“机器利用语言进行抽象并提出概念，解决当前只有人类可以解决的问题，并自我优化”。这一历史性会议启动了政府和行业对人工智能的数十年研究，包括感知、自动推理/计划、认知系统、机器学习、自然语言处理、机器人和相关领域的发展。目前，这些研究成果创造了新的经济领域，并影响了我们的日常生活，无论是地图技术、手机中的语音助手、手写识别、金融交易、智能物流、垃圾邮件过滤，还是语言翻译。在精确医疗、环境可持续发展、教育和公共福利等领域，人工智能的进步给我们社会的发展带来了巨大帮助。

过去 25 年，人工智能的重要性正在提升。这在很大程度上是由于统计和概率方法得到普及，大量数据可供利用，以及计算机处理能力的提升。过去 10 年，作为人工智能的子领域，机器学习已被证明可以带来更准确的结果，这使得人工智能的短期内发展更乐观。尽管近期的研究重点主要基于统计学方法，例如深度学习，但人工智能的影响力也在其他领域扩大，例如感知、自然语言处理、正式逻辑、知识呈现、机器人、控制论、认知系统架构、搜索和优化技术，以及许多其他方面。

人工智能近期的成就带来了重要问题，包括最终发展方向，以及这些技术的意义。在当前的人工智能技术中，最重要的技术缺陷是什么？新的人工智能进步能带来了哪些积极的、急需的经济和社会影响？人工智能技术如何被安全有益地利用？要如何设计人工智能系统，使其符合道德、法律和社会原则？这些发展对于人工智能开发者来说有何意义？

人工智能技术研发的形势正越来越复杂。美国政府以往和当前的投资带来了人工智能的突破性方法，其他领域，包括多个行业和非营利组织，也给人工智能带来了重要贡献。这种投资形势带来了问题，即美国联邦在人工智能技术开发中的投资是否恰当。联邦政府对人工智能投资的优先级应当是什么，尤其在企业不太可能投资的领域中？能推动美国优先级的行业和国际研发合作是否能带来机会？

2015 年，通过未分类研发投资项目，美国政府对人工智能相关技术的投资约为 11 亿美元。尽管这些投资带来了重要的新科技，但联邦政府内部还有机会展开更合理的协调，让这些投资发挥全部潜力。

在意识到人工智能带来的转型效果之后，白宫科技政策办公室（OSTP）于 2016 年 5 月宣布成立新的跨部门工作组，研究人工智能带来的帮助和风险。OSTP 还宣布将召开 4 次研讨会，时间为 2016 年 5 月至 7 月，而目的则是促进公众关于人工智能的讨论，研究人工智能带来的挑战和机遇。这些研讨会的结果成为了报告《准备人工智能的未来》的一部分。这份报告与本计划同时发布。

2016 年 6 月，美国国家科学技术委员会（NSTC）的机器学习和人工智能子委员会安排网络和信息技术研发项目（NIRD）全国协调办公室（NCO）去制作《美国人工智能研发战略计划》。该子委员会表示，这一计划应当传达明确的研发优先级，分析战略科研目标，专注于在行业不太可能投资的领域提供联邦投资，以及研究拓展并维护人工智能研发人才的必要性。

这一人工智能研发战略计划的输入信息来自多个来源，包括联邦机构、人工智能会议上的公开讨论、投资人工智能研发的联邦机构的 OMB 数据、OSTP 主动索取的关于美国如何以最佳方式发展人工智能的信息，以及公开刊物上关于人工智能的报道。

这一计划对人工智能的未来做出了几点假设。首先，计划假定，由于政府和行业的人工智能研发投资，人工智能技术将继续快速发展，无所不在。其次，计划假定，人工智能对社会的影响领域将继续扩大，包括就业、教育、公共安全、国家安全，以及美国经济。第三，计划假定，行业对人工智能的投资将继续增长，近期的商业成果提升了对研发投资回报的认知。与此同时，该计划假定，一些重要研究领域可能无法获得来自行业的足够投资，因为围绕公共事务常常会出现投资不足的现象。最后，计划假定，在行业内、学术界和政府部门，对人工智能专业性的需求将继续增长，这将导致政府和非政府部门的人力压力。

期望的结果

这一人工智能研发战略计划关注的不仅是短期的人工智能能力，也包括人工智能给社会和世界带来的长期转型式影响。近期人工智能的发展让人们对人工智能的潜力更乐观，推动了行业强劲的增长，以及人工智能方法的商业化。然而，尽管联邦政府可以调动业界对人工智能的投资兴趣，但许多应用领域和长期研究挑战缺乏短期盈利作为投资动力，因此可能无法得到行业足够多的重视。对于周期长、风险高的研究项目，短期内针对特定部门的开发工作，以及解决民营行业不关心的重要社会问题，联邦政府是资金的主要来源。因此，联邦政府应当加强在具有社会重要性的领域的人工智能投资，例如将人工智能用于公共卫生、城市系统和智能社区、社会福利、司法、环境可持续发展、国家安全，以及加速推进人工智能知识和技术成果的长期研究。

联邦政府部门之间的研发协调将对这些技术的发展带来积极影响，并将为政策制定者提供所需的知识，以解决人工智能应用中复杂的政策挑战。相互协调的方法将帮助美国充分

发挥人工智能技术的潜力，造福于社会。

这一人工智能研发战略规划定义了高层次框架，以识别人工智能的科学和技术缺漏，跟踪为了弥补这些缺漏而做出的联邦研发投资。人工智能研发战略规划指出了人工智能在短期和长期的战略优先级，以解决重要的技术和社会挑战。不过，该计划并未给各个联邦部门定义明确的研发日程表。计划设定了执行分支的目标。在这一框架内，各部门可以根据自身的使命、能力、主管范围和预算来设定优先级。因此，整体研发状况将与人工智能研发战略规划保持一致。

人工智能研发战略规划未讨论对人工智能技术的研究和利用政策，也没有探讨关于人工智能对就业和经济影响的广泛关切。尽管这些议题对美国至关重要，但将在经济顾问委员会的另一份报告中得到研究。本报告只专注于所需的研发投资，协助定义及推进政策制定，确保人工智能得到负责、安全、有益的应用。

推进人工智能国家优先级的目标

推进这一人工智能研发战略规划的目标是实现理想中的未来世界，即人工智能可以得到安全的利用，给社会所有成员带来巨大利益。人工智能的进一步发展将给几乎所有社会部门带来帮助，并推进国家优先级的发展，包括让经济更繁荣，改善生活质量，以及增强国家安全。可能的具体范例包括：

让经济更繁荣：新产品和服务将创造新市场，优化多个行业现有产品和服务的质量和效率。通过专家决策系统，我们可以创造出更高效的物流和供应链。通过基于视觉的司机辅助，以及自动驾驶/机器人系统，产品可以更高效地运输。通过新方法去控制制造流程，调度工作流，产品制造也将得到改进。

经济繁荣具体如何实现？

- **制造：**技术进步将引领制造业新的工业革命，覆盖工程产品的整个生命周期。对机器人更多地使用将使制造业岗位重新流回到美国。通过更可靠的需求预测，更灵活的运营和供应链，以及对制造工艺调整所带来影响更好的预测，人工智能可以加速生产力的提升。人工智能可以带来更智能、速度更快、更廉价、对环境更友好的生产工艺，提高工人效率，优化产品质量，降低成本，改善工人的健康和安全。机器学习算法可以优化对制造工艺的调度，减少库存需求。此外，商用的 3D 打印将给消费者带来帮助。

- **物流：**通过自适应调度和路线管理，民营的制造商和物流提供商可以利用人工智能去优化供应链管理。人工智能可以根据天气、交通和不可预见事件来进行自动调节，从而颠覆当前的供应链管理。

- **金融**：行业和政府可以在多个量级利用人工智能对异常的金融风险进行早期监测。通过安全控制，金融系统的自动化能减少恶意行为，例如操控市场、欺诈和异常交易的机会。人工智能可以提升效率，降低波动性和交易成本，同时预防系统性问题，例如价格泡沫或信贷风险被低估。

- **交通**：人工智能可以改善各种形式的交通方式，给安全性带来积极帮助。人工智能可以用于结构化的健康监控，管理基础设施资产，增强公众信任度，减少修理和重建的成本。通过加强对环境的感知，人工智能可以优化载客和载货车辆的安全性，向司机和其他旅客提供实时路线信息。人工智能应用还能优化网络级的交通管理，减少系统整体能耗和交通工具的排放。

- **农业**：人工智能系统可以带来新方法，发展可持续的农业，使农产品的生产、加工、储存、分销和消费更智能。人工智能和机器人可以收集特定位置关于农作物的实时数据，并只在必要的时间和地点提供养护（包括浇水、施肥和撒农药），补充农业急缺的劳动力。

- **营销**：基于人工智能的方法可以帮助商业实体更好地匹配供需，提高收入，为民营经济的持续发展提供资金。人工智能可以预测并识别用户需求，帮助他们更好地发现所需的产品和服务，同时降低成本。

- **通信**：人工智能技术可以最大化带宽利用效率，并推动信息存储的自动化。人工智能可以优化数字通信的过滤、搜索、翻译和总结，给商务和人们的生活方式带来积极帮助。

- **科技**：人工智能系统可以协助科学家和工程师阅读文献和专利，使理论与观察更一致，形成可测试的假设，通过机器人系统和模拟去完成实验，开发新的设备和软件。

教育机会和生活质量的改善：虚拟教师可以帮助人们终身学习，针对个体所遇到的挑战订制个性化学习计划，根据人们的兴趣、能力和教育需求让所有人参与其中。利用针对每个个体的个性化健康信息，人们可以更健康地生活，同时更热爱运动。智能家居和个人虚拟助手可以节约人们的时间，减少重复性劳动浪费的时间。

人工智能如何改善教育机会和社会健康？

- **教育**：由人工智能增强的学习机构将广泛普及，自动化教学将帮助学生的发展。人工智能教师可以成为人工教师的补充，专注于更高级、矫正性的学习过程。人工智能工具可以促进终身学习，帮助社会所有成员获得新技能。

- **医疗**：人工智能可以支撑生物信息系统，识别大规模基因研究中的一般风险，预测新药的安全性和效果。人工智能技术可以协助评估多个维度的数据，研究公共安全问题，提供决策支持系统，用于医疗诊断和治疗。人工智能技术可以用于为患者个人开发订制药剂。

最终，这将提升医疗的效果、病人的满意度，并减少浪费。

- **法律**：机器将可以用于分析法律案件的判例历史。流程复杂度的提升将带来更丰富的分析结果，协助信息发现过程。法律信息发现工具可以识别并总结相关证据。在复杂度达到一定水平之后，这些系统甚至可以用于形成法律论据。

- **个人服务**：人工智能软件可以利用多个来源的知识，提供更准确的信息，用于多方面用途。自然语言系统可以为真实世界、复杂环境中的技术系统提供更直观的界面。个性化工具可以带来自动化助手，用于个人和群体的日程安排。来自多个搜索结果的文字信息将可以自动得到总结，并通过多媒体进行增强。人工智能可以提供同声传译服务。

加强国家和国土安全：机器学习软件可以处理大量情报数据，对于策略快速变化的对手识别其模式。这些软件可以为容易受到攻击的关键基础设施、重要经济部门提供保护。数字防御系统可以极大地降低战场风险和伤亡。

国家和国土安全的加强如何实现？

- **信息安全和司法**：通过利用模式识别去探测个体行动者的异常行为，或是危险的群体行为，司法和信息安全官员可以发展更安全的社会。情报感知系统可以保护关键的基础设施，例如机场和发电厂。

- **安全和预测**：无论是自然原因还是人为原因引起，当基础设施受严重干扰的可能性明显上升时，分布式传感器系统，以及对通常条件的模式理解将有能力及时做出探测。这种预测能力可以协助确定，问题会发生在哪里，随后通过相应举措去解决问题，甚至在问题发生前予以解决。

然而，如果希望人工智能向积极的方向发展，那么还需要大力的研发。在人工智能的子领域，包括基础科学和应用科学领域，仍存在许多关键而困难的技术挑战。人工智能技术同时也带来了风险，例如可能改变劳动力市场，人类有可能被自动化系统增强，但也可能会被替代。此外，人工智能系统的安全性和可靠性也存在不确定。这一人工智能研发战略计划的随后章节讨论了有助于实现这些目标，减小可能风险的人工智能研发投资的高优先级、战略领域。

人工智能的现状

自诞生以来，人工智能研究经历了 3 次技术浪潮。第一次浪潮专注于人类知识。80 年代时，人工智能专注于在定义明确的领域发展基于规则的专家系统。人工智能的知识来自于人类专家，并通过“if-then”逻辑去表达，随后集成在硬件中。这些系统被成功用于解决某些严格定义的问题，然而这种系统无法处理不确定性。无论如何，这类系统引领了重要的解决方案，而这方面的技术开发目前仍然很活跃。

人工智能研究的第二次浪潮开始于 00 年代，并一直延续到现在。其标志是机器学习的兴起。数据量的增大，并行计算能力的成本降低，以及优化的机器学习技术带来了更强大的人工智能，而人工智能被用于了图像和手写识别、语音理解，以及人类语言翻译等任务。这些技术发展的成果随处可见：智能手机提供语音识别，ATM 机提供对手写支票的识别，电子邮件应用能过滤垃圾邮件，免费在线服务能完成机器学习。这些成功的关键在于深度学习的发展。

目前，人工智能系统常常能在特定任务方面胜过人类，例如象棋（1997 年）、雅达利游戏（2013 年）、图像识别（2015 年）、语音识别（2015 年），以及围棋（2016 年）。这样的发展似乎正在加快。性能最强大的系统基于机器学习方法，而不是一套人工编码的规则。

人工智能的这些成就受到了一系列基础研究的推动。这些研究正在拓展，很可能促进未来的发展。作为指标之一，从 2013 年至 2015 年，Web of Science 编目的期刊中提到关键词“深度学习”的次数增加到了 6 倍。这一趋势反映了全球研究工作的特点。而在这一方面，美国的论文发表篇数，以及被引用过至少一次的论文篇数已经不再是全球领先。

美国政府在人工智能的研究中扮演了关键角色，但商业界在人工智能研发方面也很活跃。提到“深度学习”或“深度神经网络”的专利数量正在快速增长。从 2013 年至 2014 年，风投对人工智能创业公司的投资增加到了 4 倍。人工智能应用目前也给大企业带来了明显的收入。人工智能对金融系统的影响可能更大：自动交易已占全球金融交易的约一半，交易额达到了万亿美元量级。

Deep Learning

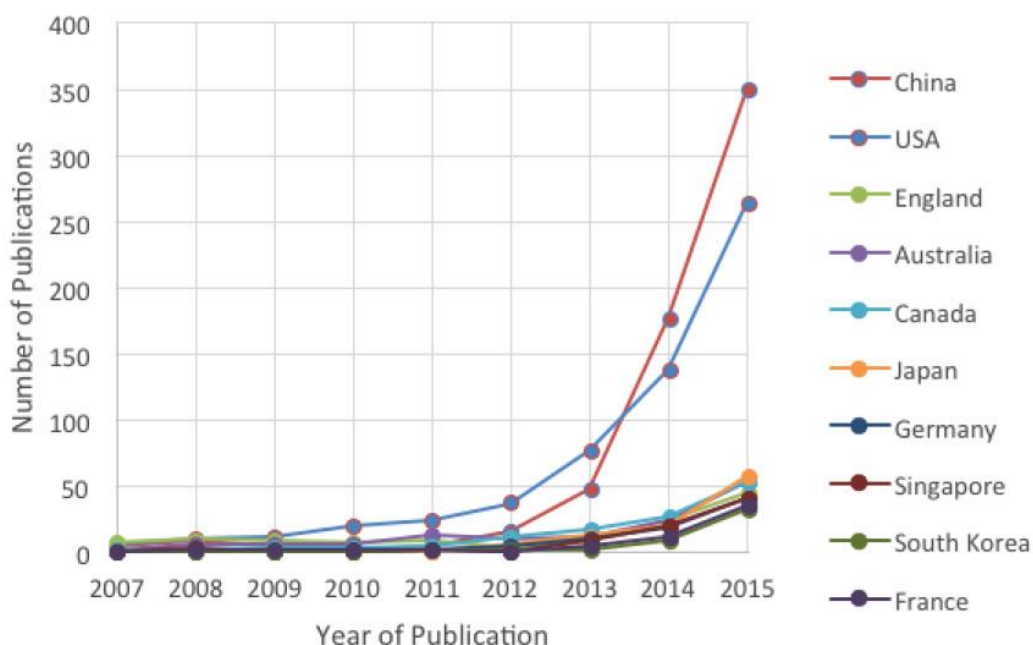


图 1：各国论文中提到“深度学习”或“深度神经网络”的篇数

Deep Learning (Cited Publications)

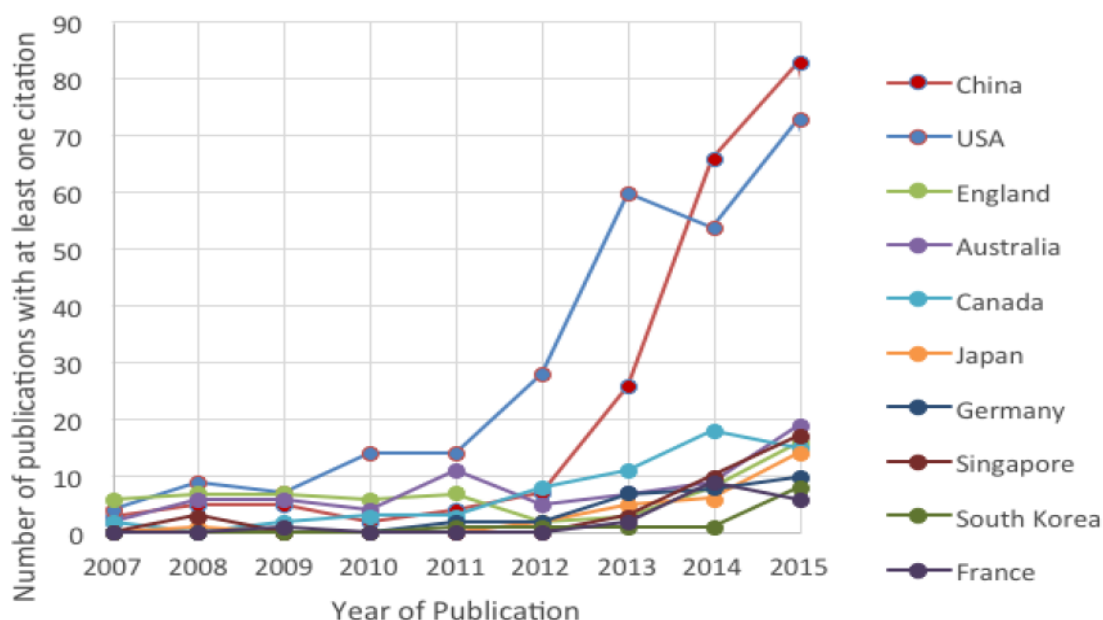


图 2：各国提到“深度学习”或“深度神经网络”的论文中，至少被引用过一次的篇数

Deep Learning in Patents

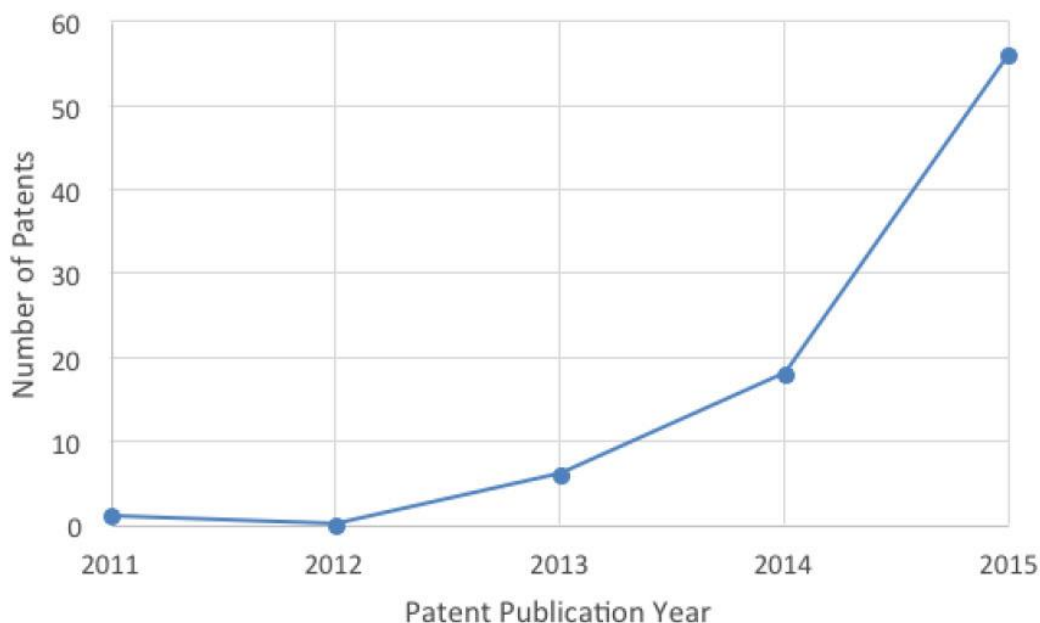


图 3：对用到关键字“深度学习”或“深度神经网络”的专利数量的分析

尽管取得了这些进展，但人工智能系统仍存在自身的局限性。几乎所有进展都发生在“狭义人工智能”领域，这样的人工智能只能完成特定任务，而“通用人工智能”，即可以胜任多种认知领域的人工智能，并未能取得太大发展。即使在“狭义人工智能”内部，发展也不平衡。用于图像识别的人工智能系统依赖于人工标记出数千个样本的正确答案。作为对比，大部分人可以通过少数几个范例完成“一点即通”的学习。大部分机器学习系统容易被复杂的场景、叠加的物体所迷惑，但即使是幼儿也能轻松完成“场景分析”。对人类来说简单的场景理解对机器来说依然困难。

人工智能目前或许正在迎来“第三次浪潮”，即专注于解释性的通用人工智能技术。这类方法的目标是通过解释和纠正界面增强学习模型，使学习基础更明确，提高输出结果的可靠性，同时加强运行的透明度，从狭义人工智能发展至有能力完成更广泛的任务。如果获得成功，那么工程师可以创造系统，利用解释模型去归类真实世界现象，参与与人类的自然交流，在遇到新任务、新情况时进行学习和推力，通过对过往经验的归纳来解决新问题。人工智能系统的这种解释模型或许可以通过先进方法来自动开发。这些模型可以帮助人工智能系统快速学习，或许可以给人工智能系统带来“理解”能力，从而使人工智能系统更通用化。

研发战略

这份《人工智能研发战略规划》中所提到的研究重点是业界不大可能重视的，因此需要联邦投资。这些研究重点设计 AI 的所有子领域，包括感知、自动推理/规划、认知系统、机器学习、自然语言处理、机器人学以及其他相关领域。人工智能涉及的范围很广，这些研发重点也涉及了整个领域，而不是只关注每一个子领域下的某些研究难点。为了让计划顺利实施，应该制定详细的线路图，指出计划实施过程中将遇到哪些能力差距。战略一部分所指出的联邦研发重点之一是，对人工智能进行长期的研究，以期有新的发现和认识。美国联邦政府多次投资高风险、高回报的基础研究，给我们带来了我们当前的生活离不开的革命性的科技进步，包括互联网、全球定位系统 GPS、智能手机语音识别、心率检测器、太阳能电池板、癌症治疗方法等。AI 将触及到社会的几乎所有层面，有可能给社会和经济带来诸多好处。因此，为了保持在这个领域的全球领先地位，美国必须将投资重点放在 AI 的基础研究上。

许多 AI 技术将跟人类肩并肩工作，因此，我们面临的挑战之一是研发出易于使用的 AI 系统。人类和 AI 系统之间的壁垒正在开始慢慢地解体，同时，AI 系统增强或者改善了人类的能力。为了找到有效的方法让人类和 AI 之间进行互动和合作，基础研究是必需的。

AI 的发展给社会带来了许多好处，还提升了美国的国家竞争力。然而，跟大部分革命性的技术一样，AI 对一些领域造成了威胁，引发了安全、伦理和法律问题，甚至对工作和经济都造成了影响。因此，在发展 AI 科学和技术的同时，联邦政府亦应该拨出资金，投资相关研究，最终研发出符合人类伦理、法律和社会秩序的 AI 系统。

目前的 AI 技术不能保证 AI 系统的安全以及可预见性，这是一个弥补的空白。保证 AI 系统的安全性是一个重大难题，因为这些系统异常复杂，且在不断变化。本计划中的若干研究重点就是旨在解决这些安全问题。首先，战略 4 强调，我们要研发出受用户信任、行为被用户接受、能按照用户意愿行动的透明系统。AI 系统拥有巨大潜能，但也十分复杂，为了保证其在跟人类和环境互动中的安全性，必须提高 AI 技术的安全性和可控制性，因此，投资相关研究就非常必要了。战略 5 呼吁联邦政府投资用于 AI 训练和测试的共享公共数据集，以推动 AI 研究，并有效地对备选方案进行比较。战略 6 讨论了标准和基准问题。标准和基准对衡量和评估 AI 系统来说是不可或缺的。

最后，AI 技术越来越普及，给 AI 研发专家带来了新的压力。对那些深刻了解 AI 技术，能提出新的理论充实该领域知识的 AI 科学家和工程师来说，机会是无限的。美国政府应该采取行动保证有足够多的 AI 人才。战略 7 讨论了这方面的挑战。

图 4 用图表说明了这份《人工智能研发战略规划》的总体组织。底部（红色部分）是横贯一切的底层基础，影响了所有 AI 系统的研发。这些基础会在战略 3-7 详细讲述。更高的层面（浅中深三种蓝色）涵盖了推动 AI 所需要的众多研究领域。这些基础研究领域将在战略 1-2 讲述。在这幅图的最顶端（深蓝部分）是将受益于 AI 发展的应用领域。整个《人工智能研发战略规划》阐释了联邦投资 AI 的高级框架。

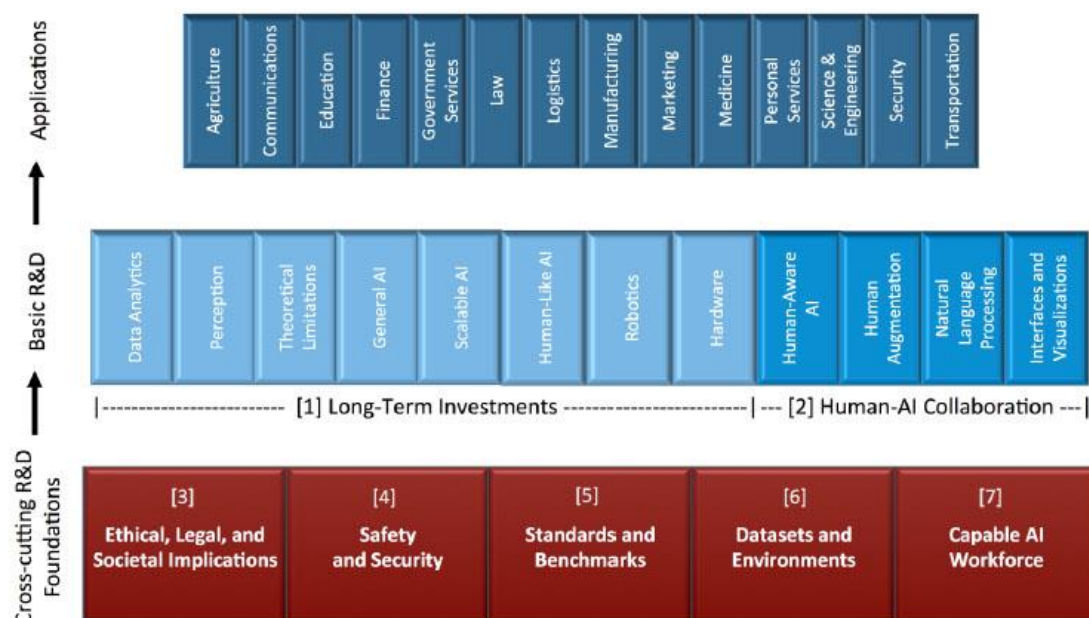


图 4：人工智能研发战略规划的组织。一系列基础性研究（最下方红色的部分）对所有人工智能研究来说都很重要。其他基本的人工智能研发领域（中间一行蓝色的部分）在此基础上进行，并将对广泛的社会领域产生影响（最上方深蓝色的部分）。（括号中的数字是这一领域在本计划中的战略编号。编号不代表重要性。）

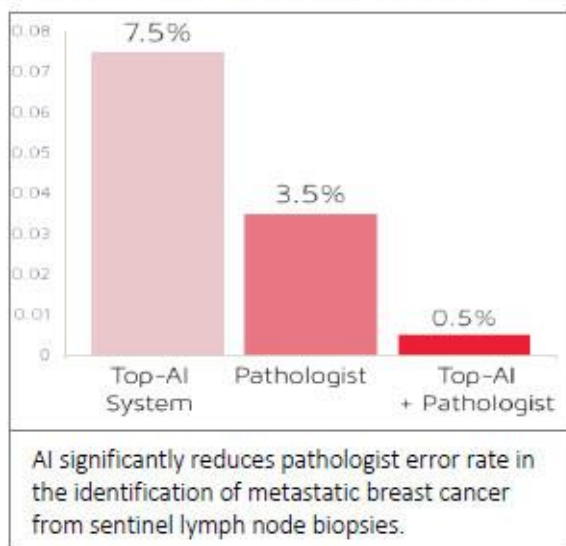
战略 1：长期投资 AI 研究

AI 研究投资需面向可获得长期回报的领域。渐进式研究容易出结果，也是长期研究的一个重要组成部分，但长期对高风险研究进行连续投资带来的回报可能会更高。这类回报可能需要 5 年、10 年，甚至更长时间才能看到。美国国家研究委员会最近的一份报告强调，联邦投资对长期研究起了关键性的作用。该报告指出，“在初步探索和商业应用之间是漫长且不可预测的孵化期，需要持续的研究和资金，”“从初步概念到成功进入市场通常需要数十年的时间。”漫长的基础研究最终带来高回报的例子包括万维网和深度学习。这两个领域的基础研究都是始于上世纪六十年代，经过研究人员 30 多年的努力之后，才从理论成为现实，改变了世界。

National Institutes of Health (NIH) grants-supported research

ARTIFICIAL INTELLIGENCE FOR COMPUTATIONAL PATHOLOGY

Image interpretation plays a central role in the pathologic diagnosis of cancer. Since the late 19th century, the primary tool used by pathologists to make definitive cancer diagnoses is the microscope. Pathologists diagnose cancer by manually examining stained sections of cancer tissues to determine the cancer subtype. Pathologic diagnosis using conventional methods is labor-intensive with poor reproducibility and quality concerns. New approaches use fundamental AI research to build tools to make pathologic analysis more efficient, accurate, and predictive. In the 2016 Camelyon Grand Challenge for metastatic cancer detection,⁶⁹ the top-performing entry in the competition was an AI-based computational system that achieved an error rate of 7.5%.⁷⁰ A pathologist reviewing the same set of evaluation images achieved an error rate of 3.5%. Combining the predictions of the AI system with the pathologist lowered the error rate to down to 0.5%, representing an 85% reduction in error (see image).⁷¹ This example illustrates how fundamental research in AI can drive the development



of high performing computational systems that offer great potential for making pathological diagnoses more efficient and more accurate.

以下小节突出了这些领域中的一些。战略 2 至 6 讨论了其他类别的重要人工智能研究。

通过数据发现知识

正如在《联邦大数据研究和发展战略计划》中所说的那样，实现对数据的理解和发现知识需要许多基础新工具和技术。为了识别所有隐藏在大数据中间的有用信息，我们需要开发更先进的机器学习算法。许多开放式的研究问题是为了解决数据的创建和使用，这些问题包括 AI 系统训练的准确性和适宜性。在处理海量数据的过程中，保证数据的准确尤其具有挑战性，因此，从海量数据中获得知识非常困难。许多研究通过保证数据质量来保证准确性，进行数据清理，发现知识。然而，我们需要进一步的研究以提高数据清理技术的效率，创造方法发现数据中的异常和不一致，并将人类反馈整合进来。研究人员需要探索新的方法同步挖掘数据和相关元数据。

AI 的许多应用从性质上来说跨学科的，且使用了异构数据。为了从各种类型的数据发现知识，还需要对多模态机器学习进行进一步的研究。AI 研究人员必须确定训练需要的数据量，正确处理大规模 vs 长尾数据需求。他们还必须学会如何识别和处理纯数据统计方法之外的小概率问题，以及跟知识来源和数据来源合作，将学习过程中的模型和本体整合在一起，在没有大数据来源时利用有限的的数据得出有效的结果。

提高 AI 系统的感知能力

感知是智能系统通向世界的窗户。感知来自传感器数据。这些数据经过处理和融合之后，用于提取跟 AI 系统任务相关的信息。感知数据被整合进来之后，AI 系统就对周围的情境有了了解，进而能够有效和安全地执行任务。硬件和算法的进步将能让 AI 系统受益，让其感知更准确，更可靠。传感器必须能够从更远的距离，以更高的分辨率，实时捕捉数据。感知系统应该能够整合来自各种传感器和其他来源的数据，以确定当前的状况，预测未来的状态。对物体的检测、分类和识别仍然具有很大挑战性，尤其是在纷杂和动态的情境下。此外，人类也需要使用传感器和算法提高感知能力，这样 AI 系统才能更有效地跟人类一起协作。在整个感知过程中，AI 系统还需要一个计算不确定性的框架，以提高置信水平和准确率。

理解 AI 的理论能力和极限

尽管许多 AI 算法的最终目标是像人类那样解决问题，但我们不是非常了解 AI 的理论能力和极限，也不知道 AI 算法在解决问题时能有多接近人类。为了更好地了解 AI 技术为何通常在实践中能够获得良好的效果，我们需要做一些理论上的工作。尽管不同的学科（包括数据、控制科学和计算机科学）都在研究这个问题，但该领域目前缺乏统一的理论模型或者框架，因此不能很好地理解 AI 系统的性能。我们还需要进一步研究计算的可解性，搞清楚 AI 系统在理论上可以解决以及不能解决哪类问题。为了搞清楚硬件将如何影响算法的性能，这样的研究必须利用现有的硬件。了解哪些问题在理论上是无法解决的，然后研究人员会研究出比较接近完美的解决方案，甚至开始研究新的硬件和 AI 系统。

研究通用人工智能

AI 可以分为“狭义 AI”和“通用 AI”。狭义 AI 系统能够执行专门的、定义明确的任务，比如语音识别、图像识别和翻译。最近，IBM 的 Watson 和谷歌的 AlphaGo 等狭义 AI 都取得了举世瞩目的成绩。这些系统被认为是“超越了人类”，因为它们在分别在问题游戏“Jeopardy”和围棋上打败了人类。但是，它们是专门为这些游戏设计的，如果要用这些系统解决更多的其他问题，就需要花很大精力重新设计。通用 AI 的长期目标是创造在众多认知领域上像人类智能那样灵活和全面的系统，这些认知领域包括学习、语言、感知、推理、创造能力等。多样化的学习能力将使得通用 AI 能够将一个领域的知识迁移到另一个领域，能够从自身经验和人类身上学到新东西。自从 AI 诞生以来，通用 AI 一直是研究人员的梦想，但当前的系统远远没有达到这个目标。研究人员正在研究狭义 AI 和通用 AI 之间的关系，寄望于一个领域的经验能被用于另一个领域。尽管目前尚没有达成共识，但大多数 AI 研究人员认为，通用 AI 需要几十年的时间才能问世。

研发可扩展的智能系统

AI 系统网络或者群组可以协调或者自动合作完成单一 AI 系统不可能完成的任务，它们还能跟人类协作完成任务，或者受人类指挥。这些系统具有规划、协作、控制和扩展能力。它们的规划能力必须足够快，能够实时适应环境的变化。它们还必须灵活应对通信带宽变化和系统故障。先前的研究专注于集中式规划与协调技术，然而，这些方法会受到单点故障的影响。分布式规划与控制技术很难通过算法实现，而且通常效率低，不完善，更容易出现单点故障。未来的研究必须找到更有效、更可靠的规划、控制和协作技术。

研究更接近人类的人工智能

实现更接近人类的 AI 需要有新一代的智能系统，比如智能教学系统、能够有效协助人类完成任务的智能助手。然而，当前的 AI 算法不能做到像人类那样学习和完成任务。人类能够通过几个例子学习，按照指导和提示完成任务，甚至通过观察其他人如何完成任务而掌握新技能。比如，医学院学生就是通过观察有经验的医生完成手术而学会做手术。围棋大师也只需完成数万棋局训练自己。相比之下，人类要花数百年时间才能完成训练 AlphaGo 所需的棋局。因此，要实现接近人类的 AI 需要更多的基础研究。

NSF-funded Framework on Game Theory for Security

Security is a critical concern around the world, whether it is the challenge of protecting ports, airports and other critical infrastructure; protecting endangered wildlife, forests and fisheries; suppressing urban crime; or security in cyberspace. Unfortunately, limited security resources prevent full security coverage at all times; instead, we must optimize the use of limited security resources. To that end, the "security games" framework—based on basic research in computational game theory, while also incorporating elements of human behavior modeling, AI planning under uncertainty and machine learning—has led to building and deployment of decision aids for security agencies in the United States and around the world.⁷⁴



For example, the ARMOR system has been deployed at LAX airport since 2008, the IRIS system for the Federal Air Marshals Service has been in use since 2009, and the PROTECT system for the U.S. Coast Guard since 2011. Typically, given limited security resources (e.g., boats, air marshals, police), and a large number of targets of different values (e.g., different flights, different terminals at an airport), security-games-based decision aids provide a

randomized allocation or patrolling schedule that takes into account the weights of different targets and intelligent reaction of the adversary to the different security postures. These applications have been shown to provide a significant improvement in performance of the different security agencies using a variety of metrics, e.g., capture rates, red teams, patrol schedule randomness, and others.⁷⁴

研发功能更强大、更可靠的机器人

机器人技术在过去十年有了重大进展，已经被应用于制造业、物流业、医学、国防、国家安全、农业和消费产品。在人们的想象中，机器人都是静静地待在工厂里，完成各种机械动作。机器人技术的最近发展使得机器人能够跟人类紧密合作。机器人技术将来能够增强、改善人类的体力和智力。然而，科学家首先需要让这些机器人系统变得更强大、更可靠、更易用。

机器人需要提高感知能力，这样才能解码各种传感器的信息，实时理解周围环境。同时机器人还需要有更好的认知和推理能力，以更好地理解物理世界，并与之互动。适应和学习能力的提高将能让机器人进行总结和自我评估，并从人类老师学到各种动作。为了让机器人能在不平坦的地形行走，避开各种障碍物，我们需要继续研究如何让机器人获得最佳的移动能力和控制能力。机器人之间还需要学习无缝协作。

推动 AI 硬件的研发

人们通常会把 AI 研究跟软件联系起来，但 AI 系统的性能也十分依赖硬件。当前深度学习复兴跟 GPU 硬件技术的发展有着直接的联系。研发针对 AI 算法优化的硬件将带来比 GPU 更高的性能。一个例子就是“神经形态”处理器。

改善硬件也能提升 AI 方法的性能。我们需要进一步研究在整个分布式系统中以受控方式打开和关闭数据管道的方法，还需要继续研究如何让机器学习算法高效地从高速数据学习，包括从多个数据流水线同时学习的分布式机器学习算法。更先进的基于机器学习的反馈方法让人工智能系统智能采样或优先考虑来自大规模模拟、实验仪器和分布式传感器系统的数据，比如智能建筑和物联网。这样的方法可能需要动态 I/O 决策，在决策的同时实时根据重要性决定是否存储数据，而不是简单地按照固定频率存储数据。

为改善的硬件创造 AI

虽然硬件的改善可以带来功能更强大的 AI 系统，AI 系统也可以提升硬件的性能。这种互利关系让硬件性能进一步提高，因为解决计算物理限制需要新的硬件设计。基于 AI 的方法对于改进高性能计算（HPC）系统的操作可能尤其重要。这样的系统消耗大量的能量。AI 被用于预测 HPC 的性能和资源使用，并进行在线优化决策以提高效率；更高级的 AI 技术可以进一步提高系统性能。AI 还可用于创建可自动重新配置的 HPC 系统，可在无人干预的情况下处理系统故障。

改进的 AI 算法可以通过减少处理器和存储器之间的数据移动来提高多核系统的性能。数据移动是百亿亿次计算系统的主要障碍。在实践中，HPC 系统的配置执行永远不会相同，并且不同的应用会同时执行，每个不同的软件代码的状态都会随时间独自发生变化。AI 算法需要设计得能为 HPC 系统在线和大规模运行。

战略 2：开发有效的人类和 AI 协作方法

虽然完全自主的 AI 系统在一些应用领域（例如水下或深层空间探索）非常重要，但在许多其他应用领域（例如灾难恢复和医学诊断），由人类和 AI 系统一起合作完成应用目标是最有效的解决方案。这种协作式的互动利用了人类和人工智能系统的互补性。虽然我们已经找到有效的方法让人类与 AI 互相协作，其中这些方法大多数是“点解决方案”，只能在特定环境中使用特定的平台来实现特定的目标。为每个可能的应用实例创造点解决方案不是可行方法，因此，我们应该需要做更多的工作来研究人类与 AI 合作的通用方法，并在特定和通用之间作出权衡。

未来的应用将因人类与 AI 系统之间的互动的性质，人类和 AI 系统的数量，以及人类和 AI 系统如何沟通和共享实时信息的不同而产生巨大差别。人类和 AI 系统之间的功能性作用通常可以分成以下几个类别：

1、AI 与人类一起执行功能：AI 系统替人类决策者执行外围任务。例如，AI 可以替人类进行工作记忆，进行短期或长期记忆检索，以及预测任务。

2、当人类遇到认知超载时，AI 帮忙执行任务：AI 系统可以实现复杂的监测功能（如飞机的近地警告系统），作出决策，当人类需要帮助时自动进行医疗诊断。

3、AI 代替人类执行任务：AI 系统可以执行人类因能力限制而不能完成的任务，例如用于复杂的数学运算。

实现人类与 AI 系统之间有效的交互需要额外的研发，以确保系统设计不会过度复杂。培训和经验可以让人更熟悉 AI 系统，确保人类对 AI 系统的能力有很好的了解，知道 AI 系统能做什么，不能做什么。为了解决这些问题，应当使用某些以人为中心的自动化原理设计和开发这些系统：

- 1、采用直观、用户友好的人机界面、控制和显示设计。
- 2、让操作者随时了解情况变化。向他们显示关键信息、AI 系统的状态以及这些状态的变化。
- 3、对操作员进行培训。反复培训一般知识、技能和能力（KSA），以及 AI 系统采用的算法和逻辑、系统可能遇到的故障。

4、使自动化变得更灵活。部署 AI 系统应被视为一种设计选择，操作员可以自己决定是否要使用 AI 系统。同样重要的是自适应 AI 系统的设计和部署，人类操作员可能在工作量过大或身体疲劳时选择使用这些系统。

研究人员在研发能有效跟人类协作的 AI 系统时会遇到许多根本性的挑战。以下的小节中列出了其中的一些重要挑战。

寻找新算法研发能感知人类存在的 AI

多年来，AI 算法已经能够解决日益复杂的问题。然而，这些算法的能力和这些系统的可用性之间存在着差距。我们需要能跟用户直观地进行互动，实现无缝人机合作的智能系统。直观交互包括浅交互，如当用户拒绝系统推荐的选项时；记录用户过去动作的基于模型的方法；或甚至基于准确的人类认知模型的深层模型。我们需要开发仅在必要和合适情况下打断人类的智能系统。智能系统也应该能够增强人类的认知，知道用户何时需要检索哪些信息，即使他们没有明确地要求系统提供信息。未来的智能系统必须能够懂得人类社会规范，并能按照这些规范进行运作。如果智能系统拥有情商，它们将能更有效地跟人类进行合作，因为这样的系统能够自动识别用户的情绪，并作出适当的反应。另一个研究目标是超越一个人对应一台机器的交互模式，研究出“系统的系统”，即多个机器组成的团队与多人同时进行交互。

人与 AI 系统的互动具有广泛的目标。AI 系统必须有理解多个目标，并能采取行动达到这些目标，还要了解这些行动遇到的束缚，以及能够适应目标的变化。此外，人类和人工智能系统必须有共同的目标，并相互理解。

开发增强人类能力的 AI 技术

虽然先前的 AI 研究重心是在执行特定任务时能跟人类并肩甚至超越人类的算法，但我们仍需要做更多的工作来开发能在静止设备（如计算机）、可穿戴设备（如智能眼镜）、植入设备（如脑接口）和在特定用户环境（例如特制的手术室）中增强人类能力的算法。例如，在认识被增强之后，医疗助理能够依靠多部设备提供的数据，指出治疗程序中的错误。其他系统可以通过帮助用户回忆过去的经历，以解决当前的问题。

人类和 AI 系统之间的另一种类型的协作需要主动学习才能实现。在主动学习中，输入来自领域专家，只有在学习算法不确定时，才会根据数据进行学习。主动学习迄今为止只在监督学习中使用，我们需要进一步的研究，才能将主动学习纳入到无监督学习中。

开发可视化和 AI 人机界面

更好的可视化和用户界面也是一个需要大力研究的领域，这能帮助人们理解大量现代数据集和来自多种来源的信息。可视化和用户界面必须清楚地呈现越来越复杂的数据和信息，而且能够为人类所理解。在对安全要求非常高的操作中，提供实时结果非常重要。在这些情况下，用户需要可以实时作出回应并提供正确信息的可视化和用户界面。

人和 AI 的协作可以应用到各种环境，甚至通信受限的场合也能用到。在一些领域，人类和 AI 的通信延迟比较低，人类与 AI 的通信非常快速且可靠。在其他领域，人与 AI 系统之间的远程通信有很高的延迟（例如，地球和火星之间的往返时间为 5-20 分钟），因此需要部署能够自主运作的平台，只有高层次的战略目标才被传达给平台。

研发更有效的语言处理系统

让人们通过口语和书面语跟 AI 进行互动一直是 AI 研究人员希望达成的目标。虽然这方面的研究已经取得了重大进展，但是语言处理研究仍然需要解决很多问题，才能让人类与 AI 的沟通跟与其他人的沟通一样有效率。语言处理方面最近取得的大部分进展归功于使用了数据驱动的机器学习方法，这种方法带来了成功的系统。例如，这些系统已经能在安静的环境中实时识别流利的英语演讲。然而，这些成就只是实现更长期目标中所迈出的第一步。当前的系统不能应对现实世界的挑战，比如它不能理解在嘈杂环境中的演讲，方音严重的话语，儿童的话语，语言障碍者的话语和手语。我们还需要研发出能够理解实时对话的语音处理系统。这样的系统将需要推断其人类对话者的目标和意图，根据不同的情景使用适当的语体、文体和修辞，并在对话产生误解的情况下使用修复策略。我们还需要进一步研究如何开发出能更容易地在不同的语言之间进行总结概括的系统。此外，我们需要进行更多的研究以获得有用的结构化领域知识，而且这种知识要以语言处理系统易于访问的形式存在。

DARPA's Personalized Assistant that Learns (PAL) Program Created the Technology that Apple Commercialized as Siri

Computing technology is critical to every aspect of modern life, but the information systems we use daily lack the general, flexible abilities of human cognition. In the Personalized Assistant that Learns (PAL) program,⁸⁵ DARPA set about to create cognitive assistants that can learn from experience, reason, and be told what to do via a speech interface. DARPA envisioned PAL technologies making information systems more efficient and effective for users. DARPA and the PAL performers worked with military operators to apply PAL technologies to problems of command and control, and PAL procedure learning technology was integrated in the U.S. Army's Command Post of the Future version Battle Command 10 (see figure) and used around the world.



DARPA was also acutely aware of the commercial potential of the PAL technology, especially for mobile applications where speech-based smartphone interaction would be required. DARPA strongly encouraged PAL commercialization and in 2007, in response to DARPA's encouragement, Siri Inc. was created to commercialize PAL technology in a system that could assist a user by managing information and automating tasks through a speech-based interface. In April 2010, Siri Inc. was acquired by Apple, which further developed the technologies to make them an integral part—and the defining feature—of Apple's mobile operating system available on the iPhone and iPad.

为了让人类与 AI 系统之间的交互更自然和直观，我们对许多其他领域的语言处理也需要获得进步。我们必须为口语和书面语模式建立起鲁棒性强的计算模型，用于确定说话者的情绪状态和立场、隐含在语音和文本中的信息。我们还需要研发出新的语言处理技术，让在物理世界中运作的 AI 系统在一定的情境中拥有共同的语言。最后，由于人们在线沟通的方式跟语言交互具有很大差异，在这些情境中使用的语言模型必须完善，这样，社交 AI 系统才能够跟人类更有效地互动。

战略 3：了解和解决 AI 的伦理、法律和社会影响

当 AI 代理人程序自主行动时，我们期望它们能够跟人类一样遵守正式和非正式社会规范，遵守法律和道德。我们需要了解 AI 对伦理、法律和社会的影响，并研究出方法让 AI 设计符合道德、法律和社会原则。我们还必须考虑隐私问题；关于这一问题的更多信息可以在《国家隐私研究战略》中找到。

与任何技术一样，AI 的用途必须符合法律和道德的原则，才能被广泛接受。我们面临的挑战是如何将这些原则应用于新技术，特别是那些涉及自主性、代理和控制的技术。

正如《稳健和有益人工智能的研究重点》(Research Priorities for Robust and Beneficial Artificial Intelligence) 所说的那样：

“为了创建表现稳健的系统，我们当然需要确定，在每个应用领域中，什么才是良好表现。这一伦理维度跟现在可用的工程技术及这些技术的可靠性，以及我们如何权衡取舍紧密相关——计算机科学、机器学习和 AI 的所有领域的专业知识都会产生影响。” 我们需要在 NITRD 相关 IT 领域（即在信息技术以及上述学科中）内部和外部进行进一步的研究，以了解对 AI 系统的研发和运用对社会造成的影响。以下小节探讨该领域 IT 研究面临的关键挑战。

改善公平性，透明度和设计问责制

许多人担心，数据密集型 AI 算法容易出错和被滥用，并可能对性别、年龄、种族或经济阶级造成不良后果。就这方面来说，正确地为 AI 系统收集和使用数据将是一个重大挑战。然而，除了纯粹跟数据相关的问题之外，更大的问题是如何设计出先天就公正、公平、透明和负责的 AI。研究人员必须学会设计出行动和决策透明、并且易于被人类理解的 AI 系统，这样我们才能排查出它们可能拥有的偏见，而不是让它们学习和重复这些偏见。如何体现和“编码”价值和信念系统的问题也应该严肃对待。科学家也必须研究，在多大程度上将正义和公平设计到系统中，以及如何在当前工程技术的条件下实现这一点。

创建符合伦理的 AI

除了公平和正义之外，我们需要关心的另一个问题是，AI 系统的行为能否遵守一般的伦理原则。AI 的进步将能如何处理跟“机器相关”的道德问题，或者 AI 的什么用途可能被认为是不道德的？伦理在本质上属于哲学问题，而 AI 技术依赖于，并受限于工程设计。因此，研究人员必须在技术可行的范围内，努力发展符合现有法律、社会规范和伦理的算法和架构。这显然是一项非常具有挑战性的任务。伦理原则通常具有不同程度的模糊性，很难转化为精确的系统和算法设计。AI 系统，特别是运用了新型自主决策算法的系统，面对各自独立，且可能冲突的价值体系时可能会面临道德困境。另外，不同的文化、宗教和信仰具有不同的伦理观念。然而，我们可以研发可接受的伦理参考框架来指导 AI 系统的推理和决策，解释和证明其结论和行动的正确性。我们需要采用一种多学科方法来生成用于训练的数据集。这些数据集应该反映正确的价值体系，里面应包含在遇到道德困境和互相冲突的价值观时如何做出选择的实例。这些例子可以包括法律或道德的“角落案件”，其结果和评判需要对用户透明。

设计符合道德的 AI 架构

我们必须在基础研究方面取得进一步进展，设计出将伦理推断整合到 AI 系统中的最佳架构。现在已经出现了多种解决方案，例如使用双层监控结构将操作层 AI 与负责评估伦理和法律问题的监视程序分开。另一种意见是安全工程是首选，用 AI 代理架构的准确概念框架保证 AI 不会伤害人类。第三种方法是制定伦理架构，并结合对 AI 系统行为的逻辑约束。随着 AI 系统变得越来越通用，它们的架构将可能包含负责问题的子系统。研究人员应该专注于如何最好地解决 AI 系统的整体设计，以实现伦理、法律和社会目标。

战略 4：保证 AI 系统的安全

要想将 AI 系统广泛地应用，还需要保证这个系统可以在有效地控制之下安全地运行。要想打造可靠、可信任的 AI 系统，对其的安全性的研究势在必行，AI 系统和其他任何复杂的系统一样，在至关重要的安全领域面临着很多挑战，原因在于：

- 复杂而不确定的环境：在很多情况下，AI 系统需要在复杂的环境下运行，其中，有很多潜在的情形是无法完全在事先进行检查和测试的，有时候，AI 系统会面临设计人员完全无法预料到的情境。
- 突然出现的行为：对于那些在部署后可以自行学习的 AI 系统，其行为很大程度上可能是由系统在无人监管的情形下学习的那段时间决定的。在这种情况下，AI 系统的反应和行为将很难预测。
- 目标设定的偏差：将人类的目标“翻译”成计算机指令是很困难的，因此，AI 系统接收到的程序指令的目标，可能和程序员想要其实现的目标出现偏差。
- 人机交互：在很多情况下，AI 系统的表现会被与人类的互动所影响，此时，人机交互时人类的不同反应也会影响 AI 系统的安全性。

为了解决以上以及其他问题，需要对更先进的 AI 系统的安全性进行额外的投资，其中包括其解释能力、透明性、信任度、可验证性以及防攻击的安全性、AI 长期的安全性和与价值观的调整等方面。

提升 AI 解释能力、透明度

在实际研究中，提升 AI 的“解释能力”（explainability）和透明度（transparency）是很大的一项挑战。很多算法，包括那些基于深度学习的算法，对于用户来说都是含糊的，

目前 AI 领域很少有对结果进行解释的机制，在一些领域，比如医疗，这是一个非常大的问题，因为医生在进行诊断和对治疗方式下结论时，需要——进行解释。一些 AI 技术，比如决策树推理(decision-tree induction)，自带的解释功能经常提供不太准确的解释。因此，研究人员需要开发更加透明的、出自本能可以向用户解释其结果和行为的系统。

建立信任

为了赢得信任，AI 系统的设计者需要创造准确的、可靠的系统，以及具有指导性的、用户体验高的界面，同时，AI 系统的操作者也需要足够的时间，获得关于系统运行方面的训练，同时了解系统表现方面的局限性。那些获得用户广泛信任的复杂系统，比如汽车的手动操控系统，一般都是透明（系统运行的情况用户都可以看见）、可信赖（系统的反应都是用户接受的）、可审查（系统可以被评估）、可靠（系统的行为反应都是用户想要的），以及可恢复（需要时用户随时可以拿回掌控权）。现在和未来的 AI 系统都面临一个显著的挑战，那就是软件生产技术在质量方面的不稳定性。技术的发展使得人类和 AI 的联系更加紧密，在如何建立信任方面，最大的问题是如何追上不断变化和提升的 AI 的性能、预测 AI 在技术方面的长期的发展和应用，以及建立指导性的原则和政策，以保证在 AI 研究能从设计、打造到使用都遵循最好的实践，包括对操作人员的安全运行方面的适当的培训。

增强可验证和可确认性

在验证性和确认性方面，AI 系统需要新的方法。“验证性”（verification）是建立一个满足很多正式的要求的系统，而“确认性”（validation）则是建立满足用户在操作方面的需求的系统。安全的 AI 系统可能需要在评估（确定系统是否在不正常运行，特别是在非预期的环境下）、诊断（确认系统非正常运行的原因）、修复（调整系统，解决非正常运行的问题）方面需要新的手段。对于那些超时自主运行时间的系统，设计者可能无法考虑系统可能遇到的所有情形，这样的系统要想保持可靠并且具有活力，需要带有自行评估、自行诊断和自行修复的能力。

美国宇航局艾姆斯研究中心 (NASA Ames Research Center) - 在故障发生前进行预测

由于模型驱动故障检测方法存在不足之处，美国宇航局艾姆斯研究中心在2003年开发了数据驱动的故障检测方法，名为“归纳式监测系统”（Inductive Monitoring System, IMS）。从那以后，美国宇航局就在内部部署了这一检测系统，包括对航天飞机和国际空间站（ISS）的监测，一些美国宇航局之外的领域也采用了这一系统。



2014年，猎户座载人飞船（Orion Crew Vehicle）试飞，在这个过程中，IMS被用于电子系统的监测。



C-130“大力士”（C-130 Hercules）军用运输飞机使用了预测式监测软件，用来预测引擎之间气流交换阀的故障。

2012年，综合性工程管理解决方案（Comprehensive Engineering Management Solutions, CEMSol）购买了IMS系统并对其进行改进，与美国宇航局艾姆斯研究中心以及洛克希德·马丁公司（Lockheed Martin）一起合作，在洛克希德C-130“大力士”军用运输飞机上对其进行测试。在此次测试中，洛克希德·马丁公司投入7万美元，但是，该系统帮助洛克希德·马丁在维护成本和项目延期方面节省的费用达到了投入数额的十倍。

对攻击的防护

在重要系统种部署的 AI，必须具有非常强的活力才能应对意外，需要能够应对一系列的国际网络攻击。安全方面的工程包括理解系统的弱点，以及理解潜在的攻击者的行为背后的意图。尽管 NITRD（The Networking and Information Technology Research and Development，美国政府的网络与信息科技研究与发展项目）发布的网络安全研发战略计划（Cybersecurity R&D Strategic Plan）已经对网络安全方面进行了详细的计划，但是对于 AI 来说，仍然需要面对一些特定的风险，比如，一项关键的研究是“敌对机器学习”（adversarial machine learning），也就是在 AI 的训练数据中加入恶意的数据，或者改变算法、对一项事物进行微小改变以便 AI 无法正确对其进行身份确认（比如对面部进行一些伪装来防止 AI 的面部识别）等，由此来探索 AI 系统将受到多大程度的影响。在高度自主的网络安全领域部署 AI，也是需要更深入的研究的领域，这方面一项比较近期的例子是，美国国防部先进研究项目局（Defense Advanced Research Projects Agency, DARPA）在其“网络大挑战”（Cyber Grand Challenge）项目中使用 AI 代理，对网络攻击进行分析和反击。

获得 AI 的长期安全，对其价值观保持调整

最终，AI 系统可能可以实现递归式自我提升 (recursive self-improvement)，也就是说，对软件的改进将由软件自己完成，而不是由人类的程序员，要保证自我改善的 AI 系统的安全性，仍需要更多的研究：可以时刻自我监测其系统行为是否与人类设计员最初的目标一致的架构；在系统被评估的时候防止其释放的限制措施；价值学习，也就是说系统应具备推测用户的价值观、目标和意图的能力；价值观方面可以防止自我改变的观框架。

战略 5：开发可供 AI 培训和测试的公共共享数据集和环境

AI 能够使人类获益的能力将会继续增加，但前提是行业开发出 AI 领域的培训和测试资源。培训数据集和其他资源的多样性、深度、质量、准确性将显著影响 AI 的表现。很多不同的 AI 技术都需要有高质量的培训和测试数据、动态和互动的测试台和模拟环境。这方面并不是一个简单的技术问题，而是一个“大众福利”的问题，因为如果 AI 培训和测试的资源仅仅掌握在少数的机构手中，那么 AI 的发展就会受到限制，我们必须同时尊重 AI 数据的商业价值和个人对 AI 的兴趣及权利。我们需要针对 AI 的很多应用开发出高质量的数据集和测试环境，以保证公众能够享有合理的数据集以及测试和培训资源，还需要更多的开源软件库和工具，来帮助 AI 的研发加速发展，接下来的小节将列出重要的几大领域。

开发一系列数据集以满足 AI 领域的兴趣和应用的请求

AI 培训和测试的数据集的开发对于获得科学可靠的结果至关重要，而在数字化时代，对技术以及社会-技术 (socio-technical) 基础设施的打造以支持可持续的研究是非常具有挑战性的——同时对于 AI 技术的发展又不可或缺。缺乏经验证的、对公众开放的数据集，将对 AI 的发展产生限制。与其他数据密集型科学领域一样，经过数据的验证对于 AI 无比重要。研究人员必须要能够在相同和不同的数据的验证下获得同样的结果，用于研究的数据集必须要能够代表现实世界应用中的挑战，而不是只采用简化了的版本。为了缩短这个过程，政府和政府资金支持的机构，甚至行业内的机构，必须将现有的数据集开放。

AI 在机器学习方面遇到的挑战往往都与“大数据”分析有关。考虑到数据集的多样化，从中获取具有代表性的数据集，并且对这些非结构性或半结构性的数据进行分析是非常困难的。如何从绝对和相对 (看具体情境) 的角度挑选具有代表性的数据？目前，现实世界的数据库非常不连续、不准确而且具有很大的干扰性。因此，一些数据预处理的技术 (比如数据净化、整合、转化、缩减和取样) 对于建立实用的针对 AI 应用的数据集非常重要。数据预处理技术如何影响数据的质量？尤其是在对数据进行额外的分析的时候？

鼓励政府资金支持的研究和机构分享 AI 数据集很有可能将能够刺激 AI 领域的创新发展,但是,保证数据分享的安全性也需要相应的技术,因为数据的拥有者在将他们的数据与 AI 研究社区进行分析的时候会承认一定的风险。另外,数据集份额开发和分析也必须遵循相关的法律和规定,同时也必须遵循相关的道德规范。风险可能会来自多个方面:比如对数据集的不正当使用、不准确或者不正当的泄露、委保证隐私和保密需求而进行的数据去身份化技术的限制性等。

分配培训和测试资源以满足商业和公众利益

随着数据持续地爆炸式发展,数据的来源、世界范围内的信息技术不管是在数量还是规模方面都在增长,与此同时,对数据进行分析的技术却无法赶上信息规模的增长速度。数据获取、整理、分析和视觉化技术在研究中都面临很多的挑战,而从海量数据中提取有价值的信息所需要的科技的发展已经滞后。虽然针对数据的资料库已经存在,但是它们往往无法应对数据集的规模化发展、具有非常有限的经验证的有效信息以及无法支持富语义方面的搜索,我们需要动态的、敏捷的数据资料库。

支持 AI 研究所需的这种开放、分享式基础设施项目方面的一个例子是由国土安全部开发的 IMPACT 项目(网络风险&信任的政策和分析方面的信息市场, Information Marketplace for Policy and Analysis of Cyber-risk & Trust)。这个项目通过协调和发展现实世界的数据和信息分享能力(包括工具、模型和方法)来支持全球网络安全方面的研究。IMPACT 项目也支持国际网络安全研发社区、关键基础设施提供商、政府支持者等之间的实验性的数据分享,类似的项目能够在各个应用领域帮助 AI 研发的加速。

开发开源软件库和工具箱

开源软件库和工具箱的存在为所有能够联网的开发者提供了接触顶级 AI 技术的渠道。类似 Weka 工具箱、MALLET 和 OpenNLP 等资源,帮助加速了 AI 的发展和应用。开发工具,包括免费和低价的开发语言(比如 Octave 和 Python)降低了使用和扩展这些数据库的门槛。此外,对于那些不希望直接整合这些数据库的人,任何基于云的机器学习服务都可以通过低延时网络协议来进行比如图像分类等任务,这个过程中不需要任何的编程工作。最后,很多这样的网络服务还提供特定的硬件的使用,包括基于 GPU 的系统。从这个角度来说,我们基本可以认定那些针对 AI 算法的特定软件,比如仿神经处理器等,也可以通过类似的服务提供给公众。

这些资源联合起来可以提供 AI 技术的基础设施,以通过企业家在无需使用昂贵的硬件或软件、无需成为 AI 领域的专家的前提下,开发出解决狭义领域的问题的解决方案,最终

促使市场的创新发展。对于狭义的 AI 领域，市场创新的门槛相比于很多科技领域来说是非常低的。

为了保持这个领域的高度创新性，美国政府可以支持对开放 AI 技术的开发和使用，特别是那些使用标准的开放格式和使用代表语义信息的开放标准的资源，将会对这个领域的创新发展极为有利。

政府还可以通过加速政府自身内部对于开放 AI 技术的使用，来鼓励社会对于开放 AI 技术的接纳，以此为创新者提供比较低的门槛。政府应该尽可能地对开放软件项目的算法和软件做贡献。由于政府也会有特定的担忧，比如更为强调数据的隐私和安全性等，所以可能政府也需要开发一些机制，为政府自身接纳 AI 技术铺平道路，比如，可能需要在各个政府部门之间开展一次“地毯式扫描”，以寻找各部门特定的 AI 应用的领域，然后再确定为了可以使用 AI 技术，这些部门和机构需要解决的问题。

战略 6：制定标准、参照来衡量和评估 AI 技术

在指导和推广 AI 技术的研发时，标准、参照、测试台，以及 AI 社区对其的接纳至关重要。以下小节将介绍那些仍需要更多进展的领域：

开发广泛的 AI 标准

标准的制定必须追上 AI 应用的快速发展。这些标准可以提供要求、配置、指导原则、特点等能够保证 AI 技术能够满足很多重要目标的要求，同时保证很多功能的实现以及 AI 技术之间的互通，只有这样 AI 系统才能够安全可靠的运行。标准的实行能够为技术的发展和互通铺平道路。AI 相关的标准方面的一个例子是由电气和电子工程师协会（Institute of Electrical and Electronics Engineers，IEEE）制定的 P1872-2015（机器人和自动化的标准本体，Standard Ontologies for Robotics and Automation）标准。该标准为知识的代表、常用的数据和定义提供了一个系统的方法，因此，人、机器人和其他人工智能系统之间，无歧义的信息得以传递，同时该标准也为机器人领域 AI 技术的应用奠定了基础。AI 的细分领域都需要这样的标准，具体需要提供的标准包括：

- 软件工程：对系统复杂性、可持续性、安全性进行管理，同时监控和控制突然出现的行为；
- 表现：保证准确性、可靠性、稳健性、可访问性，保证规模可扩大；
- 数据：对影响表现以及是否符合标准的因素进行量化；

- 安全：评估风险控制和系统、人机交互、控制系统、是否合规等安全方面的情况进行分析；
- 实用性：保证交互界面和控制系统有效、易于操作；
- 互通性：通过标准和兼容的交互界面对可以互通的部件、数据和交易模型进行定义；
- 保密：解决保密、信息安全和网络安全方面的问题；
- 隐私：对信息在处理、传送和存储之时的安全保护的掌控；
- 可追踪：对所有发生的事件提供记录（执行、测试和完成），以及对数据的整理；
- 领域：对具体领域的标准和相应的框架进行定义

制定 AI 技术的基准

针对测试和评估而制定的基准，可以为 AI 发展的开发标准和评估是否符合标准提供可量化的衡量工具。行业基准可以通过一些战略性的情景推广业内取得的新进展来推动行业创新；此外还能提供客观的、可用于检测 AI 技术的进展的数据。为了有效地评估 AI 技术，必须开发相关的测试方法和衡量数据并将之标准化。标准的测试方法将提供相关的协议和流程，以对 AI 技术的表现进行衡量、对比和管理。为了描述 AI 技术的特点，必须有量化的衡量指标，这些指标包括但不限于：准确性、复杂性、信任度和能力、风险、不确定性；解释能力；意外偏差；与人类表现的比较；经济影响等。需要注意的是，基准的制定需要靠数据驱动，战略五讨论了培训和测试的数据集的重要性。

作为 AI 相关的技术的一个成功例子，美国国家标准和技术协会（National Institute of Standards and Technology, NIST）开发了一套全面的测试方法和表现衡量数据，来评估应急机器人的表现，目的是通过在机器人进行标准测试时获取的数据对不同型号的机器人的表现进行量化的比较。这样的比较结果可以用来支持采购决定，以及帮助开发者了解部署的性能。美国材料与试验协会的国际标准委员会正在对这样的测试方式是否应用用作国土安全应用的机器人操作设备的标准（被称为 E54.08.01 标准）进行评估，而这样的测试方法的各种版本正被研究社区应用于 RoboCup 救援机器人联盟的比赛之中，后者强调的是自主运行的能力。另一个例子是 IEEE 的工业自动化敏捷机器人大赛（Agile Robotics for Industrial Automation Competition, ARIAC），这是由 IEEE 和 NIST 联合举办的，目的是利用 AI 和机器人策划领域的最新进展来推广机器人的敏捷性，这项大赛的一项核心内容就是测试工业机器人系统的敏捷性，目的是促进那些工厂厂房里的机器人变得更高效、更自主，对工人的时间要求更低。

尽管这些措施为制定 AI 基准提供了很好的基础，但是它们仍限制于特定的领域。更广阔的 AI 领域需要更多的标准、测试台和基准，来保证 AI 解决方案得到广泛地应用。

提供更多的 AI 测试台

测试台的重要性在《未来的网络实验》(Cyber Experimentation of the Future) 报告中有详细介绍，“测试台的重要性在于，其可以使研究人员使用可以运行的数据到现实世界中的系统中去实验...也可以在好的测试环境中进行实验。” AI 所有的领域的发展都需要足够的测试台。政府拥有海量的敏感信息，很多数据都无法向研究社区之外公布。为了能够让学界和业界人员在安全的测试环境下工作，我们需要一些合适的措施和项目。通过这些测试台，AI 模型和实验方法可以在研究社区进行分享，为 AI 科学家、工程师和学生提供独特的研究机会。

AI 社区需要参与标准和基准的制定

为了促进 AI 标准在政府、学术机构和行业的应用，政府的领导和协调对于 AI 标准的制定至关重要，AI 社区——由用户、业内人士、学术人员和政府人员组成——必须积极参与到这些标准和基准的制定过程中。根据角色和职责的不同，不同的政府机构在吸引 AI 社区的参与工程中扮演着不同的角色，这种协调作用在收集用户驱动的需求、预测开发者驱动的需求和推广教育机会方面是很重要的。用户驱动的需求决定了 AI 技术发展过程中遇到的挑战和发展目标，而社区的基准着眼在确定进展、拉小差距和针对特定问题促使创新解决方案方面。这些基准必须包括定义和分配技术事实的方法，创造基准模拟和分析工具也将加速 AI 的发展。这些基准带来的结果还可以将适当的技术与特定的用户需求相匹配，为标准的合规、合格产品名单和潜在的供货商选择提供客观标准。

行业和学术机构是 AI 技术发展的主要来源，推广和协调他们对标准和基准的制定的参与至关重要。随着解决方案的诞生，越来越多的机会也在诞生，比如通过分享大家在技术机构、开发中青睐的标准、测试以确定高质量和互通的解决方案等方面的共同的愿景来预测开发者驱动和用户驱动的的标准。

高影响力、基于社区的、AI 相关的基准项目方面一个成功的例子是文本检索会议 (Text Retrieval Conference , TREC)，是由 NIST 从 1992 年开始的、未大规模文本检索的发展提供所需的技术设施的大会，超过 250 个小组参与了 TREC，包括学界和业界各种规模的各种机构，制定的标准被广泛应用，其数据是由 TREC 制定，对于信息检索的研究发展的重新激活做出了贡献。第二个例子是 NIST 的阶段性的机器学习在生物测定学领域的应用基准的项目，特别是面部识别机会，始于 1993 年的面部识别技术大会 (Face Recognition Technology , FERET)，为面部识别的发展提供了标准的数据集，帮助支持面部识别算法的

开发，并制定了一项评估协议，这个大会经过这些年的发展，已经变成了面部识别供应商测试大会（Face Recognition Vendor Test，FRVT），包括了数据集的分发、难题的解决和进行接下来的技术的评估等。这项基准制定的大会为面部识别技术的发展做出了很大贡献。不管是 TREC 还是 FRVT，都是 AI 相关的社区参与基准制定活动的成功例子，AI 其他领域的发展还需要类似的努力。

需要注意的是，标准的开发和应用，以及参与基准制定的活动，是要付出代价的。研发机构只有在可以看到显著的好处的时候才有参与的动力。集合各个部门和机构参与到 AI 标准的制定中，能够进一步促使他们参与标准的制定和推广。基于社区制定的基准，比如 TREC 和 FRVT，也通过提供培训和测试数据降低了行业的门槛，培养健康的竞争，促使最好的算法的但是，同时也为行业内参与者提供客观和可比较的指标。

战略 7：更好地理解国家 AI 研发劳动力需求

实现本战略所列出的 AI 研发进展需要很大的研发力量，在 AI 研发领域具有强大的布局的国家，将能够引领未来世界自动化的发展，成为算法创造和开发、性能展示、商业化等领域的领跑者。技术方面的发展将能够为这些进展打下基础。

尽管还没有官方的关于 AI 从业人员方面的统计，无数来自商业和学术机构近期的报告显示，AI 技术人才出现短缺，而需求正在急速上涨。据报道，高科技公司大量投资，招聘 AI 领域的学术专家和学生，高等学府和行业机构之间展开了针对 AI 人才的争夺大战。

要想更好地理解国家未来在 AI 研发方面的劳动力需求，我们还需要更多的研究，也需要更多的数据来对目前 AI 研发人才市场的情况进行总结，包括学界、政府和商业领域各自的需求。我们需要探索 AI 劳动力市场的需求和供应的研究，以便对未来的需求进行预测，同时未来对 AI 研发市场的未来供应情况进行预测，我们也需要更多的理解这个市场，需要同时考虑教育的途径和潜在的再教育的机会，鉴于研究显示多样化的 IT 人才组合可以提升生产力，多样化的问题也需要更多的探索。一旦我们更好地理解眼前和未来的 AI 研发劳动力需求，就可以制定合理的计划和行动，去解决现有的或者未来潜在的 AI 研发劳动力市场面临的问题。

建议

联邦政府整体将支持本计划中列出的 7 大战略优先级，并通过支持以下建议来实现其目标：

建议 1：根据本计划中的战略 1 至战略 6，开发人工智能研发执行框架，识别科技机会，支持高效的人工智能研发投资合作。

联邦政府部门应当通过 NITRD 展开合作，制定研发执行框架，就本计划列出的研发挑战推进合作和发展。这将帮助联邦政府部门更便捷地计划、协调和合作，支持这一战略计划。执行框架应当基于使命、能力、主管范围和预算考虑每个部门的研发优先级。根据执行框架，可能需要成立投资项目，协调国家人工智能研究计划的执行。为了执行这一战略计划，NITRD 应当考虑成立专注于人工智能的跨部门工作组，并与现有工作组展开协调。

建议 2：根据本计划中的战略 7，研究美国的行业形势，培育并维持健康的人工智能研发工作者团队。

在解决本报告提出的研发战略挑战时，健康、有活力的人工智能研发团队非常重要。尽管有报告指出，人工智能研发专家的短缺可能越来越严重，但目前还没有任何正式的就业数据，描述人工智能研发工作者的状态、预期的劳动力储备，以及人工智能劳动力的供需情况。考虑到在解决本计划提出的战略优先级的过程中，人工智能研发工作者扮演的角色，我们需要更好地理解这些情况，培育并维持健康人工智能研发工作者团队。NITRD 应当研究，如何更好地描述和定义，当前和未来人工智能研发工作者队伍的需求，展开新的研究或作出建议，确保有充足的研究工作者去解决美国的人工智能需求。正如研究结论所说，适当的联邦机构应当采取行动，确保建立并维护健康的全国人工智能研发工作者团队。

首字母缩略词

| | |
|----------|---|
| 3-D | Three Dimensional |
| AI | Artificial Intelligence |
| ANNs | Artificial Neural Networks |
| ARIAC | Agile Robotics for Industrial Automation Competition |
| ARMOR | Assistant for Randomized Monitoring over Routes |
| ASTM | American Society of the International Association for Testing and Materials |
| ATM | Automated Teller Machine |
| BRAIN | Brain Research through Advance Innovative Neurotechnologies |
| CEMSol | Comprehensive Engineering Management Solutions |
| COMPETES | America Creating Opportunities to Meaningfully Promote Excellence in Technology Education and Science |
| CoT | Committee on Technology |
| DARPA | Defense Advanced Research Projects Agency |
| DHS | Department of Homeland Security |
| DoD | Department of Defense |
| DOE | Department of Energy |
| DOT | Department of Transportation |
| FERET | Face Recognition Technology |
| FRVT | Face Recognition Vendor Test |
| GPS | Global Positioning System |
| GPU | Graphics Processing Unit |
| HPC | High Performance Computing |
| I/O | Input/Output |
| IBM | International Business Machines Corporation |
| IEEE | Institute of Electrical and Electronics Engineers |
| IMPACT | Information Marketplace for Policy and Analysis of Cyber-risk & Trust |
| IMS | Inductive Monitoring System |
| IoT | Internet of Things |

| | |
|--------|--|
| IRIS | Intelligent Randomization in International Scheduling |
| ISS | International Space Station |
| IT | Information Technology |
| KSA | Knowledge, Skills and Abilities |
| LAX | Los Angeles World Airports |
| MALLET | Machine Learning for Language Toolkit |
| NASA | National Aeronautics and Space Administration |
| NCO | National Coordination Office for NITRD |
| NIH | National Institutes of Health |
| NIST | National Institute of Standards and Technology |
| NITRD | Networking Information Technology Research and Development |
| NLP | Natural Language Processing |
| NRL | Naval Research Laboratory |
| NSF | National Science Foundation |