

AI 时代，我们要像大模型一样进化

▶ 刘嘉

大模型的成功并非偶然——从早期符号主义 AI 的失败，到深度学习的崛起，再到 Transformer 的成功，每一次进化都是从无数被淘汰的算法、模型中艰难诞生。在这艰难曲折的探索中，人类智慧的金块无疑是 AI 头上的一盏明灯。反过来，大模型的进化经验，能否成为我们人类认知进化的营养？由此，我们破茧成蝶，与 AI 时代同频共振，开启认知与智慧的跃迁。

为人生定义目标函数

所有的机器学习，在开始训练前，都必须明确一个目标函数（又称损失函数或成本函数）。这个函数定义了模型希望达到的理想状态，而训练的全部意义就在于不断优化参数，让模型越来越接近这个目标。正所谓学习未动，目标先行。

作为机器学习的一个分支，神经网络从一开始就是另类，因为它的目标函数太宏大、太有野心，以至于当辛顿请求其所在的多伦多大学校长再招收一名神经网络的研究者时，该校校长是如此

回答的：“一个疯子就足够了。”的确，神经网络的开创者都有一个在外人眼里近似疯狂的目标函数：1943 年麦卡洛克和皮茨提出的“简陋”神经元是要模拟“神经活动内在观

刘嘉
清华大学基础科学讲席教授，清华大学心理与认知科学系主任、人工智能学院教授，长期从事脑科学与人工智能交叉研究。



念的逻辑演算”，1958 年罗森布拉特提出的第一个真正意义上的人工神经网络——感知机，是要模拟“大脑信息存储和组织”。OpenAI 训练 GPT 的目标函数，就是要用一个巨大的神经网络去容纳所有的人类知识，从而实现 AGI（Artificial General Intelligence，通用人工智能）。

虽然疯狂却是唯一可行之路。GPT-4 把几乎全部的人类知识压缩进了 1.8 万亿个参数，在通用认知任务上的表现卓越，从此 AGI 不再是科幻且遥不可及的。神经网络宏大的目标函数的背后是规模化法则：参数规模越大，优化空间越广，最终实现目标的可能性越大。

人类学习也遵循同样的道理，如果我们把目标函数设定为短期、狭隘的目标，如考取某个证书、通过某次考试，那么这个目标函数的确容易实现。但是，我们得到的只是一个线性模型，目标只要稍微复杂一点、稍微变化一点，这个线性模型就再无用处之地。这在机器学习中也称为“局部最优”陷阱。当一个模型陷入局部最优的舒适区，就不再



* 本文摘自刘嘉教授撰写的《通用人工智能：认知、教育与生存方式的重构》一书（中信出版集团，2025 年 6 月），原书章节标题为“像大模型一样进化”

演化，最终错过了更广阔、更深远的优化空间以抵达“全局最优”。同样，人生的发展也会出现局部最优——在人生某个阶段取得了看似不错的成就，实际上却限制了后续的发展空间。所以，短期看是目标达成，长期看则是机会丧失。

人本主义心理学家亚伯拉罕·马斯洛曾经问他的学生：“你们当中，谁将成为伟大的领导者？”

学生只是红着脸，咯咯地笑，不安地蠕动。马斯洛又问：“你们当中，谁计划写一本伟大的心理学著作？”学生结结巴巴地搪塞过去。马斯洛最后问道：

“你们难道不想成为一个心理学家吗？”这时，所有学生都回答“想”。这时，马斯洛说道：“难道你们想成为平庸的心理学家？这有什么好处，这不是自我实现。”马斯洛解释道，我们其实不仅仅害怕失败，也害怕成功。

在这现象的背后，是与自尊纠缠在一起的自卑：

我们对伟大的人和事物都有一种敬畏感——在面对他们时，会感到不安、焦虑、慌乱、嫉妒甚至敌意，因为他们会让我们产生自惭形秽的卑微感。于是，当我们试图获得荣誉、成功、幸福等美好的事物时，还未行动，我们却产生了“这是真的吗”“我不行”“我不配”的自我质疑，因为陌生的阳光如同黑暗一样可怕。

萨姆·奥尔特曼在一次接受采访时，回忆起刚创业时遭到的嘲讽：

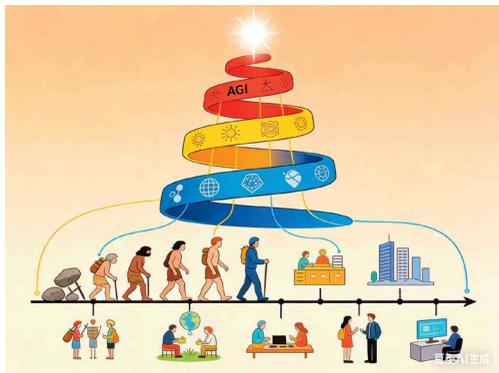
“回想起来，一件非常成功的事情是，我们从

一开始就确定了 AGI 的目标，而当时在业内，你甚至不能谈论这个目标，因为它听起来太疯狂了，近乎痴人说梦。所以这立即引起了苏茨克维的注意，也吸引了所有优秀年轻人的注意，当然，也引来了不少前辈的嘲笑。不知何故，我觉得这是一个好兆头，它预示着某种强大的力量。我们当时是一群乌合之众，我的年龄是最大的，大概 30 岁，所以当

时大家觉得我们这群不负责任的年轻人什么都不懂，净说些不切实际的话。但那些真正感兴趣的人会说，‘让我们放手一搏吧！’”

这就是 OpenAI 的目标函数，所以才有今日之 OpenAI。

唯有那些真正基于深层理解、价值判断和创造性思维的能力，才属于人的不可替代的能力



作为个人，我们的目标函数应该是什么？

在我看来，那就是构建属于我们自己的、特立独行的“个人知识体系”。我们的知识体系是我们认知世界的“眼睛”，正如色盲者无法正确分辨这个世界的颜

色，而一个知识体系有缺陷的人不可能触摸到这个世界的本质。

进入 AGI 时代，个人知识体系的重要性被无限放大，这是因为 AI 正在逐渐接管那些标准化、结构化的信息处理任务，而唯有那些真正基于深层理解、价值判断和创造性思维的能力，才属于人的不可替代的能力。而这些能力，恰恰植根于独特的个人知识体系之中。所以，不断拥抱新的经验、新的知识，更新推理思维链，打破认知边界，都是在构建一个能与世界深度对话、与自我持续共鸣的个人

知识体系。

“兰叶春葳蕤，桂华秋皎洁。欣欣此生意，自尔为佳节。”马斯洛说，这才是“奔放的人生”，而不是“枯萎的人生”，因为“如果你总是想方设法掩盖自己本有的光辉，那么你的未来注定暗淡无光”。

使用随机梯度下降优化人生

在机器学习中，随机梯度下降（stochastic gradient descent, SGD）是被广泛使用的优化算法之一。其原理简单而高效：每一步都在当前的位置基础上，找到一个大致的方向，然后往那个方向迈进一小步。而这个大致正确的方向，来自当前的误差——算法通过不断迭代调整模型参数，沿着矫正误差最陡梯度前进，逐步找到使损失函数最小的参数值。所以，正是因为存在误差，我们才能判断前进的方向。

在大模型的预训练过程中，输入的数据首先被表示为一系列的 token，这些 token 逐层穿过神经网络的各个隐藏层，并最终在输出层生成下一个 token 的预测值。模型根据上下文生成的预测值与实际语料中的真实 token 之间往往存在一定差异，这个差异就是模型的预测误差（error），具体可表示为误差函数： $error = diff(\text{预测值} - \text{实际值})$ 。大模型正是利用这个误差信息进行学习，通过反向传播算法将误差逐层传递回网络中的每个神经元，以确定每个参数的优化方向与幅度，再使用随机梯度下降等优化算法，逐步调整和更新网络的权重参数，以持续减小损失函数的数值，提升模型预测的准确性。

由此，大模型的学习过程就构成了一个不断循

环的优化流程：预测下一个 token → 计算误差 → 反向传播误差 → 利用梯度下降优化参数 → 更新模型权重 → 预测下一个 token。大模型的所有知识和能力，便是通过反复地执行上述循环、不断根据误差进行参数调整而逐渐获得的。

大模型只能从错误中学习，人也不例外。这是因为梯度下降的优化算法与大脑的预测编码（predictive coding）机制有异曲同工之妙。预测编码理论认为，大脑是一个主动预测外部世界的系统，它不断根据已有的经验形成预测，随后将这些预测与现实中接收到的信息进行对比。当预测与实际感知之间出现差异时，大脑就会产生误差信号（prediction error）。

这种误差信号会激活大脑中与奖赏和纠错机制相关的多巴胺系统，从而重塑大脑神经元之间的连接。

换句话说，错误为大脑提供了一种清晰的、明确的反馈信号，帮助我们快速地发现原有知识或方法的不足，迫使我们重新审视自己原有的信念或行为模式，并尝试新的、更加准确的做法。与之相反，当我们的预测正确、表现良好时，大脑获得的反馈信号是弱而模糊的。所以成功的体验非常美好而错误让人痛苦，但是我们的成长来源于如何应对、修正错误，因为错误本质上并非失败，而是一种推动我们持续更新认知结构、增强适应能力的动力源泉。

但是，人是追求奖励、逃避惩罚的动物，“少犯错、不犯错”是我们所接受的教育的核心，所以主动试错对我们而言是知易行难。随机梯度下降则为此提供了解决之道。

随机梯度下降的核心魅力之一，在于它能从不确定中找到确定性——目标函数清晰，但是通向目

我们需要做的，就是“强行起飞，粗糙开始，空中加油”

标函数的路径不确定。也就是说，我们不要执着于精确地规划未来的每一步，因为这样反而可能陷入过度分析而迟迟无法行动。我们需要做的，就是“强行起飞，粗糙开始，空中加油”——找一个大致正确的方向（梯度），然后向前走一步（下降）。

不必在乎当下的这一步是否最优，做时间的朋友，能多走几步就多走几步。

因为对于梯度下降这件事，起点不重要，终点才重要。起点都是初始化的随机参数，众生平等；终点则是损失函数的能量最小值。所以，家境是否优渥不重要，是不是名牌大学毕业不重要，年龄太大也不重要，因为这些都只是起点，或者最多只能算是“中点”而非终点。梯度下降算法能保证的是：不管起点在哪里，最后得到的解都差不多，当然前提是一直按照梯度的方向走下去。所以，坚持走。

然后，四处走走（随机），因为每一个方向都是你对世界的新认识。包容性和灵活性是随机梯度下降的核心魅力之二。如果只是沿着熟悉的道路前进，虽然容易并且安全，却可能会让你陷入认知的局部最优陷阱——你以为自己已经理解了整个世界，实则只是固守在一个狭窄的角落。

正如随机梯度下降强调随机抽样是为了避免陷入局部最优，人生也需要随机性的探索，这样才能发现没有见过的风景。阅读陌生领域的书籍，与不熟悉的人交谈，尝试未知的可能性，正是利用了随机性带来的认知增益。它引领我们遇到新的误差、新的意外，并因此而激发新的学习过程，推动认知结构的重新构建。正是在随机探索中，我们不断修正对世界的理解，逐渐接近真实。

随机，不仅是算法优化的策略，更是我们深入认识世界、走向自我更新的重要方法。

奥特曼曾经谈到他的一次“四处走走”：

我在26岁时卖掉了我的初创公司，然后中间

空了一年。在那个年代，在硅谷这是很难想象的行为，因为那是一个根据你的职位和你所做的工作确定社会地位的地方。但是如果你真的可以在两份工作之间空出一年，我是非常推荐的，我甚至觉得这是我职业生涯里做得最对的事情。在那一年里，我读了很多书，在很多感兴趣的领域有所涉猎。……我学到了核工程知识；AI时代开始了，我学习了关于AI的理论；我学习了生物制造的相关知识。……我到很多地方旅行，从某种程度上讲，我感受到了这个世界其他部分真实的样子，我见了从事各行各业的人，并与之交谈……我有充足的时间，所以如果我遇到了有意思的看起来不错又需要帮助的人，我会帮助他们。……我没有安排自己的时间表，所以我可以立刻飞到其他国家参加会议。我开始做这些随机的事情。几乎所有的事情都没有开花结果，但是对之后事情产生深远影响的种子已经种下了。

这个种子，最终发芽成长为OpenAI。

人生所需不过一份注意

GPT的T，指的是Transformer，其最核心、最精妙之处就是“注意力机制”。它会对一段文本中每个词语与其他所有词语之间的关系进行评估，计算出它们之间的关联强弱程度，从而捕捉信息之间的相互关系，以实现高效而精准的信息处理。所以，学习的本质也是注意力分配的艺术。

我们所处的世界彼此相连，而非孤立随机。在物理层面，世界由物质和能量组成，它们之间不断地相互作用，形成复杂而稳定的秩序。在生命层面，物种之间通过复杂的生态网络连接起来，生态链中每个环节互依互存，任何个体的变化都可能引发连锁反应。在人文社会层面，每个人看似独立，但无时无刻不在通过沟通、情感联结与社会网络交织在一起。文明的存续与演化，来源于人与人之间频繁而有序的互动。

英国诗人约翰·多恩说：“没有人是一座孤岛，可以自全。……任何人的死亡都是我的损失，因为我是人类的一员，因此不要问丧钟为谁而鸣，它就为你而鸣。”美国行为科学家阿莫斯·特沃斯基也说：“人不复杂，复杂的是人与人之间的关系。”

应当如何分配注意力来认识我们所在的这个世界呢？

第一，注意高质量的数据和人。在机器学习领域，有一个广为人知的第一性原理：“垃圾输入，垃圾输出。”再多的参数，再强大的算力，如果输入的数据质量低下，最终训练出来的大模型也必然表现糟糕。所以，OpenAI 在训练初期便严格把控数据质量，选用了维基百科、经典书籍、科研论文、优秀代码和高质量互联网内容作为注意力处理的信息。这些精心挑选的材料构成了 GPT 的认知基座。

截至 2024 年 6 月，我国短视频用户数量达到 10.5 亿，占整体网民的 95.5%，人均每天观看时长约 151 分钟。而阅读用户只有短视频用户的一半，人均每天阅读时长只有 23 分钟。AI 在学习，人类却在沉迷。

真正与注意力门当户对的是高质量的数据集和人。在进入某个领域前，首先精心构建你的数据集：谁是这个领域的权威，哪些书、线上课程是这个领域的经典，哪些工具能让这个领域的抽象知识变得具象清晰？之后，阅读入门材料快速建立对这个领域的基本认知；接下来，对经典或权威的书籍或教材进行深度学习，建立完善的知识框架；最后，通过专业研究文献并与专家或 AI 互动交流，拓宽和深化自己的认知边界。

第二，注意实例而非规则。符号主义给 AI 以规则：“如果一个动物有尖尖的耳朵，胡须明显，并且眼睛在夜间能反光，那么它是猫。”这时，狐狸、猓狍、浣熊和狼也会被符号主义 AI 识别成

猫。而联结主义只给 AI 猫的图片，各种各样猫的图片，让注意力在海量的数据中主动探寻其中蕴含的模式和规律。前者是授人以鱼——人类先提取特征，然后把特征喂给 AI，即人类向 AI 输入人类学习的结果，AI 只需要记忆，正所谓前面有多少智能，背后就有多少人工。后者是授人以渔——没有工程师总结的规则，只有精心挑选的实例，让神经网络自己学习，让它自己去充分挖掘全部可能，因为“足够大的神经网络当然无所不能”（计算软件 Mathematica 的创造者史蒂芬·沃尔弗拉姆语）。学会放手，效果反而惊人。

孩子的大脑，也如一个刚刚初始化的大模型，有极大的参数空间等待优化。与其告诉他人生道理，不如给他精选的样例，让他通过自己的探索得到答案。这就是认知心理学家和教育心理学家杰罗姆·布鲁纳在其经典著作《教育过程》中提出的范例教学，又称归纳式教学。

在数学教学中，教师给出一系列完整解题步骤的例题，学生通过分析示例主动理解数学概念和方法，而不是教师直接讲解抽象的数学公式；在语文教学中，教师让学生通过反复接触大量语言样例归纳语法规则，而非直接灌输语法规则。这种方法不仅能加深理解，还更易于将其迁移到新的问题或情境中。

所以，孩子在成长过程中碰到的每一个难题，都不妨看作一次有意义的训练样例，父母无须立刻给出结论或答案，要让孩子自己去观察、体验、比较、反思，从中找到自己的道。放弃说教，“给予注意，学会陪伴”，这才是养育孩子的黄金法则。

成人也是如此。初等教育和高等教育赋予我们的道理如同预训练阶段的基础知识，它们在大脑中构建了认知的底层模型，却不足以直接指导我们应对真实复杂的生活场景。生活真正考验我们的是具体情境中的决策能力，而这种能力恰恰来自后续不

断的微调和强化学习。

例如，面对亲密关系中的冲突，书上说“要理解对方，包容不同观点”，但这样的抽象道理并不能让我们解决冲突；只有去倾听、去表达、去调节情绪，然后根据对方的反馈微调和优化我们“人际交往专家模块”的参数。所谓“纸上得来终觉浅，绝知此事要躬行”，这样，我们才不会陷入“懂得了很多道理，依旧过不好这一生”的局部最优陷阱。

第三，注意也是遗忘。学习的本质，是对知识体系的优化。大模型像一个捡破烂的拾荒者，无差别地记忆所有接触的信息。而人超越大模型的，是其所独有的“选择性遗忘”：有意识地强化对重要知识和场景的记忆，同时主动遗忘那些低效甚至有害的信息。所以，积极的遗忘并非失败，而是一种认知优化的策略，它可以让宝贵的注意力聚焦于那些真正有价值的信息和故事。《洛丽塔》的作者弗拉基米尔·纳博科夫说：“你所领悟的人生真理，皆是你曾付出代价的往事。”

在学习过程中，选择性遗忘就是“先做加法，再做减法”的思维模式。为策划一个项目，我们会收集大量的信息，做大量的调研，努力将各种可能性都纳入考虑范围。这是必要的第一步，即先做加法。越接近决策阶段，就越需要精准地做减法，选择性遗忘。比如，关于一款新产品，我们最初想法无数：既要满足市场需求，又要成本可控；既要功能强大，又要操作简单；既想满足年轻人的需求，又不愿放弃中年人市场。但是，真正的产品设计者，

要敢于主动“遗忘”那些充满吸引力但干扰产品核心定位的冗余信息，从而将注意力分配给真正的核心。著名设计师迪特·拉姆斯曾说：“好的设计不是堆砌更多的功能，而是敢于删去多余的东西。”遗忘，也是注意力分配的艺术。

生活中，我们有时会情绪低落，这可能是因为过去一些不愉快的经历：或许是一次失败的考试，一次刻骨铭心的分手，甚至是朋友无意中的伤害。

这些不愉快持续侵占和消耗着我们的注意力，不断地唤起痛苦的记忆，让我们陷入“身在当下，心在过去”的困境而无法自拔。选择性遗忘不是强迫忘记这些不愉快，或者逃避甚至否认它们曾经发生。选择性遗忘是承认，是接纳——承认它们

确实已经发生，无法更改，接纳它们曾给自己带来的伤害。但是需要明白的是，它们并不必然定义我们现在以及未来的人生。

心理学家卡尔·荣格说：“我们无法改变过去的事实，但我们可以改变看待这些事实的态度。”只有当我们真正接纳了这些痛苦的经历，允许自己放下情绪上的执着与执念，过去的负面经历才会与我们握手言和，逐渐淡去；唯有这样，注意力才会回归当下，回归我们能掌控的事情上。于是，我们重获内心的平静与自由。

遗忘，既是告别，也是起航。🌊

【本文配图由 AI 生成】

人超越大模型的，是其所独有的“选择性遗忘”

