

AI大模型技术分享

从基础原理到产品实践

分享人：黄康德

2025年12月12日

目录

1. 第一部分：理解AI时代的技术基础
2. 第二部分：提示词优化与上下文工程
3. 第三部分：AI产品开发新范式

第一部分：理解AI时代的技术基础

从底层原理到产品形态，建立完整的知识地图

大模型到底是什么

大模型本质上是一个超级大脑，通过阅读海量文本学会了理解和生成语言。2022年11月ChatGPT发布后，AI突然像人一样会聊天、写代码、解题，背后就是大模型的功劳。

核心原理非常简单：大模型的基础是"预测下一个词"。输入"今天天气很"，它会预测下一个词最可能是"好"、"差"、"不错"。这个预测不是查词典，而是基于看过的几千亿个词，学会了语言规律、知识和逻辑。

重要认知

大模型不是搜索引擎，不是数据库，它是一个**概率性语言生成器**。这意味着它可能"一本正经地胡说八道"（幻觉），也可能创造性地解决问题。

AI产品的三大类型

Chatbot（聊天机器人）

能力边界是"单次对话"。你问一句，它答一句，每轮对话相对独立。代表产品是最开始的ChatGPT网页版、Kimi网页版。适合问答、咨询、简单内容生成。

Copilot（智能助手）

能力边界是"人机协作"。AI在旁边辅助，人主导任务。GitHub Copilot（代码补全）、Office Copilot（文档生成）都是典型。特点是高频人机互动：AI生成，人审核修改。人是主驾驶，AI是副驾驶。

Agent（智能体）

能力边界是"自主执行"。你给AI一个目标（"帮我调研新能源汽车市场并写份报告"），AI自己拆解任务、搜索信息、分析数据、生成报告。人只负责验收结果，不需要每一步参与。Manus、Claude Code、AutoGPT都是Agent。

Agent是怎么工作的：从理解目标到完成任务的完整闭环

Agent不是简单的API调用，它是一个能自主思考、规划、执行、反思的系统。

Agent工作的四个核心步骤

第一步：理解目标并拆解任务 说"帮我做个竞品分析"，Agent会先理解目标，自动拆解成：1)确定竞品列表 2)收集每个竞品信息 3)整理对比表格 4)分析优劣势 5)写总结报告。这个过程叫"任务规划"，底层是模型的推理能力。

第二步：调用工具执行子任务 Agent知道自己能做什么、不能做什么。对于不能做的事，它会调用工具：搜索引擎查信息、计算器做数学运算、代码解释器运行代码、API查询数据库。工具调用能力让Agent突破了纯文本生成的限制。

第三步：验证结果并纠正错误 这是Agent最核心的能力。执行完子任务后，Agent会验证：结果合理吗？数据完整吗？有没有遗漏？发现错误会调整思路重新执行。比如搜索信息时发现结果太少，会换关键词重新搜索。

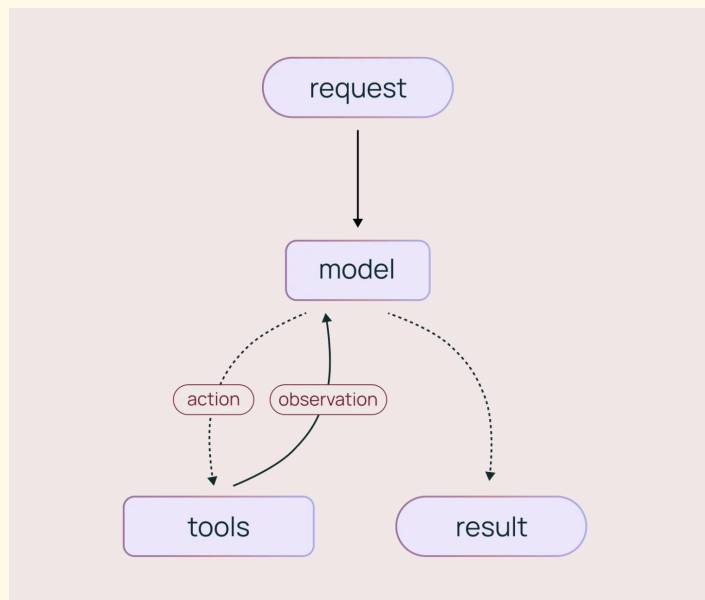
第四步：整合输出最终成果 所有子任务完成后，Agent把碎片信息整合成完整报告，可能还会反思：报告结构清晰吗？重点突出吗？用户会满意吗？有问题会优化后输出。

为什么Agent（智能体）是下一个金矿

Agent的关键特点是能调用工具，从而能和真实世界交互。

另一个核心变化是**反馈循环**，Agent可以借助工具实现反馈循环，AI自己验证、修正、迭代，减少一步错步步错的情况，从而能够独立完成复杂任务。

工具调用+反馈循环，是Agent能充分发挥大模型的能力，真正替代人工的关键。



Agent的技术支撑

关键技术

- **工具调用**：大模型被训练成能生成特殊的"工具调用指令"，系统收到后会调用相应API，再把结果返回给模型
- **MCP协议**（模型上下文协议） Agent和外部工具对话的统一标准。像USB-C接口，任何设备都能插。MCP让Agent能调用任何工具（搜索、数据库、API等），极大扩展能力边界
- **记忆系统**：Agent需要记住做了什么、结果是什么、下一步做什么。通过上下文窗口或外部数据库实现
- **模型能力**：强模型（GPT-5、Claude-4）能一次性完成复杂任务的规划和执行，具备可靠的工具调用能力

RAG是什么：大模型的"查资料"能力

RAG全称是Retrieval-Augmented Generation，中文翻译"检索增强生成"。也就是先查资料，再把查到的资料塞进提示词，让大模型根据资料回答问题。没有RAG时大模型是闭卷考试，只能凭记忆答题；有了RAG是开卷考试，可以现场查书。

RAG解决三个核心问题

第一，知识时效性 大模型的知识停留在训练数据截止时间（比如2024年4月），不知道最近发生的事。RAG让大模型能查最新信息，回答"今天的新闻"这种时效性问题。

第二，领域专业性 大模型是通才不是专才，不懂你公司的内部规章制度、产品手册、技术文档。RAG把这些专业资料塞进提示词，大模型就能基于这些资料回答专业问题。AI客服、企业内部知识库都是这个原理。

第三，回答可控性 大模型容易"一本正经胡说八道"（幻觉），RAG强制它只能基于提供的资料回答，显著降低胡说八道的概率。这对企业应用至关重要——让AI的回答有据可查，不是凭空编的。

RAG的实现

检索质量是前提

RAG的核心是"检索", 检索系统可以是搜索引擎、向量数据库、关系型数据库、API接口, 只要是"先查后答"都算RAG。但如果检索回来的资料不准、不全、不相关, 大模型再强也答不好。

实现RAG的三要素

RAG是一个复杂的检索-理解-生成系统。真正有效的RAG需要: 精准的**检索算法** (找到相关的信息)、高质量的**知识库** (资料准确完整)、**提示词**工程 (让大模型正确使用资料)。三缺其一, 效果都会大打折扣。

实现方式

最简单的RAG实现 调用大模型之前, 先根据问题检索相关信息, 在提示词里加上检索到的信息, 让大模型根据这些信息回答问题, 这种方式简单有效, 能解决很多问题。

效果更好的方式 用Agent实现, Agent可以调用工具检索信息, 如果信息不够, Agent可以自主决定检索更多信息, 直到信息足够完成回答。这种多轮检索更像人类的行为, 可以应对更复杂的问题。

模型即Agent (2025)

趋势：大模型API不只是文本生成，而是会思考、能调工具的Agent系统

思考：模型能展示完整"思考过程"，使推理透明化 新一代大模型具备"推理链"输出能力。面对复杂问题，它会像人类一样先在内部推演步骤（例如展示数学题的解题过程），再给出最终答案。这不仅提升了结果的可信度，更让开发者能够观察、调试模型的"思路"，为构建可靠Agent提供了关键基础。

工具：模型能自主调用工具，从"思考"走向"行动" 模型清楚自身局限（如知识过时、无法计算），因此学会了在需要时生成指令，调用搜索引擎、代码解释器或业务API等工具。这使其突破了纯文本的界限，能够处理实时、动态、专业的任务，真正成为能连接现实世界的"行动者"。

进化：工具调用进入"并行"与"穿插"模式，大幅提升执行速度 早期工具调用是线性的：思考-调用-等待。现在，模型支持：**(1) 并行调用**：一次性发起多个无依赖关系的工具请求（如同时读多个文件）**(2) 思考中调用（边想边做）**：在推理中途随时插入工具调用获取关键信息。这极大优化了复杂任务的执行效率。

大模型全景图（2025）

大模型分类

维度	类别	特点与适用场景
尺寸	小模型 (<10B)	响应快、成本低，轻量任务，手机、IoT设备端侧部署
	中模型 (10-100B)	平衡性能与效率，企业通用场景
	大模型 (>100B)	智能程度高，复杂任务，MoE架构激活参数约10-70B
模态	纯文本模型	LLM基础形态，对话问答，Agent应用
	多模态模型	支持图像、音频、视频，跨模态理解与生成

国外主流大模型

模型

所属公司

核心优势

GPT-5

OpenAI

能力强大且均衡，生态成熟

Claude 4.5

Anthropic

编程能力顶尖，逻辑推理强

Gemini 3 Pro

Google

多模态能力全面，性价比高

Grok 4

xAI

内容过滤很少，响应速度快

国内主流大模型

模型	类型	特点
豆包 Seed	闭源	综合能力强，多模态齐全
千问 Qwen	开源/闭源	开源生态活跃，多模态齐全
文心 ERNIE	开源/闭源	中文写作优势，支持多模态
腾讯 混元	开源/闭源	多模态和3D能力强
DeepSeek	开源	数学推理顶尖，综合能力强
Kimi K2	开源	长文本、Agent/编程能力强，兼容Claude

专业视觉模型（专注文档/特定场景）

- **MinerU 2.5** (上海AI Lab) - SOTA级文档解析, 1.2B参数
- **PaddleOCR-VL 3.0** (百度) - 顶尖OCR模型, 0.9B参数
- **MiniCPM-V 4.5** (面壁智能) - 端侧最强视觉理解, 8B参数

大模型的能力边界 - AI擅长的

文本处理与生成

- 内容创作：文案、文章、代码、创意写作
- 语言转换：翻译、改写、风格调整、格式转换
- 信息提取：总结、关键词提取、结构化整理

知识问答与推理

- 基础推理：逻辑分析、因果关系推断
- 知识整合：跨领域信息综合、概念解释

代码相关任务

- 代码生成：根据需求编写各种语言的代码
- 调试协助：错误分析、优化建议

产品设计启示：这些是大模型的"舒适区"，可以作为产品的核心功能点。

大模型的能力边界 - AI难做好的

1. 需要"绝对精确性"与"可验证性"的任务 - 大模型的"幻觉"会导致致命错误。

大模型输出是概率性生成（而非逻辑推导），即使"看起来正确"，也可能包含编造的事实（虚构法条、假数学公式、错误代码）。这类任务的核心要求是"100%准确"，一旦出错会引发严重后果，因此大模型只能做"辅助建议"，不能作主。

2. 需要"高隐私保护"的任务 - 大模型的"数据泄露风险"不可控。

大模型的推理通常需要将数据发送到云端（除非是本地部署的小模型，但小模型能力有限），用户的敏感数据（医疗记录、财务信息）一旦上传，可能因"模型记忆"或"数据泄露"被滥用。

3. 需要"低成本、高频次"的简单任务 - 大模型的"性价比"太低。

大模型的推理成本很高（GPT-4调用成本约0.03美元/1k tokens，传统API可能只需0.001美元/次），对于"简单、重复性"任务，用大模型是资源浪费。

第二部分：提示词优化与上下文工程

一套让AI更听话的实战方法

提示词工程的进化史

第一代：原始提示时期（2021-2022）

GPT-3刚发布，AI能力弱，经常答非所问。核心挑战是输出不稳定、逻辑混乱。提示词策略是越简单越好、避免歧义。ChatGPT问世后，AI能理解自然语言，但输出仍不稳定。

规律：AI能力弱，提示词优化无从谈起。

第二代：结构化提示时期（2022-2023）

2022年初"Chain-of-Thought"（思维链）论文发布，突破是让AI"一步一步思考"，大幅提升复杂任务准确率。产生了角色设定、任务拆解、示例参考等方法。第一批"提示词工程师"岗位出现。

规律：提示词优化变成可复用的技巧。

提示词工程的进化史

第三代：系统化提示时期（2023-2024）

GPT-4发布，AI理解能力质变，能识别多种指令模式。催生了ICIO、CRISPE等框架。提示词升级为系统化的"工程学科"。

规律：AI变聪明，能读懂复杂的框架结构。

第四代：上下文工程时期（2024-2025）

Claude 3.5、GPT-4o支持20万tokens上下文，理念彻底转变：从"琢磨完美提示词"转向"提供充分上下文"。旧思维是"出考题"，新思维是"给教科书"。企业"知识库"成为核心竞争力。

规律：AI足够强，上下文足够大，不再依赖于提示词技巧。

理解AI的脑回路：7个底层原理

原理1：概率预测本质

AI输出"最可能"而非"最正确"。大模型本质是"预测下一个词的概率分布"，不是"寻找正确答案"。AI会输出"看起来最合理"而非"事实最准确"的答案，这就是AI"一本正经胡说八道"（幻觉）的根本原因。

应对策略：

- 不要问AI不知道的事（它会瞎编）
- 给AI提供事实锚点（把相关资料给它）
- 要求AI标注引用（"让它提供回答依据"）

理解AI的脑回路：7个底层原理

原理2：注意力机制

AI的"眼睛"会疲劳。一次只能关注有限信息，有"首尾效应"：开头结尾最容易记住；重复信息会被加强；特殊标记更显眼。

实践技巧：

- 重要信息放开头或结尾
- 重复关键指令2-3次（如"记住，必须输出JSON格式"）
- 用 `<task></task>` 、 `<output></output>` 等标签分隔信息

理解AI的脑回路：7个底层原理

原理3：少样本学习

给AI看例子是最好的教学。AI擅长模仿，给2-3个例子就能学会规律。要点：例子要覆盖主要场景和边界情况；例子格式必须与任务一致。

原理4：思维链

让AI慢思考。AI的"潜意识"反应快但易出错，强迫它"一步一步想，给出思考过程"，准确率大幅提升，这就是"Chain-of-Thought"（思维链）。现在的推理模型（如deepseek-r1）本身就支持深度思考，可通过参数控制思考长度。

理解AI的脑回路：7个底层原理

原理5：角色代入

AI是个演员。AI会根据设定的角色和语境调整语气、知识侧重点和回答风格。你问"什么是AI"，小学生与大学教授的说法会截然不同。

原理6：SOP

把复杂规则变成流程。复杂提示词涉及过多规则和逻辑，AI难以遵循。改造成SOP，按流程办事就简单多了。SOP帮助人类用户明确传达期望，也让AI应用更流程化、规范化。

原理7：结构化数据格式

大模型学习了海量JSON数据，支持JSON Mode确保格式正确，可以用JSON Schema描述输出结构。同时大模型也擅长识别XML标签，如 `<tag></tag>`，可自定义、结构化、嵌套，用成对标签标记上下文能提高准确性，避免歧义。

ICIO框架（最通用的基础框架）

ICIO = Instruction（指令）+ Context（上下文）+ Input Data（输入数据）+ Output Indicator（输出指示）

结构模板

```
<instruction>【明确指令】请完成什么任务？一句话说清楚</instruction>
```

```
<context>【背景】为什么要做这个？相关背景是什么？</context>
```

```
<input_data>【输入数据】需要处理的具体内容</input_data>
```

```
<output_indicator>【输出要求】什么格式？什么风格？多长？有什么限制？</output_indicator>
```

ICIO框架实战案例

用户评论情感分析

```
<instruction>分析以下用户评论的情感倾向，分类为正面、中性或负面</instruction>
```

```
<context>这些评论来自我们新上线APP的用户反馈，将用于改进产品体验</context>
```

```
<input_data>
```

评论1：这个APP太难用了，到处都是BUG！

评论2：界面挺简洁的，但功能太少

评论3：非常流畅，比竞品好太多了！

```
</input_data>
```

```
<output_indicator>输出JSON格式：{"评论1":"负面",...}。不要添加任何解释或额外文字</output_indicator>
```

精髓：

指令清晰（AI知道要做什么），背景充分（知道数据用途，会更准确），输入明确（知道要处理哪些内容），输出可控（知道格式要求，不会输出废话）。

CRISPE框架（适合创意和内容生成）

CRISPE = Capacity and Role（角色）+ Insight（洞察）+ Statement（陈述）+ Personality（性格）+ Experiment（实验）

这个框架适合写文章、做方案、创意输出等需要"思考深度"的任务。

结构模板

`<capacity_and_role>`你是一个什么角色？具备什么能力？`</capacity_and_role>`

`<insight>`提供用户不知道的洞察、内幕、背景信息`</insight>`

`<statement>`需要做什么？明确陈述任务`</statement>`

`<personality>`用什么语气风格？（简洁、幽默、专业、严谨...）`</personality>`

`<experiment>`要求AI提供多个选项或不同角度的答案`</experiment>`

CRISPE框架实战案例

撰写产品推广文案

<capacity_and_role>你是一位有15年经验的知名品牌营销专家，在快消品行业有丰富经验</capacity_and_role>

<insight>我们这款产品是"无糖气泡水"，目标用户是25-35岁都市白领，关心健康但也追求生活品质。竞品主要有元气森林、喜茶气泡水。我们产品的独特优势是：采用天然果汁调味（不是人工香精）</insight>

<statement>撰写一篇小红书风格的推广文案，突出健康+美味的双重卖点</statement>

<personality>语气轻松活泼，像朋友在聊天，适当使用emoji和网络流行语</personality>

<experiment>提供3个不同风格的版本：版本1强调健康，版本2强调口感，版本3强调生活方式</experiment>

效果：

CRISPE框架产出的内容有深度、有洞察、有多样性，不只是机械地完成任务。

六个可直接套用的提示词模板

模板1：让AI问清楚需求

任务：帮我做XXX

先复述一遍你的理解和规划，有不懂的和我沟通。

模板2：让AI主动纠错

任务：帮我写一段Python代码实现XXX功能

写完后，请：

1. 仔细review代码，不要有任何坏味道
2. 确保测试100%通过，不要有任何bug

六个可直接套用的提示词模板

模板3： 专家会诊

请从3个不同角色的视角分析这个问题：

角色1：技术专家（关注可行性）

角色2：产品经理（关注用户价值）

角色3：运营人员（关注推广和成本）

问题：[你的问题]

每个角色给出分析和建议，最后总结共识和分歧。

模板4： 对抗性思维

任务：审查以下方案的风险

请你扮演"魔鬼代言人"，专门挑毛病：

1. 找出5个可能失败的原因
2. 找出3个被忽略的关键问题
3. 找出2个过度乐观的假设

方案：[你的方案]

六个可直接套用的提示词模板

模板5：费曼学习法

主题：[你想理解的概念]

请用费曼学习法的三个层次解释：

1. 5岁小孩能懂的版本（用比喻和简单语言）
2. 高中生的版本（有专业术语但解释清楚）
3. 专家的版本（深入技术细节）

模板6：第一性原理

问题：[你要解决的问题]

请用第一性原理分析：

1. 这个问题的本质是什么？回到最基本的事实
2. 哪些是大家默认假设？这些假设一定对吗？
3. 抛开现有方案，从0开始会怎么设计？
4. 现有方案的哪些部分可以优化？

上下文工程：从"一句话"到"一本书"

上下文工程是提示词工程的进化版。不再关注"怎么写提示词"，而是关注"怎么给AI提供**充分且必要**的上下文信息"。

传统做法是200-2000字提示词，用户需要是"提示词专家"。新做法是20万字全套资料（文档、数据、历史对话、规范、案例），Agent从上下文中自主学习，用户不需要是专家。

比喻

提示词工程=考官出考题，要出得巧妙才能考出学生水平；

上下文工程=给学生一本教科书和工作手册，学生自己看书就能完成任务。

上下文的要素：构建充分且必要的信息

要素1：任务上下文——这次要解决什么具体问题

包括：目标（SMART原则：具体、可衡量、可实现、相关性、有时限）、范围、约束、验收标准。

```
<task_context>
```

项目名称：新版App用户注册流程优化

目标：3个月内将注册转化率从45%提升到60%

范围：只优化手机号注册，不包括第三方登录

约束：保持现有品牌形象、加载时间不超过3秒、符合隐私法规

验收标准：A/B测试显著性>95%，样本量>10000

```
</task_context>
```

上下文的要素

要素2：背景知识上下文——相关的领域知识、业务逻辑、行业规范

包括：产品文档、技术文档、用户手册、行业报告。关键原则是相关的、结构化的、最新的。

```
<knowledge_context>
```

产品文档：

- 产品定位：面向25-35岁都市白领的健康管理App
- 核心价值：AI个性化饮食建议
- 用户画像：注重健康但时间有限，愿意为便利付费

技术规范：

- 技术栈：React Native + Node.js
- 设计系统：使用Material Design
- API标准：RESTful，返回JSON格式

```
</knowledge_context>
```

上下文的要素

要素3：历史经验上下文——过去成功或失败的经验、案例、数据

包括：历史项目文档、会议记录、用户反馈、测试数据。关键原则是真实、具体、可复用。

```
<history_context>
```

成功案例：

- 2024年3月：简化注册流程（从5步减到3步），转化率提升12%
- 关键做法：移除邮箱验证（只保留手机号），使用一键填充

失败教训：

- 2024年1月：增加更多注册项（包括职业、收入），转化率下降18%
- 原因分析：用户觉得隐私泄露风险大，放弃注册

用户反馈（来自客服记录）：

- "希望注册更快一点，现在填的东西太多了"
- "担心手机号会泄露，收到骚扰电话"

```
</history_context>
```

决定AI效果的关键因素：大模型+上下文+提示词

大模型：聪明的大脑

大模型系统的大脑，决定AI效果的基础。目前顶尖的模型对提示词的理解和遵循非常好，对上下文的理解和工具使用也更好，而比较弱的模型（例如私有化部署的小模型、中等模型）无法理解太复杂的提示词和上下文。

上下文：完备的地图

AI回答问题或执行任务，通常需要很多背景知识和任务上下文，RAG用提示词变量给AI提供上下文，Agent则主动调用工具获取上下文。例如AI写代码之前需要理解已有的代码，回答问题需要了解相关的业务知识。好的上下文就像一张完备的地图，让AI能更快到达目的地。

提示词：清晰的需求

现阶段模型能力足够强，只要把事情交代清楚，AI就能很好的理解和遵循。但是，需求不明确是经常发生的事情，需要和AI多讨论沟通，确保AI的理解和你实际需求是一致的。好的提示词就是清晰的需求。

第三部分：AI产品开发新范式

从智能的进化到商业价值重构

大模型的智能进化：指数级增长才刚刚开始

核心观察

各大厂商一边喊瓶颈，一边持续刷新记录。OpenAI、Google、Anthropic、DeepSeek、月之暗面在2024-2025年不断推出更强模型。从GPT-3到GPT-5用了不到4年，每个版本都是质的飞跃。现在的模型已经能通过图灵测试、完成复杂推理、写代码、做科研。

行业趋势

不要相信"大模型遇到瓶颈"的说法。全球顶尖的AI人才都集中在少数几家头部公司，大量的未公开研究成果在持续突破。公开论文显示的技术已经落后6-12个月。从2018年BERT的3亿参数，到2025年DeepSeek V3突破6000亿参数，增长曲线仍在加速。

对AI从业者的启示

短期（1-2年）

不要怀疑大模型会继续变聪明。现在讨论的很多问题（幻觉、推理能力、成本），很可能在2026-2027年通过新架构被解决。你的产品架构要预留升级空间，不要把当前模型的局限当作永久性问题的设计解决方案。

长期（5-10年）

大模型终将全面超越人类通用智能，只是时间问题。可能不是5年，可能是10-15年，但方向是确定的。这意味着我们现在做的所有产品，未来都会被更强大的模型重构。不要把产品价值建立在"模型不够强"的假设上，而要建立"模型会越来越强，我怎么利用它"的思路。

大模型的价格雪崩：开源模型主导的成本革命

市场现实

大模型API价格以每月10-30%的速度下降。GPT-4刚发布时，每百万tokens要60美元，现在降到了个位数。开源模型更是把价格打到谷底：Qwen、DeepSeek的性能已经接近顶级闭源模型，而部署成本只有1/10。

价格雪崩的三大原因

1. **市场竞争白热化**：模型厂商疯狂的价格战，API价格无限接近成本。这不是烧钱补贴，而是算力成本确实在快速下降
2. **技术进步压缩成本**：模型架构优化、量化技术、推理加速让单次推理成本降低90%+
3. **开源模型崛起**：DeepSeek、Qwen等开源模型性能已接近GPT-4，算力供应商都可以部署，没有模型研发成本

价格雪崩对AI应用开发者的影响

根本性影响

成本不再是核心制约因素。过去调用GPT-4会心疼钱，现在你可以大胆地给AI塞10万字的上下文，让它做复杂任务。开发者应该把注意力放在"怎么用好大模型"而不是"怎么省大模型的钱"。

长期趋势

大模型最终会成为像水、电、网络一样的廉价基础设施。但要注意的是，顶级闭源模型（GPT-5、Claude-4）可能仍然保持高价，因为它们代表最前沿能力。未来的市场可能是"开源模型占80%（通用需求），闭源模型占20%（最高端需求）"的双层结构。

从提示词工程到上下文工程：AI应用设计范式的转变

核心转变

大模型变强后，提示词优化没那么重要了。2023年你需要研究各种提示词技巧：Chain-of-Thought、Few-Shot、角色设定、思维链，甚至花钱上李一舟的课。现在，只要你把需求说清楚，AI基本都能理解。

更重要的是"上下文工程"

不是怎么写提示词，而是怎么给AI提供"充分且必要"的背景信息。

为什么上下文比提示词重要？

强模型是一点就通的，你不需要教它怎么思考，只需要告诉它"你是谁、要解决什么问题、有什么约束条件、历史对话是什么"。比如让AI做一个用户运营方案，与其琢磨提示词技巧，不如把"用户画像数据、历史活动效果、预算限制、渠道资源"这些上下文塞给它。

弱模型和强模型的用法差异：弱模型（比如7B模型）脑子转不过来，任务要简单，一次只做一件事；强模型（比如GPT-5、Claude 4）可以一次性给复杂指令，甚至多个任务并行。

AI产品的核心价值：10倍效率提升是临界点

AI产品值钱的根本原因：它能带来10倍甚至100倍的效率提升。而且这种提升是跳跃式，不是渐进式。没有这个量级，用户不会改变自己的习惯。

效率提升的三个维度

第一，时间压缩 以前写100行代码要1小时，现在AI写完你改15分钟，效率提升4倍；以前分析1小时会议录屏要2小时人工听，现在AI自动总结10分钟，效率提升12倍；以前读10份竞品报告要3天，现在AI提炼要点30分钟，效率提升50倍。时间压缩是用户最能感知到的价值。

第二，能力平民化 以前只有资深程序员能写复杂算法，现在初级程序员用AI也能写；以前只有专业设计师能做品牌设计，现在运营同学AI生成就能用；以前只有数据分析师能做用户行为分析，现在产品经理用AI也能做。AI让几十年积累的专业技能，一夜之间"贬值"。这对中小企业是重大利好。

第三，成本重构

以前开发一个App需要5个人3个月，现在1个人用AI编程工具1个月就能做（成本从50万降到5万）；以前需要雇佣客服团队24小时轮班，现在AI客服成本是人力成本的1/10；以前需要10个律师助理审合同，现在AI审一遍，律师看重点就行（成本从100万降到20万）。成本重构直接改变了商业模式。

为什么10倍效率是临界点？

3倍提升：大家都会用AI，但不会改变核心工作流程（原来怎么干现在还怎么干，只是快一点）

10倍提升：工作流程必须重构，AI成为核心工具，人变成AI的管理者

100倍提升：整个行业被重塑，旧岗位消失，新岗位诞生，商业模式彻底变了

举个例子：为什么AI编程这么火？

因为它同时满足三个维度：

- **时间**上从小时级降到分钟级（10倍提升）
- **能力**上初级程序员能写出高级程序员水平的代码（能力民主化）
- **成本**上开发成本从百万级降到十万级（成本重构）

这种三重效率革命的产品，就是最大的机会。

AI编程的范式转移：从技术实现者到需求架构师

旧范式（2023年以前）

主要人写代码，AI辅助生成碎片代码。AI只能完成函数级别的任务，需要频繁人机协作，人得盯着AI输出，AI不懂整体架构，经常出错需要纠正。人既是架构师又是码农，AI只是个智能输入法。

新范式（2025年）

人做需求和方案沟通，AI负责实现。AI能独立完成复杂功能模块，甚至整个项目。人更像技术管理者：拆解需求、定义接口、验收成果。AI能自主Debug、优化、迭代，不需要每一步都干预。

对人的能力要求根本性改变

Cursor/Claude Code的使用体验:

1. 你告诉AI: "我要做一个用户注册功能, 包括手机号验证、邀请码系统、防刷机制"
2. AI自动规划: 设计数据库表、写API接口、实现短信服务集成、做Rate Limiting
3. 10分钟后, 完整的功能交付, 包括单元测试
4. 你负责Review代码、做Code Review、确保符合架构规范

程序员的价值从"写代码"转向"指导AI写代码":

- 不需要记具体语法, 但要懂架构设计
- 不需要手写每一行代码, 但要能判断代码质量和潜在风险
- 沟通能力和需求拆解能力, 比编码能力更重要

软件价值的重构：从代码资产到解决方案的回归

AI让软件开发成本降低了10倍，这对软件行业意味着什么？不是简单的降价，而是价值逻辑的重新定位。

短期内会出现的现象

- 大量同质化软件爆发：ChatBot、AI助手、低代码平台，千篇一律
- 价格战：因为开发成本低，大家都打价格战，利润被压缩到极限
- 开源项目爆炸：GitHub上每天新增AI项目超过1000个，选择困难症加重

软件价值的重构

长期来看，价值会回归到正确的地方

代码不再是护城河 复制一个Cursor可能只需要3个月（开源社区已经有几十替代品），技术本身不能形成壁垒

真正的价值在"理解用户问题"的能力 为什么同样用AI编程，有人能做出爆款，有人做不出来？差别在于对用户痛点的理解深度。技术平权后，洞察力成为稀缺资源

交付完整解决方案的能力 客户要的不是代码，而是能解决业务问题的完整服务，包括部署、运维、培训、定制开发。这些都是AI替不了的

软件行业的范式转移

从"卖代码"转向"卖解决方案和服务"

从"产品公司"转向"咨询公司+产品公司"的混合模式。纯SaaS模式可能被"SaaS+服务"模式取代。

定价逻辑的改变

过去按功能收费（功能越多越贵），未来可能按"AI替代了多少人工"收费（替代价值越大越贵）。

对创业者的建议

如果你有行业Know-How（懂某个垂直领域），现在是用AI放大价值的最好时机。你不需要组建庞大的技术团队，1-2个工程师+AI工具就能做出之前需要20人团队才能做的产品。关键是把行业理解转化为产品能力。

开源碾压一切：速度决定生死的新规则

开源软件已经多到用不完，大厂全部下场，更新频率极高。

这不是开源社区的小打小闹，而是商业规则的改变。

几个震撼的事实

- Manus（号称第一个通用Agent）还没正式发布，GitHub上已经有10+开源替代方案
- Cursor刚火3个月，VSCode开源社区就出了Cline、Aider、Continue等十几个替代品
- Claude Code本身就是开源的，但开源社区又针对它做了优化版

对创业者的启示

- **不要重复造轮子**：你想做的功能，开源社区可能已经做好了，直接用就行。不要自己研发，浪费时间。
- **套壳是最佳选择**：在技术开源、更新飞快的环境下，只有套壳才能跟上节奏。自研底层技术永远追不上开源社区的速度。
- **价值在服务层**：既然技术开源，你的竞争力不在代码，而在部署、运维、用户支持、行业Know-How。客户愿意为"能用起来"付费，而不是"技术先进"付费。

开源 vs 商业的新平衡点：未来可能是"开源解决80%的通用需求，商业公司解决20%的定制和服务"。商业公司的价值不在技术独占，而在把开源技术包装成稳定、可靠、易用的产品。

警惕开源的陷阱：不是所有开源项目都能存活。90%的开源项目半年后就没人维护了。选择开源方案时，要看社区活跃度、Star数、贡献者数量、大厂背书。不要选太小众的。

"套壳"的价值重估：产品化能力才是真护城河

业界说法："AI应用就是套壳"，把套壳做到极致就是牛逼。这句话带点贬义，但暗含深刻真理。

什么是真正的套壳？

OpenAI的API是核弹头，但你不能直接把核弹头卖给用户，你需要导弹、发射系统、瞄准器——这就是套壳。套壳是把复杂的AI能力包装成普通用户能用的产品，解决"最后一公里"的用户体验问题。

套壳的三个层次（从浅到深）：

1. 对话层套壳：最早的套壳，做聊天界面，写提示词模板。ChatGPT刚出来那会的各种套壳网站，现在基本死光了。壁垒最低，最容易被取代
2. 工作流套壳：把AI嵌入业务流程，像Coze、Dify、n8n这种，把AI和API、数据库、Webhook连接起来。壁垒中等，需要理解业务流程
3. Agent套壳：像Manus、Claude Code、Cursor，不是单次调用AI，而是让AI自主规划、执行、调试，完成复杂任务。壁垒最高，需要系统架构能力

套壳的本质和启示

为什么套壳比底层技术更难、更有价值？

- 技术开源了，大家都会用，但知道"用在哪、怎么用、怎么让用户爽"是另一回事
- 底层模型每3个月迭代一次，你的产品架构能不能快速跟上？

关键认知：套壳的本质是"产品化能力"。这不是贬义，而是AI时代最稀缺的能力。全世界都在研究模型，但全世界懂用户、懂场景、懂产品的人还是少数。技术可以复制，但对用户的理解、对场景的洞察、对细节的追求，无法复制。

对创业者的启示：不要自责"我们只是套壳"。承认自己是套壳，并把套壳做到极致，这是正道。与其花1000万美金训练一个中等模型，不如花100万美金把用户体验做到90分。用户不关心你是不是套壳，用户关心好不好用。

AI产品的生存法则：绑定在指数级增长的曲线上

大模型及其应用生态是指数级进化的，像火箭一样快速上升。

核心原则

- **模型会变强**：你的产品价值应当随之自动提升
- **成本会下降**：你的服务成本随之下降，利润自然提升
- **生态会进化**：你的技术栈应当建立在活跃的社区之上

如何正确"坐火箭"?

✅ 做这些（坐火箭）:

- 抽象掉具体模型，随时替换
- 选择最活跃的开源社区和生态
- 轻量开发，快速验证价值

❌ 别做这些（造轮子）:

- 围绕当前模型能力上限设计复杂系统
- 自己实现已经有成熟方案的通用模块
- 过早优化和过度抽象，"完美的抽象"会变成最重的包袱

数据是数字石油：但核心在于炼化能力

公开数据+持续更新的爬虫数据，是互联网时代最值钱的资产之一。但原始数据只是原油，价值有限。

为什么数据像石油？

- **原始数据=原油**，值钱但价值有限。大家都会爬，没有门槛
- **清洗、结构化、关联后的数据=成品油**，价值指数级提升。炼油需要技术、工艺、时间
- **时序数据+关联数据+真实验证数据=高纯度化学品**，价值连城。需要深度加工和提纯

几个赚钱的例子：

- **天眼查/企查查**：爬取全国工商数据、法院判例、招投标信息，清洗关联后卖给B端用户，年费几千到几万不等。核心壁垒不是爬虫，而是数据清洗和关联算法
- **5118/爱站**：爬取搜索引擎关键词数据，做SEO分析工具。客单价虽然不高但用户量大。
- **新榜/西瓜数据**：爬取抖音、公众号、B站数据，做榜单和数据分析，广告主和品牌方买单。

数据的深层价值

1. 时序价值 不只是静态快照，而是变化趋势。比如"某公司过去3年的招聘趋势"比"该公司现在的员工数量"更有价值。趋势能预测未来。

2. 关联价值 单一数据源容易被复制，但多维度关联很难。比如"把招聘数据+工商变更+融资信息+专利信息关联起来，预测哪家公司要扩张"。这些关联需要行业Know-How。

3. 真实性价值 AI生成内容泛滥后，能验证数据真实性的源数据会更值钱。比如"哪些内容是真人发的，哪些是AI批量生成的"，这种验证数据有巨大价值。

数据不是护城河，"炼化能力"才是：数据采集谁都会，但清洗规则、打标体系、关联逻辑、更新机制，这些才是核心技术。一个爬虫工程师可能花3个月就能爬全所有工商数据，但要做出天眼查那样的产品，还需要3年的数据清洗和关联优化。

对创业者的建议：如果你有数据采集+清洗的能力，现在是用AI放大价值的最好时机。你可以爬公开数据，用AI自动清洗、分类、打标、生成关联，价值。

风险警告：数据合规风险越来越大。**爬虫有可能违法**，特别是抓取个人信息、竞品数据。一定要注意合规，最好只抓公开数据，或者和客户合作，抓客户授权的数据。

谢谢聆听

Q & A

AI大模型技术分享