

# LangChain 项目主体结构与深度解析（基于当前仓库实际代码）

## 说明

本文档面向“读代码/二次开发/集成落地”的工程视角，对本仓库的包划分、真实入口（public API）、关键运行机制（尤其是 langchain\_v1 的 Agent 图式执行）、中间件体系、以及集成生态做系统化说明。

本文档只描述仓库结构与逻辑，不包含任何训练相关内容，也不包含需要在 Linux 服务器执行的命令。

## 一、仓库结构树（校正后的关键视图）

```
langchain-master/
|
├── README.md          # 项目主文档 (生态介绍)
├── LICENSE            # 许可证
├── AGENTS.md          # monorepo 开发规范/贡献指引 (面向开发者)
├── CLAUDE.md          # monorepo 开发规范/贡献指引 (面向开发者)
└── CITATION.cff       # 引用信息
|
└── libs/              # Python monorepo: 多个独立发布包 (每包独立
    pyproject.toml/uv.lock)
|
    ├── core/             # langchain-core: 基础抽象与协议 (不含第三方集成)
    │   ├── pyproject.toml
    │   ├── uv.lock
    │   └── langchain_core/
    │       ├── __init__.py      # 版本/弃用警告 surface (公共入口很薄)
    │       ├── agents.py       # 旧式 Agent schema (主要用于兼容)
    │       ├── caches.py       # 缓存抽象
    │       ├── env.py          # 环境读取工具
    │       ├── exceptions.py    # 异常
    │       ├── globals.py       # 全局开关 (debug/verbose 等)
    │       ├── prompt_values.py # PromptValue 体系
    │       ├── rate_limiters.py # 速率限制抽象
    │       ├── retrievers.py    # Retriever 抽象
    │       ├── stores.py         # Store 抽象
    │       ├── structured_query.py # 结构化查询抽象/工具
    │       ├── sys_info.py      # 系统信息
    │       └── version.py
    |
    ├── _api/               # 内部 API: beta/弃用/内部路径工具
    ├── callbacks/          # callbacks: 回调与 run manager
    ├── tracers/            # tracers: 运行追踪 (Console/LangSmith 等)
    ├── language_models/     # BaseChatModel/BaseLLM 等抽象
    ├── messages/           # BaseMessage/AIMessage/ToolMessage 等
    └── prompts/            # PromptTemplate/ChatPromptTemplate 等
```

```
|   └── tools/      # BaseTool 与工具定义/转换
|   └── runnables/    # Runnable 协议与 LCEL (链式组合)
|   └── output_parsers/ # 输出解析器 (JSON/Pydantic 等)
|   └── outputs/     # LLMResult/ChatResult 等结果对象
|   └── documents/    # Document 抽象与转换
|   └── document_loaders/ # Loader 抽象
|   └── embeddings/    # Embeddings 抽象
|   └── vectorstores/   # VectorStore 抽象
|   └── indexing/      # 索引 API 抽象与内存实现
|   └── load/          # 序列化/反序列化协议 (Serializable 等)

|── langchain_v1/      # langchain: 当前活跃维护的“主包实现层”
|   ├── pyproject.toml
|   ├── uv.lock
|   └── langchain/
|       ├── init.py      # 仅版本号 (入口很薄)

|       ├── agents/      # v1 Agent: 基于 LangGraph 的图式执行
|           ├── init.py    # 对外导出 create_agent/AgentState
|           ├── factory.py  # create_agent: 编译 StateGraph (核心入口)
|           └── structured_output.py
|               └── middleware/ # middleware: 在图节点/模型调用/工具调用插入逻辑
|                   ├── types.py   # AgentState/ModelRequest/ModelResponse + decorators
|                   ├── _execution.py # shell middleware 执行策略 (Host/Docker/CodexSandbox)
|                   ├── _redaction.py
|                   ├── _retry.py
|                   ├── context_editing.py
|                   ├── file_search.py
|                   ├── human_in_the_loop.py
|                   ├── model_call_limit.py
|                   ├── model_fallback.py
|                   ├── model_retry.py
|                   ├── pii.py
|                   ├── shell_tool.py
|                   ├── summarization.py
|                   ├── todo.py
|                   ├── tool_call_limit.py
|                   ├── tool_emulator.py
|                   ├── tool_retry.py
|                   └── tool_selection.py

|       ├── chat_models/   # init_chat_model: 统一初始化各 provider chat model
|           ├── init.py
|           └── base.py

|       ├── embeddings/    # embeddings 统一入口 (简化层)
|       ├── tools/          # 注意: 此处 tool_node.py 为兼容 re-export (非自研 ToolNode)
|       ├── messages/
|       └── rate_limiters/
```

```

├── langchain/          # langchain-classic: 遗留实现层 (兼容为主)
│   ├── pyproject.toml
│   ├── uv.lock
│   └── langchain_classic/
│       ├── init.py      # 通过 getattr 做旧入口兼容与弃用提示
│       ├── agents/      # 旧式 agents/chains/tools 等大量实现
│       ├── chains/
│       ├── chat_models/
│       ├── llms/
│       ├── embeddings/
│       ├── vectorstores/
│       ├── document_loaders/
│       ├── tools/
│       └── ...           # 规模大, 更多为历史兼容
|
├── partners/          # 官方维护的一小部分第三方集成包 (每个目录独立发布)
│   ├── openai/          # langchain-openai
│   ├── anthropic/       # langchain-anthropic
│   ├── huggingface/     # langchain-huggingface
│   ├── ollama/          # langchain-ollama
│   ├── deepseek/        # langchain-deepseek
│   ├── xai/             # langchain-xai
│   ├── chroma/          # langchain-chroma
│   ├── qdrant/          # langchain-qdrant
│   └── ...              # 还有部分 provider (但并不覆盖全生态)
|
├── text-splitters/    # langchain-text-splitters: 文本切分独立包
├── standard-tests/    # langchain-tests: 集成包的标准化测试库
├── model-profiles/    # langchain-model-profiles: 模型能力 profile 的 CLI/数据生成
├── cli/               # langchain CLI: 脚手架/模板/serve 等命令
├── Makefile            # libs/ 级别跨包 lock/check-lock
└── README.md

```

=====

## 二、包与版本关系、依赖边界 (核心理解)

### 1. langchain-core (libs/core)

langchain-core 是所有上层包与第三方集成共同依赖的“底座”。它提供统一抽象与协议，但刻意不引入任何 provider 依赖，保持依赖轻量、API 稳定、可被大量集成复用。

### 2. langchain (libs/langchain\_v1)

这里的包名是 langchain (当前活跃维护)。它在工程上属于“实现层/整合层”，主要提供：

- (1) Agent 工厂 create\_agent 与 middleware 系统 (基于 LangGraph 的图式执行)。
- (2) init\_chat\_model：对不同模型 provider 做统一初始化与动态加载。
- (3) 少量“薄封装”模块，用于将 core 的抽象组合成高层可用接口。

关键事实：langchain\_v1 的 Agent 不是“传统 while-loop”，而是编译成 LangGraph 的 StateGraph，节点/边代表执行流程与条件跳转。

### 3. langchain-classic (libs/langchain)

这是遗留兼容包 langchain-classic。它包含大量历史实现 (chains/agents/tools/...)，并在包入口中通过 `getattr` 做旧式导入的兼容提示与迁移提醒。它更像“历史兼容与迁移缓冲层”，原则上不再添加新能力。

### 4. partners (libs/partners)

partners 目录只包含 LangChain 团队直接维护的一部分第三方集成包（每个 provider 一个独立包）。生态中大量集成并不在本仓库，而是外置独立仓库（例如 google/aws 等）。

### 5. standard-tests (libs/standard-tests)

这是一个“测试库包” langchain-tests，用于给集成包提供统一测试基类与规范，使不同 provider 的实现能对齐 LangChain 的接口契约。

### 6. text-splitters (libs/text-splitters)

文本切分工具被拆成独立发布包 langchain-text-splitters，避免主包依赖膨胀，也方便单独升级发布。

### 7. model-profiles (libs/model-profiles)

模型能力 profile 维护工具包 langchain-model-profiles。它的目标是从外部数据源同步“模型能力数据”，并生成/更新各集成包里的 profile 数据文件，供 chat model 暴露 `.profile` 使用。

### 8. cli (libs/cli)

这是 langchain CLI 包，与“应用脚手架、模板、开发工具”相关。它不是 langchain-core 或 langchain\_v1 的运行时核心，但对开发体验很重要。

=====

## 三、真实对外入口 (Public API Map: 入口到实现的对应关系)

说明：下面列的是“用户实际会 import 的入口”与“仓库内真实实现文件”。这比目录树更能定位问题与扩展点。

### 1. v1 主包 langchain (libs/langchain\_v1)

#### (1) langchain.agents.create\_agent

实现：libs/langchain\_v1/langchain/agents/factory.py:create\_agent

作用：创建并编译一个可运行的 Agent 图 (CompiledStateGraph)。

#### (2) langchain.agents.AgentState

实现：libs/langchain\_v1/langchain/agents/middleware/types.py:AgentState (TypedDict)

作用：定义 Agent 的最小状态结构 (messages/jump\_to/structured\_response)。

#### (3) langchain.chat\_models.init\_chat\_model

实现：libs/langchain\_v1/langchain/chat\_models/base.py:init\_chat\_model

作用：统一初始化不同 provider 的 BaseChatModel (支持“字符串模型标识”和“运行时可配置模型”)。

#### (4) langchain.tools.tool\_node (注意)

实现：libs/langchain\_v1/langchain/tools/tool\_node.py

重要更正：该文件主要用于“兼容性 re-export”，实际 ToolNode 等核心实现来自 LangGraph (langgraph.prebuilt.tool\_node)。

### 2. core 包 langchain\_core (libs/core)

#### (1) Runnable 协议 (LCEL 组合与统一调用)

实现：libs/core/langchain\_core/runnables/base.py

(Runnable/RunnableSequence/RunnableParallel 等)

## (2) Messages

实现: libs/core/langchain\_core/messages/\*  
(AIMessage/HumanMessage/SystemMessage/ToolMessage 等)

## (3) Tools

实现: libs/core/langchain\_core/tools/\* (BaseTool + 工具 schema/转换)

## (4) Prompts

实现: libs/core/langchain\_core/prompts/\* (PromptTemplate/ChatPromptTemplate 等)

## (5) Callbacks/Tracers

实现: libs/core/langchain\_core/callbacks/\* 与 libs/core/langchain\_core/tracers/\*

### 3. classic 包 langchain\_classic (libs/langchain)

入口: langchain\_classic.\* (大量历史模块)

说明: 包入口会对部分旧式 root import 给出弃用提醒 (用于迁移, 不建议新功能基于 classic 开发)。

=====

## 四、LangChain v1 Agent: 图式执行机制 (从“循环图”升级为“可定位的图模型”)

这一节对应你原文档中的“Agent 执行流程”，但用仓库真实实现方式描述: create\_agent 编译 StateGraph。

### 4.1 Agent 的状态 (AgentState)

AgentState 是一个 TypedDict，最关键字段如下:

- (1) messages: 对话消息列表 (Required)。是整个图的主要“累积状态”。
- (2) jump\_to: 可选的流程控制字段 (tools/model/end)。它被标记为 EphemeralValue/PrivateStateAttr，表示更偏“运行期控制”而非业务输入输出。
- (3) structured\_response: 当启用结构化输出时，最终解析出的结构化结果 (输入通常省略，输出可选)。

### 4.2 create\_agent 的核心节点 (Node)

create\_agent 会按 middleware 的实现情况，动态组装下列节点 (节点名有明确约定，方便定位与调试)：

- (1) before\_agent 链: {MiddlewareName}.before\_agent (可多个, 按传入顺序串联)
- (2) before\_model 链: {MiddlewareName}.before\_model (可多个, 按传入顺序串联)
- (3) model: 固定节点名 "model" (封装模型绑定、系统提示注入、结构化输出处理、模型调用)
- (4) after\_model 链: {MiddlewareName}.after\_model (可多个, 按传入顺序串联)
- (5) tools: 固定节点名 "tools" (当存在 client-side tools 时会加入)
- (6) after\_agent 链: {MiddlewareName}.after\_agent (可多个, 结束前运行一次)

### 4.3 图的入口、循环入口、循环出口、结束节点 (Entry/Loop/Exit)

create\_agent 在构图时会计算:

- (1) entry\_node: 图开始运行的第一个节点 (优先 before\_agent, 其次 before\_model, 否则 model)
- (2) loop\_entry\_node: 每轮循环的起点 (通常是 before\_model 或 model)
- (3) loop\_exit\_node: 每轮循环的出口 (通常是 after\_model 链首或 model)
- (4) exit\_node: 图结束前最后的节点 (存在 after\_agent 则为 after\_agent 链, 否则 END)

### 4.4 工具调用与路由条件 (Conditional Edges)

这部分决定“何时去 tools、何时结束、何时回到 model”，是 Agent 行为的核心。

### (1) 从 model (或 after\_model) 到 tools / end / model\_destination

逻辑要点:

- 1) 若 messages 中最后一条 AIMessage 没有 tool\_calls，则直接走向 end\_destination (退出循环)。
- 2) 若有 tool\_calls，则会过滤掉已经有对应 ToolMessage 的调用，以及用于 structured output 的“结构化工具调用”。
- 3) 若仍存在 pending tool\_calls，则发送到 tools 节点执行 (可并发发送多条 ToolCallWithContext)。
- 4) 若 state 中已经出现 structured\_response，则走向 end\_destination (结构化输出完成即退出)。
- 5) 若 AIMessage 有 tool\_calls 但没有 pending (可能由 middleware 注入工具消息)，则回到 model\_destination 再次让模型处理。

### (2) 从 tools 到 model\_destination / end\_destination

逻辑要点:

- 1) 如果本轮执行的所有 client-side tools 都设置 return\_direct=True，则可以直接结束 (走 end\_destination)。
- 2) 如果执行到了 structured output tool (ToolStrategy 产物)，则结束 (走 end\_destination)。
- 3) 否则回到 model\_destination 开启下一轮。

## 4.5 模型绑定与结构化输出 (ResponseFormat)

create\_agent 支持 response\_format，用于结构化输出。关键点在于“策略”：

- (1) ToolStrategy：通过工具调用产出结构化结果 (会把结构化工具加入 tools 列表，并可能强制 tool\_choice)。
- (2) ProviderStrategy：使用 provider 原生 structured output 能力 (通过 to\_model\_kwargs/strict=True 等绑定方式)。
- (3) AutoStrategy：用户给 raw schema 时，先保留“自动检测意图”，运行时根据 model 能力决定用 ProviderStrategy 还是 ToolStrategy。

判定依据包含：

- (1) BaseChatModel.profile 中是否标注 structured\_output 支持。
- (2) 在 profile 不可用时，有一组 fallback 名称前缀名单 (如 gpt/o3/grok 等) 用于推断。
- (3) 部分模型对“structured output + tool calling”有组合限制 (例如代码中对 gemini 做了特殊处理)。

结构化输出的异常处理：

- (1) 若 ToolStrategy 解析失败，可根据 handle\_errors 配置决定是否向 messages 注入错误 ToolMessage 并重试。
- (2) 若同一轮产生多个结构化工具调用，会触发 MultipleStructuredOutputsError，并可走同样的重试/失败策略。

## 4.6 middleware 的两类插入点与组合顺序

middleware 的插入点分两类：

- (1) 图节点型 hook: before\_agent/before\_model/after\_model/after\_agent  
它们表现为图上的“中间节点”，返回 dict 表示 state 更新，或通过 jump\_to 控制下一跳 (需声明 can\_jump\_to 才会建立 conditional edge)。
- (2) 调用包裹型 hook: wrap\_model\_call / wrap\_tool\_call  
它们不表现为独立图节点，而是“包裹住 model 调用/工具执行”的 handler 链。组合规则为：  
传入 middleware 列表的第一个是最外层 (outermost)，会包裹住后续所有 handler。

同步/异步约束 (非常关键)：

如果 middleware 只实现了 awrap\_model\_call，而你用同步方式运行 agent (invoke/stream)，会触发 NotImplementedError；反之亦然。原因是该系统不会自动把 async hook 适配为 sync hook (或相反)，以避免隐藏的线程/事件循环语义问题。

=====

# 五、Middleware 目录深解（按能力分组，而不是按文件名堆叠）

---

目录位置：libs/langchain\_v1/langchain/agents/middleware/

## 5.1 middleware 的形态

(1) 类形态：继承 AgentMiddleware，实现任意 hook 方法。

(2) 函数形态：通过 decorators 动态生成 middleware 实例：

before\_agent / before\_model / after\_model / after\_agent / wrap\_model\_call / wrap\_tool\_call / dynamic\_prompt

这些 decorators 定义在 types.py 中，能够把普通函数包装为 AgentMiddleware 的实例，并可设置 state\_schema/tools/can\_jump\_to/name 等属性。

## 5.2 安全与合规类（输入输出治理）

(1) \_redaction.py：脱敏与内容清洗相关逻辑（通常插入 model 调用前后或 tool 调用前后）。

(2) pii.py：个人信息保护策略（常见做法包括检测、掩码、阻断或提示）。

## 5.3 可靠性与稳定性类（重试、回退、限流、上限）

(1) \_retry.py：通用重试策略基元（供其他模块复用）。

(2) model\_retry.py：模型调用重试（常在 wrap\_model\_call 层实现）。

(3) tool\_retry.py：工具执行重试（常在 wrap\_tool\_call 层实现）。

(4) model\_fallback.py：模型回退（失败时切换备用模型）。

(5) model\_call\_limit.py：限制模型调用次数（防止无限循环或成本失控）。

(6) tool\_call\_limit.py：限制工具调用次数（防止工具循环或滥用）。

## 5.4 人在回路与交互治理类（可控性、可审查）

(1) human\_in\_the\_loop.py：在人需要确认/干预的节点中断与恢复（通常依赖 jump\_to 或图中断机制）。

(2) todo.py：面向 Agent 的待办/计划辅助逻辑（用于任务分解或执行轨迹）。

(3) summarization.py：对 messages 做摘要压缩（控制上下文长度与成本）。

(4) context\_editing.py：对上下文做“编辑/重写/裁剪”的通用能力。

(5) file\_search.py：文件检索类能力（通常与工具/上下文增强一起用）。

## 5.5 工具策略与工具治理类（选择、模拟、shell）

(1) tool\_selection.py：工具选择策略（可能在 before\_model 或 wrap\_model\_call 中影响 tools/tool\_choice）。

(2) tool\_emulator.py：工具模拟（在无真实工具环境时模拟 tool 行为或回放）。

(3) shell\_tool.py：shell 工具相关能力（通常包含安全边界与执行策略）。

(4) \_execution.py：shell 工具执行策略与安全限制（Host/Docker/CodexSandbox 等），说明了“为何要把执行策略抽象出来”。

提示：上述文件的具体 hook 点需要结合 factory.py 的构图方式理解。最关键的判断是：该 middleware 是“图节点型”（before/after）还是“包裹型”（wrap\_model\_call/wrap\_tool\_call），以及它是否声明 can\_jump\_to 从而允许流程跳转。

=====

# 六、Core 层关键概念补全（与 v1/classic/partners 的关系）

---

## 6.1 Runnable 协议与 LCEL (runnables/)

Runnable 是 LangChain 的统一“可执行单元”协议。关键方法包括：

invoke/ainvoke: 单输入 -> 单输出

batch/abatch: 多输入 -> 多输出 (默认可并行加速)

stream/astream: 流式输出

这些统一语义使得 prompt、model、parser、tool、router 都能组合成链式结构。

组合语义 (工程价值) :

(1) 使用 | 构建 RunnableSequence (顺序执行, 前一步输出给后一步) 。

(2) 在 sequence 中放入 dict 形成 RunnableParallel (并行执行, 输出字典) 。

(3) stream 能否贯穿整条链取决于组件是否实现 transform/atransform, 某些组件会阻塞流式直到完成。

## 6.2 Message 系统 (messages/)

Message 是对话态应用的“统一数据结构”。其中 ToolMessage/FunctionMessage 承担“模型调用工具后的结果回填”语义, 是 Agent 循环与工具执行的粘合层。

## 6.3 Tools (tools/)

BaseTool 抽象了工具的 schema、参数、执行与结果返回方式。v1 的 Agent 通过模型的 bind\_tools 将工具 schema 暴露给模型, 并通过 tools 节点执行 client-side 工具, 把结果封装为 ToolMessage 回填到 messages。

## 6.4 Prompts (prompts/ + prompt\_values.py)

PromptTemplate/ChatPromptTemplate 等提供“可组合的提示构建”。它们往往作为 Runnable 链的一段出现 (prompt | model | parser) 。

## 6.5 Callbacks 与 Tracers (callbacks/ + tracers/)

这两类模块提供运行期观测、调试与记录能力。它们与 Runnable 的 config 参数协同工作, 用于把运行轨迹输出到 console 或外部观测系统 (例如 LangSmith) 。

## 6.6 序列化与兼容 (load/)

Serializable 等能力是跨进程、跨环境传递 LangChain 组件信息的基础。由于历史包 classic 负担较大, core 的 load/serializable 更偏“通用底座”, classic 里仍会存在更多历史兼容路径。

## 6.7 关于 langchain\_core.agents.py (常见误解纠正)

langchain\_core.agents.py 主要提供旧式 agent schema (AgentAction/AgentStep/AgentFinish 等), 并在 docstring 中明确这是“为了向后兼容”。v1 的 create\_agent 体系并不以这些 schema 为核心, 而是基于 messages + LangGraph StateGraph。

---

# 七、集成生态与标准测试 (为什么 partners 只是“部分”)

---

## 7.1 partners 的定位

partners 是“官方维护的一小部分集成包集合”, 每个目录都是独立发布包 (例如 langchain-openai、langchain-anthropic) 。本仓库并不包含全部 LangChain 生态集成。

## 7.2 外置集成仓库

大量集成被拆到独立仓库（例如 google、aws 等），常见原因包括：

- (1) 依赖复杂度高，不适合让主包背负大量可选依赖。
- (2) 版本与发布节奏不同，需要独立维护。
- (3) 更便于与第三方共同协作与测试。

## 7.3 standard-tests (langchain-tests)

这是一个“测试库”，为集成包提供标准化测试基类（例如

ChatModelUnitTests/ChatModelIntegrationTests）。它帮助保证不同 provider 的实现满足相同接口契约，减少“实现可用但边界行为不一致”的问题。

=====

# 八、开发与构建：uv + Makefile (monorepo 研发视角)

## 8.1 独立发布包的工程事实

libs/ 下每个包目录都有：

pyproject.toml：该包依赖与构建配置

uv.lock：该包锁定依赖版本

因此它是一个“多包独立版本”的 monorepo，而不是单包。

## 8.2 libs/Makefile 的跨包能力

libs/Makefile 提供跨包的 lock/check-lock，用于一次性更新或校验 core/text-splitters/langchain/langchain\_v1/model-profiles 等包的 lockfile 一致性。

提示：这类 Makefile 主要用于维护者的依赖一致性与 CI 校验，不属于运行时核心逻辑。

=====

# 九、文件统计（工程量级参考）

以下数字为粗略量级，便于把握复杂度：

1. langchain-core：约 200+ 文件（底座抽象与协议）
2. langchain (v1)：约 30+ 文件（精简整合层，但依赖 LangGraph 承担图执行）
3. langchain-classic：约 1500+ 文件（历史实现与兼容负担）
4. partners：每个集成 15-100 文件不等（依 provider 复杂度变化）

总计：约 3000+ 核心代码文件（不含测试数据与部分模板资源）