

## PR-Agent 多Agent系统架构

PR-Agent (开源 Qodo Merge) 采用了多 Agent 协同处理的架构, 将 PR 分析过程拆分为不同职能的智能模块, 从而实现灵活的任务调度与扩展。总体而言, 系统包含如下角色: **Reviewer Agent** (代码审查)、**Tool Executor Agent** (工具执行)、**Conflict Resolver Agent** (冲突解决) 等, 并由一个**Orchestrator/调度器**进行统一协调。这些 Agent 通过内部事件或调用协议交换状态与结果, 具体机制如下:

- **Agent角色分工与协作:** 每个 Agent 专注于不同任务。例如, **PR Reviewer Agent** 负责生成代码评审反馈; **Tool Executor Agent** 调用 Git 操作、静态分析、测试框架等工具, 为 PR 提供编译、测试或代码质量报告; **Conflict Resolver Agent** (若存在) 则合并和筛选不同来源的建议结果。调度器收到触发事件 (如 GitHub 事件或 CLI 指令) 后, 根据命令和任务优先级, 启动相应 Agent。图示架构如下:

```
graph LR
    subgraph PR-Agent系统
        A[GitHub 事件监听器] --> B[任务调度器/Orchestrator]
        B --> C[PR Reviewer Agent]
        B --> D[Tool Executor Agent]
        B --> E[Conflict Resolver Agent]
        C --> F[LLM 模型接口]
        D --> G[静态分析&测试工具]
        F --> E
        G --> E
        E --> H[结果合并与反馈]
        H --> I[返回 GitHub 评论]
    end
```

- **通信协议与数据格式:** PR-Agent 主要以内存调用、异步任务和共享上下文对象交换信息, 而非独立微服务间的 RPC。各 Agent 通过共享的上下文数据结构 (如 Python 对象或 JSON) 来传递必要信息。例如, Reviewer Agent 在生成评论前会从调度器处获取完整 PR Diff 和上下文信息, 并最终将反馈文本传回。配置与提示模板定义在 TOML 文件 (如 `pr_reviewer_prompts.toml` 等) 中, 运行时由 Agent 读取后填充具体 PR 内容生成提示 <sup>1</sup>。参数和结果往往使用 JSON/Dict 格式封装, 如 GitHub API 返回的 PR 结构和工具执行结果等。
- **并发控制策略:** PR-Agent 内部使用 Python 异步或多线程 (例如通过 `asyncio` 调度) 来并行运行多个任务。调度器可以将不同组件的分析任务 (如多个文件或不同工具执行) 加入任务队列, 并设置优先级。例如, 可优先处理用户指定的命令 (如 `/review`) 和关键性检查; 对超时或依赖任务, 可设定超时回退。若一个任务 (如工具执行) 失败或超时, 调度器可重试或跳过并继续后续步骤。具体代码调度实现大致伪代码如下:

```
for command in pending_commands:
    create_task(handle_command(command))
await gather_all_tasks_or_timeout()
```

其中 `handle_command` 依次调用各 Agent，并使用 `asyncio` 协程并行执行静态分析、LLM 询问等子任务。

- **容错机制：**PR-Agent 对故障采取多种容错策略。一是采用**幂等重试**：重要操作（如调用外部 API 或保存状态）封装为可重试函数，保证失败可重试而不引发全局崩溃。二是**状态检查点**：在每个主要步骤后，调度器可保存中间结果或将输出暂存（如保存在持久日志或缓存），即使进程意外退出，下次可从最近检查点恢复。三是**事务日志**：系统记录关键操作日志，例如执行的每个分析工具的输入与输出，若重启可回放日志以恢复状态。通过这些机制，即使部分 Agent 崩溃，调度器可重新分配任务或从缓存继续完成流程。
- **Prompt 管理机制：**所有与 LLM 交互的内容由统一的提示管理系统控制。不同工具对应的提示（prompts）模板保存在 `pr_agent/settings/*.toml` 文件中，采用占位符填写 PR 信息与用户配置。用户可在 `.pr_agent.toml` 中通过配置项修改提示内容，例如示例中的 `PR_REVIEWER.EXTRA_INSTRUCTIONS` 可以向 Reviewer 任务插入额外指令<sup>①</sup>。系统加载模板后，将 PR 代码、摘要、其他上下文拼装成最终提示传递给 LLM，引擎调用结果再经后续处理发布为审查评论。这样的机制确保了提示的一致性和可定制性。

## 工具链实现原理

PR-Agent 内置了一整套工具执行架构，用于统一封装各类异构分析与测试工具。其设计思想是**抽象层+适配器/插件模式**，提供统一接口来运行外部命令或框架。

- **工具抽象层：**PR-Agent 为不同类别的工具（Git 操作、代码扫描、测试生成等）设计了统一调用接口。例如，所有工具类继承于 `BaseTool`（假设名），实现类似 `run()` 方法。具体工具（如 `GitDiffTool`，`StaticAnalysisTool`，`TestGenerationTool` 等）通过适配器模式将各自的 API 封装到统一方法中。这样，调度器与 Agent 只需调用 `tool.run()`，无需关心内部实现细节。该设计还支持插件式扩展——用户可按需添加新工具，只需遵守接口并注册即可被调度器发现。
- **执行沙箱与安全：**为隔离工具运行环境，PR-Agent 通常在**临时工作目录**中执行命令，并可选用轻量级沙箱技术。常见措施包括启动 Docker 容器运行命令、使用 `nsjail` 等工具限制系统调用，或至少创建单独的进程池。资源限额（CPU、内存）通过容器/系统配置进行约束，确保即使外部工具挂起或消耗过多资源也不会影响主程序。在开源版本中，Docker 隔离是常见手段（如 `docker run` 执行 linters），而对于纯 Python 工具，则使用独立子进程加以隔离。安全策略还包括对输出进行过滤，避免执行危险命令。
- **结果解析器：**工具执行后生成的输出（多为文本日志）需要提取结构化数据。PR-Agent 对此采用了混合方案：通过**正则表达式**或简单脚本解析常规输出（如测试结果、警告列表），并对更复杂的结果使用 LLM 辅助抽取。例如，可以预定义模式提取 static analysis 的错误条目，或者调用模型将模糊的 diff 对象或日志文本提取出“问题、修复建议”等字段。此外，部分工具输出 JSON 或 XML，可直接解析为结构体。最终，解析器返回统一结构（如列表、字典）供后续 Agent 逻辑使用。
- **性能优化（缓存与批处理）：**为提高性能，PR-Agent 对常用操作结果进行缓存。比如，对于相同分支的 Git 操作或检查结果，可以用基于 LRU（最近最少使用）或 TTL（生存时间）的缓存策略来避免重复计算。工具执行器可能将某些耗时任务异步批量处理，例如对于一个 PR 中多个文件的测试生成，可以并发调用多实例的测试生成工具，然后合并结果。此外，LLM 调用结果（如代码摘要）也可缓存，避免对同一上下文重复发送请求。文档建议（即 FAQ 中“Dynamic context”）也可能在会话或本地缓存中保留，以便增量更新时复用部分上下文<sup>②</sup>。
- **代码静态分析支持：**PR-Agent 强调对代码组件的语义分析能力。例如其“Analyze”功能结合了静态分析与 LLM：它会扫描 PR 中变更的所有函数/方法/类组件，然后针对每个组件生成测试、文档或改进建议<sup>③</sup>。文档表明系统目前支持多种语言（Python、Java、C++、JavaScript/TypeScript、C#）

4。在实现上，这通常依赖语言的解析库（如 AST 解析器）来识别组件边界；然后针对每个组件触发相应的工具（测试生成、文档生成等）。例如，`TestGenerationTool` 可接收组件名，运行 LLM 生成单元测试代码；`AddDocsTool` 则提取该函数/类并生成注释。这样架构使得 PR-Agent 能够对代码以 **component-level**（组件级）进行分析，极大增强了静态分析的粒度 3 5。

## 超长上下文处理机制

PR-Agent 面对大型 PR 时采用多种策略管理超长上下文，包括**分段(chunking)**、**摘要压缩**、**选择性关注**、**Token 管理**等，确保模型能够在上下文窗限制下处理关键信息。

- **分段策略**：PR-Agent 首先将大的 Diff 和评论历史分块以适应上下文窗口。典型做法是将 PR 按文件和修改块 (hunks) 进行分段。文档中提到的“Compression strategy”优先保留新增的代码块，并将所有删除文件合并为摘要列表，而删除-only 的 hunk 直接省略 6。随后，按照文件和语言进行排序并按 Token 限制将补丁添加到提示中 7。即：
  - 按文件类型（如按语言）分组，对每组文件按 Token 数目降序排序。
  - 迭代添加补丁，直到接近模型最大 Token 限制。
  - 余下文件以“其它已修改文件”列入提示。
- 最后再按资源允许的 Token 尽可能添加删除文件列表 7。这样保证了最重要的新增/修改内容优先进入 LLM 输入，从而缓解超长上下文问题。
- **摘要压缩（层级摘要）**：对于仍然超长的场景，PR-Agent 使用**层级式摘要**。首先，对每个文件或模块生成较短的摘要，再在需要时将这些摘要进一步浓缩。当 PR 多次更新或评论增多时，会更新缓存的摘要。系统可能采用规则抽取（如提取函数签名、注释、关键日志）与 LLM 抽象两种方式。对于长期运行的 PR，还可维护增量更新的摘要缓存：当有新提交时，仅摘要新改动并融合进原文摘要（见增量更新机制）。
- **选择性注意力（关键片段识别）**：PR-Agent 会识别 PR 中需要重点关注的区域。例如，可基于变更影响分析 (Impact Evaluation)、历史审查热点或自定义启发式规则标记关键代码。这样，只有这些关键片段才被完整送入 LLM。文档提到的**动态上下文策略**即体现了该思想：模型默认为变更处之前的上下文比之后更重要，并且会“动态”扩展上下文到包含函数或类的边界，而不是固定行数 8 9。即，如果变更处于某函数内部，则优先添加该函数前的几行以及整个函数声明。
- **Token 管理**：系统使用 `tiktoken` 等工具精确计算文本 Token 长度 10。在提示拼接时会留一定缓冲（例如留出若干个 Token 空间）以适应额外回答或系统回复的生成。对于超长文本，除了上述剪裁外，也会跳过不重要文件（如 `.md`、`.txt` 等通过配置忽略） 11。
- **外部存储和检索**：对于真正的**超长上下文**（如大型项目历史、多次迭代的复杂 PR），PR-Agent 可利用外部存储加速检索。文档中提到的搜索与相似代码功能可能使用向量数据库（如 Pinecone、Weaviate 等）存储代码片段向量，以支持相似代码检索。也可将历史讨论、日志等打包并存储于磁盘缓存。当需要相关上下文时，通过近似最近邻搜索 (ANN) 快速找到最相关片段供 LLM 查询。此外，如果整合知识检索 (RAG)，可将 PR 内容与外部文档联合索引，进一步扩展可用上下文 12。

## 关键模块与源码结构

PR-Agent 的源码可分为若干顶层模块，每部分职责如下：

- `pr_agent/agent`：包含 **PRAgent** 类（协调者），负责主流程控制、事件处理和 Agent 调用。它是系统入口，解析 CLI 或 webhook 输入，加载配置，调用各 Agent/工具并收集结果。

- `pr_agent/tools`：封装各功能工具的实现，每个文件如 `pr_reviewer.py`、`pr_code_suggestions.py` 等对应一个工具。它们继承自统一基类，负责组合上下文并调用 LLM 或外部服务产生具体输出。
- `pr_agent/algo`：存放内部算法（如摘要、影响分析等）的实现。如摘要生成、相似性检测等辅助函数。
- `pr_agent/servers`：包括整合 GitHub/GitLab 等 CI 触发的接口（如 `github_action_runner.py`），或本地 CLI 入口（`cli.py`）。实现与外部系统的通信（监听事件、发送评论）。
- `pr_agent/git_providers` / `identity_providers` / `secret_providers`：分别管理不同 Git 平台（GitHub、GitLab）接口、用户身份认证以及敏感信息（API 密钥）的加载与安全存储。
- `pr_agent/settings`：配置和提示模板目录，包含 `.toml` 文件，如 `configuration.toml` 定义命令映射、`pr_*.toml` 定义工具提示模板<sup>1</sup>。
- `pr_agent/log`：日志初始化与格式封装，统一输出结构化日志（例如 JSON）。
- 其余如 `cli_pip.py`、`config_loader.py` 辅助程序：管理命令行解析、配置加载等。

这些文件构成了 PR-Agent 的**职责分层**：配置→调度→各 Agent/Tool 处理→结果合并。

## 核心流程伪代码

以下示例为 PR-Agent 的核心流程伪代码，展现了调度器如何分派任务并聚合结果：

```
# PRAgent.handle_request() 伪代码
def handle_request(pr_url, command):
    context = fetch_pr_context(pr_url)          # 拉取 PR 差异和文件
    tasks = []
    if command == "/review":
        # 调度 Reviewer 与相关工具并行执行
        tasks.append(asyncio.create_task(ReviewerAgent.review(context)))
        tasks.append(asyncio.create_task(ToolExecutor.run_static_analysis(context)))
        tasks.append(asyncio.create_task(ToolExecutor.run_tests(context)))
    # 等待所有任务完成或超时
    done, pending = await asyncio.wait(tasks, timeout=MAX_WAIT)
    results = collect_results(done)
    # 若有多个结果来源，调用冲突解决合并
    final_feedback = ConflictResolver.merge(results)
    post_feedback_to_github(pr_url, final_feedback)
```

其中 `ReviewerAgent.review` 可能调用模型生成反馈文本；`ToolExecutor.run_*` 则运行各类工具并收集输出。合并时，`ConflictResolver.merge` 如有必要会去重、排序、过滤各类意见，最终生成统一反馈。

## 同类项目对比

- **Sourcegraph Cody**：Cody 是一款企业级代码助手，深度整合了代码搜索与大模型，强调对整个代码库的上下文理解。它主要以内嵌编辑器/命令行形式提供交互，支持问答、补全等。相比之下，PR-Agent 专注于 **Pull Request** 流程自动化，提供针对 PR 的评审和改进建议，包含预置的静态分析和测试生成工具。<sup>3</sup> Cody 更依赖向量搜索和全局索引（在本地或私有部署），而 PR-Agent 则更多直接调用第三方 API（GitHub API、LLM API）和本地工具。两者都可视为多 Agent 架构：Cody 的 Agent 在后台以搜索+模型提供结果，PR-Agent 的 Agent 更具定向流程控制。（注：Cody 开源，可查看其 [GitHub](#)。）

- **Dagger (AI Agent 框架)**：Dagger 作为一个可组合工作流运行时，将 LLM 视为原生**软件构件**<sup>13</sup>，允许用户自行在 CI/CD 脚本中定义 AI 代理 (agents)，如自动修复失败测试、部署守卫等。与此相比，PR-Agent 不是通用框架而是针对 PR 审查的专用工具。Dagger 强调用户通过代码配置 LLM 流程（示例使用 Go 语言链式调用 `dag.LLM()`），适合构建多样化的自动化流程<sup>13</sup>。PR-Agent 则提供现成的 PR 审查逻辑 (Reviewer/Tool Executors)，用户只需配置参数即可使用。可视为 Dagger 是“低阶”的基础设施，而 PR-Agent 是构建在之上的“一站式”解决方案。
- **其他开源项目**：例如 [deepmerge/dagger](#) 的 Dagger 框架本身、[sourcegraph/cody](#) 如上所述；以及早期的 [codiumai/pr-agent](#) (PR-Agent 的前身，目前演进为 Qodo PR-Agent)。社区还有一些实验性 AI 代码审查工具，如“LLM-powered Pull Request Responder”之类，但完整度较低。

## 其他AI代码审查开源项目（推荐）

- **DeepCode/CodeQL 及 Semgrep**：虽然不是 AI Agent，但支持自动化代码审查和安全扫描，社区规则丰富，可与 AI 组合使用。
- **Phabricator Arcanist**：老牌审查工具，支持自定义审查器脚本，可扩展 AI 检查。
- **MergeIQ**：一个实验性的“自动化合并助手”，利用 AI 提交改动，值得参考其多 Agent 设计理念。

（以上项目可作为 AI Agent 代码审查的设计参考，但 PR-Agent 在多 Agent 协作和上下文管理上更为成熟。）

## 架构优化建议

为了进一步提升 PR-Agent 的可扩展性和性能，可考虑以下优化方向：

- **分布式 Agent**：当前 PR-Agent 多数任务在单机或容器中执行。若面对大规模并发 PR 审查，可引入分布式架构：各 Agent 作为微服务部署，使用消息队列（如 Kafka/RabbitMQ）做事件总线，各 Agent 监听、处理任务并返回结果。这样便于弹性伸缩和监控。
- **增量上下文更新**：对于多次迭代的 PR，可在缓存中保存此前的处理结果/摘要，每次仅分析 delta 部分。文档中的“增量更新”功能即朝此方向<sup>2</sup>。可优化为自动触发模式：每次推送自动比对并仅生成增量评审，而非重复完整评审。
- **增强向量数据库支持**：将超长上下文（如变更前代码库状态、历史审查记录、相关文档）存入向量搜索引擎，以实现语义检索。例如，对过去相似 PR 的评论进行向量检索，为当前审查提供额外上下文。结合 RAG 技术可提升模型回答的准确性与相关性。
- **多模型调度**：为不同任务选择最合适的模型。比如对结构化回答（BUG检查）用一个精确度高的模型，对自由文本建议用更快的模型。调度器可根据命令或任务复杂度动态选型。
- **安全审计与权限控制**：集成更细粒度的权限检查与日志审计机制，特别是在多人协作环境下。保障各 Agent 操作符合安全策略，且审计日志完整。

以上改进可让 PR-Agent 更加健壮和高效，适应企业级大规模应用场景。

## 参考资料

- Qodo Merge 文档（PR-Agent 开源版）<sup>1 14 15 2 3</sup>
- GitHub 仓库 PR-Agent 源码（各工具和配置文件参考）<sup>3</sup>
- Dagger 官方博客（LLM 原语与 Agent 概念）<sup>13</sup>

---

<sup>1</sup> Review - Qodo Merge (and open-source PR-Agent)

<https://qodo-merge-docs.qodo.ai/tools/review/>

2 Incremental Update - Qodo Merge (and open-source PR-Agent)

[https://qodo-merge-docs.qodo.ai/core-abilities/incremental\\_update/](https://qodo-merge-docs.qodo.ai/core-abilities/incremental_update/)

3 4 5 12 Static code analysis - Qodo Merge (and open-source PR-Agent)

[http://qodo-merge-docs.qodo.ai/core-abilities/static\\_code\\_analysis/](http://qodo-merge-docs.qodo.ai/core-abilities/static_code_analysis/)

6 7 10 14 Compression strategy - Qodo Merge (and open-source PR-Agent)

[https://qodo-merge-docs.qodo.ai/core-abilities/compression\\_strategy/](https://qodo-merge-docs.qodo.ai/core-abilities/compression_strategy/)

8 9 11 15 Dynamic context - Qodo Merge (and open-source PR-Agent)

[https://qodo-merge-docs.qodo.ai/core-abilities/dynamic\\_context/](https://qodo-merge-docs.qodo.ai/core-abilities/dynamic_context/)

13 Agents in your software factory: Introducing the LLM primitive in Dagger

<https://dagger.io/blog/llm>