

deepseek-dify

dify-ollama AI



- <https://iovhm.com/book/books/cee63/page/9872e>
- <https://iovhm.com/book/books/k8s/page/harbordockerdocker>
- <https://github.com/docker/compose/releases>
- <https://iovhm.com/book/books/bbcbf/page/deepseek>
- <https://github.com/langgenius/dify>

ollama

```
version: "3"
services:
  ollama:
    image: harbor.iovhm.com/hub/ollama/ollama:0.5.12
    container_name: ollama
    restart: always
    privileged: true
    ports:
      - "11434:11434"
    volumes:
      - ./ollama:/root/.ollama
# deploy:
#   resources:
#     reservations:
#     devices:
#       - driver: nvidia
#         capabilities: [gpu]
#         count: all
networks:
  - vpclub-bridge
```

2

```
# ollama pull deepseek-r1
# ollama pull bge-m3
```

```
diff docker docker
```

```
docker-compose.yml docker dify
```

```
10
```

```
docker-compose up -d docker-compose --profile=xxx up -d
```

```
docker-compose.yml
```

- api
- worker
- web
- db
- redis
- sandbox
- plugin_daemon
- ssrf_proxy
- nginx
- weaviate

```
env
```

```
.env.example env .env .env8080 .env
```

<https://docs.dify.ai/zh-hans/getting-started/install-self-hosted/environments>

```

CONSOLE_API_URL=
CONSOLE_WEB_URL=
SERVICE_API_URL=
APP_API_URL=
APP_WEB_URL=
FILES_URL=

#
EXPOSE_NGINX_PORT=80
EXPOSE_NGINX_SSL_PORT=443

```

```
# [REDACTED], [REDACTED] false, [REDACTED] https://updates.dify.ai [REDACTED]
# [REDACTED] CloudFlare Worker [REDACTED]
# [REDACTED], [REDACTED]
CHECK_UPDATE_URL=

# [REDACTED]
VECTOR_STORE=weaviate

# Weaviate [REDACTED], [REDACTED] http://weaviate:8080
WEAVIATE_ENDPOINT=http://weaviate:8080

# [REDACTED] Weaviate [REDACTED] api-key [REDACTED]
WEAVIATE_API_KEY=WVF5YThaHlkYwhGUSmCRgsX3tD5ngdN8pkih
```

ollama

dify

设置

- 工作空间
- 模型供应商**
- 成员
- 数据来源
- API 扩展
- 通用
- 语言

模型供应商 Q 搜索 ESC

系统模型设置

模型列表

- Ollama** (LLM, TEXT EMBEDDING) 显示模型 > 添加模型
- 深度求索** (LLM) API-KEY 设置

安装模型供应商 发现更多就在 Dify 市场 >

- OpenAI** (langgenius · 15,369) 如果没有增加模型供应商，在这里会出来ollama 我已经安装了，所以就不出来了
- Anthropic** (langgenius · 4,260)
- AWS** (langgenius · 933)
- Azure OpenAI** (langgenius · 2,207)
- Azure AI Studio** (langgenius · 714)
- Cohere** (langgenius · 1,163)

添加 Ollama

模型类型 *

LLM
 Text Embedding

模型名称 *

deepseek-r1:7b

基础 URL *

http://ollama:11434

模型类型 *

对话 ×

模型上下文长度 *

4096

最大 token 上限 *

4096

是否支持 Vision

是
 否

最多支持 10 个模型

[如何集成 Ollama](#)



word markdown -> wps ai



选择数据源



导入已有文本



同步自 Notion 内容



同步自 Web 站点

上传文本文件



拖拽文件至此，或者 [选择文件](#)

已支持 TXT、MARKDOWN、MDX、PDF、HTML、XLSX、XLS、DOCX、CSV、MD、HTM，每个文件不超过 15MB。



威海智慧谷知识库.docx
DOCX · 0.02MB



下一步 →

[创建一个空知识库](#)



分段设置



通用

通用文本分块模式，检索和召回的块是相同的



父子分段

选择父子分段

使用父子模式时，子块用于检索，父块用作上下文

父块用作上下文



段落

此模式根据分隔符和最大块长度将文本拆分为段落，使用拆分文本作为检索的父块

分段标识符 [ⓘ]

\n\n

分段最大长度

500

tokens



全文

整个文档用作父块并直接检索。请注意，出于性能原因，超过10000个标记的文本将被自动截断。

子块用于检索

分段标识符 [ⓘ]

\n

分段最大长度

200

tokens



文本预处理规则

替换掉连续的空格、换行符和制表符

删除所有 URL 和电子邮件地址

预览块

重置



索引方式



高质量 推荐

调用嵌入模型处理文档以实现更精确的检索，可以帮助LLM生成高质量的答案。



经济

每个数据块使用10个关键词进行检索，不会消耗任何tokens，但会以降低检索准确性为代价。

⚠ 使用高质量模式进行嵌入后，无法切换回经济模式。

Embedding 模型

bge-m3

检索设置

[了解更多](#)关于检索方法，您可以随时在知识库设置中更改此设置。



向量检索

通过生成查询嵌入并查询与其向量表示最相似的文本分段



Rerank 模型

Top K

3

Score 阈值

0.5



全文检索

索引文档中的所有词汇，从而允许用户查询任意词汇，并返回包含这些词汇的文本片段



混合检索 推荐

同时执行全文检索和向量检索，并应用重排序步骤，从两类查询结果中选择匹配用户问题的最佳结果，用户可以选择设置权重或配置重新排序模型。



创建空白应用

选择应用类型

新手适用



聊天助手

简单配置即可构建基于 LLM 的对话机器人



Agent

具备推理与自主工具调用的智能助手



文本生成应用

用于文本生成任务的 AI 助手

进阶用户适用



BETA

Chatflow

支持记忆的复杂多轮对话 workflow



BETA

工作流

面向单轮自动化任务的编排 workflow

应用名称 & 图标

威海小智



描述 (可选)

输入应用的描述

没有想法? 试试我们的模板 →

取消

创建





编排

提示词 生成

在这里写你的提示词，输入 '{' 插入变量、输入 '/' 插入提示内容块

提示词，好的提示词能得到意想不到的效果

0

变量 添加

变量能使用户输入表单引入提示词或开场白，你可以试试在提示词中输入 {{input}}

知识库 召回设置 添加

您可以导入知识库作为上下文

模型选择和参数调整

deepseek-chat CHAT 发布

发布

调试与预览

保存并发布应用

调试窗口

和机器人聊天 发送

功能已开启

更多设置

管理



董 董列涛's W... 探索 工作室 / 威海小智 知识库 工具 插件 董 董列涛

编排

提示词 + 生成

威海智慧谷智慧园区智能问答助手小智提示词

定位
我叫小智，我是威海智慧谷智慧园区智能问答助手，是一个专为园区管理、企业员工及访客设计的智能交互平台。旨在通过自然语言处理技术，提供即时、准确的园区相关信息和服务支持。

能力
1. **信息查询**：能够快速响应关于园区设施、服务、活动等信息的查询。
2. **园区公告**：提供园区公告信息，包括新建建筑、园区活动等。
3. **园区服务**：提供园区服务信息，包括园区班车、园区保洁等。
4. **园区设施**：提供园区设施信息，包括园区会议室、园区健身房等。
5. **园区安全**：提供园区安全信息，包括园区门禁、园区监控等。
6. **园区环境**：提供园区环境信息，包括园区绿化、园区空气质量等。
7. **园区交通**：提供园区交通信息，包括园区公交线路、园区停车场等。
8. **园区生活**：提供园区生活信息，包括园区餐饮、园区购物等。
9. **园区文化**：提供园区文化信息，包括园区文化活动、园区展览等。
10. **园区其他**：提供园区其他信息，包括园区新闻、园区公告等。
794

变量 + 添加
变量能使用户输入表单引入提示词或开场白，你可以试试在提示词中输入 {{input}}

知识库 召回设置 + 添加

iovhm-com 高质量·向量检索

调试与预览

你服

和机器人聊

功能

增强 web app 用户体验

- 对话开场白
你好，我是小智，是威海智慧谷的智能客服，我将竭诚为您服务，请问有什么可以帮到您的？
- 下一步问题建议
设置下一步问题建议可以让用户更好的对话。
- 引用和归属
显示源文档和生成内容的归属部分。
- 内容审查
您可以调用审查 API 或者维护敏感词库来使模型更安全地输出。
- 标注回复
启用后，将标注用户的回复，以便在用户重复提问时快速响应。



 代码改写 对代码进行修改，来实现纠错、注释、调优等。	 代码解释 对代码进行解释，来帮助理解代码内容。
 代码生成 让模型生成一段完成特定功能的代码。	 内容分类 对文本内容进行分析，并对齐进行自动归类
 结构化输出 将内容转化为 Json，来方便后续程序处理	 角色扮演 (自定义人设) 自定义人设，来与用户进行角色扮演。
 角色扮演 (情景续写) 提供一个场景，让模型模拟该场景下的任务对话	 散文写作 让模型根据提示词创作散文
 诗歌创作 让模型根据提示词，创作诗歌	 文案大纲生成 根据用户提供的主题，来生成文案大纲
 宣传标语生成 让模型生成贴合商品信息的宣传标语。	 模型提示词生成 根据用户需求，帮助生成高质量提示词
 中英翻译专家 中英文互译，对用户输入内容进行翻译	

[[[]]]curl [[[]]]postman[[[]]]deepseek api[[[]]]deepseek[[api key]]

```
curl --location 'https://api.deepseek.com/chat/completions' \  
--header 'Content-Type: application/json' \  
--header 'Authorization: Bearer sk-xxxxxxxxxxxxxxxxxxx' \  
--data '{  
  "model": "deepseek-chat",  
  "messages": [  
    {"role": "system", "content": "[[[]]]L. [[ Markdo  
    {"role": "user", "content": "[[[]]] \\' [[[]]] \\' [[[]]]"}  
  ],  
  "stream": false  
}'
```

[[[]]]deepseek, [[[]]]

[[[]]]deepseek <https://www.deepseek.com/>

DeepSeek-R1 已发布并开源，性能对标 OpenAI o1 正式版，在网页端、APP 和 API 全面上线，点击查看详情。

deepseek

探索未至之境

API key

deepseek 开放平台

- 用量信息
- API keys
- 充值
- 账单

API keys

列表内是你的全部 API key，API key 仅在创建时可见可复制，请妥善保存。不要与他人共享你的 API key，或将其暴露在浏览器或其他客户端代码中。为了保护你的帐户安全，我们可能会自动禁用我们发现已公开泄露的 API key。我们未对 2024 年 4 月 25 日前创建的 API key 的使用情况进行追踪。

名称	Key	创建日期	最新使用日期
iovhm-com	sk-cf29e*****db46	2025-03-06	2025-03-06

创建 API key

10 5-10 key

deepseek 开放平台

- 用量信息
- API keys
- 充值
- 账单

充值

在线充值 对公汇款

支付金额

¥10 ¥20 ¥50 ¥100 ¥300 ¥500 自定义 查看价格

【错峰优惠活动】北京时间每日 00:30-08:30 为错峰时段，API 调用价格大幅下调：DeepSeek-V3 降至原价的 50%，DeepSeek-R1 降至 25%，在该时段调用享受更经济更流畅的服务体验。【查看价格详情】

支付方式

- 支付宝 ALIPAY
- 微信支付

去支付

deepseek

#

编排

提示词 + 生成

- 3. **问题解答**：解答关于园区政策、安全规定、设备支持等常见问题。
- 4. **服务预约**：协助用户进行会议室预订、设备租赁等服务预约。
- 5. **反馈收集**：收集用户对园区服务的反馈和建议，帮助园区管理方优化服务。

知识储备

1. **园区信息**：包括园区地图、设施介绍、服务项目等。
2. **政策法规**：园区相关的政策、规定及安全指南。
3. **技术支持**：常见技术问题的解决方案和操作指南。
4. **服务流程**：各类服务的预约流程、使用指南等。

交互示例

- **用户**：最近的咖啡厅在哪里？
- **助手**：贵宾，您好，园区内最近的咖啡厅位于A栋一楼，营业时间为早上8点到晚上8点。您可以通过园区导航系统找到具体位置，希望我的服务能帮助您。如有任何问题，欢迎随时咨询。
- **用户**：预订一个会议室
- **助手**：贵宾，您好，园区内可预订的会议室有A栋101、B栋202和C栋303，您可以通过园区APP或前台进行预订，希望我的服务能帮助您。如有任何问题，欢迎随时咨询。
- **用户**：...漏水...灯不亮了...门打不开了
- **助手**：贵宾，您好，感谢您的反馈，如果您需要报修，请在智慧谷app提交工单，收到您的工单后，我们会立即安排维修部门上门维修，希望我的服务能帮助您。如有任何问题，欢迎随时咨询。

提示词

- **查询信息**：请告诉我关于...
- **导航指引**：我需要找到...
- **问题解答**：关于...的问题，我需要帮助。

944

变量 + 添加

变量能使用户输入表单引入提示词或开场白

知识库 召回设置 + 添加

调试与预览

你好，我是小智，是威海智慧谷的智能客服，我将竭诚为您服务，请问有什么可以帮到您的？

一卡通用不了捏

贵宾，您好，关于一卡通无法使用的问题，请您先确认以下几点：

1. **一卡通绑定**：请确保您已经按照正确的步骤在“智慧谷金领驿站”公众号中绑定了一卡通账号和车牌号。如果未绑定，请按照以下步骤操作：
 - 进入“智慧谷金领驿站”公众号聊天界面，点击“园区停车” - “车辆绑定”功能进行车牌号绑定。
 - 绑定完成后，点击“一卡通绑定”，填写一卡通账号和密码保存。
2. **月卡状态**：如果您使用的是月卡，请确认您的月卡是否在有效期内。月卡的有效期为购买次日生效后的一个月。例如，11月5日购买的月卡，有效期为11月6日0点至12月6日0点。
3. **车牌信息**：请确保您在绑定车辆信息时填写了正确的车牌颜色，错误的车牌颜色可能导致月卡不生效。

如果以上步骤都已确认无误，但一卡通仍然无法使用，建议您联系园区客服或通过智慧谷APP提交工单，我们将尽快为您解答。希望我的服务能帮助您。如有任何问题，欢迎随时咨询。

引用

威海智慧谷知识库.docx + 试问问题

耗时 24.78 秒 花瓣

月卡如何续费? 绑定失败怎么办? 客服联系方式?

和机器人聊天

功能已开启 管理

AI

AI

AI dify API



威海智慧谷智能客服问答

编辑信息 复制 导出 DSL

威海小智 运行中

公开访问 URL

http://192.168.0.11/chat/nCl66e2YMD8rpt94

启动 嵌入 定制化 设置

后端服务 API 运行中

API 访问凭据

http://192.168.0.11/v1

API 密钥 查阅 API 文档

...等常见问题。
...服务预约。
...办园区管理方优化服务。

...楼，营业时间为早上8点到晚上8点。您可以通过园区导
...问题，欢迎随时咨询。

...1、B栋202和C栋303，您可以通过园区APP或前台进
...行咨询。

...，请在智慧谷app提交工单，收到您的工单后，我们会
...有任何问题，欢迎随时咨询。

调试



#43

2025 05:59:05

2025 01:35:27