

算法设计文档

项目名：华为云智能问答系统

版本：1.0

编订：王飞鸿

日期：2018-07-20

部署地址：<http://58.87.125.79:5000/chat>

1. 引言

1.1 编写目的

本文档的目的是详细地介绍华为云智能问答系统的算法及设计思路，对本系统的各模块、程序、子系统进行了实现层面上的要求和说明。使系统开发人员清楚认识本系统所需要实现的功能以及功能模块的划分、数据库的表结构、设计思想。

1.2 背景

本文档介绍的产品是华为云智能问答系统，该软件面向所有需要使用华为云产品的人群。该软件基于华为云网站页面作为原始数据由都柏林科技有限公司团队开发。主要方便用户咨询华为云产品的使用，并节省了客服人力。

1.3 软件功能

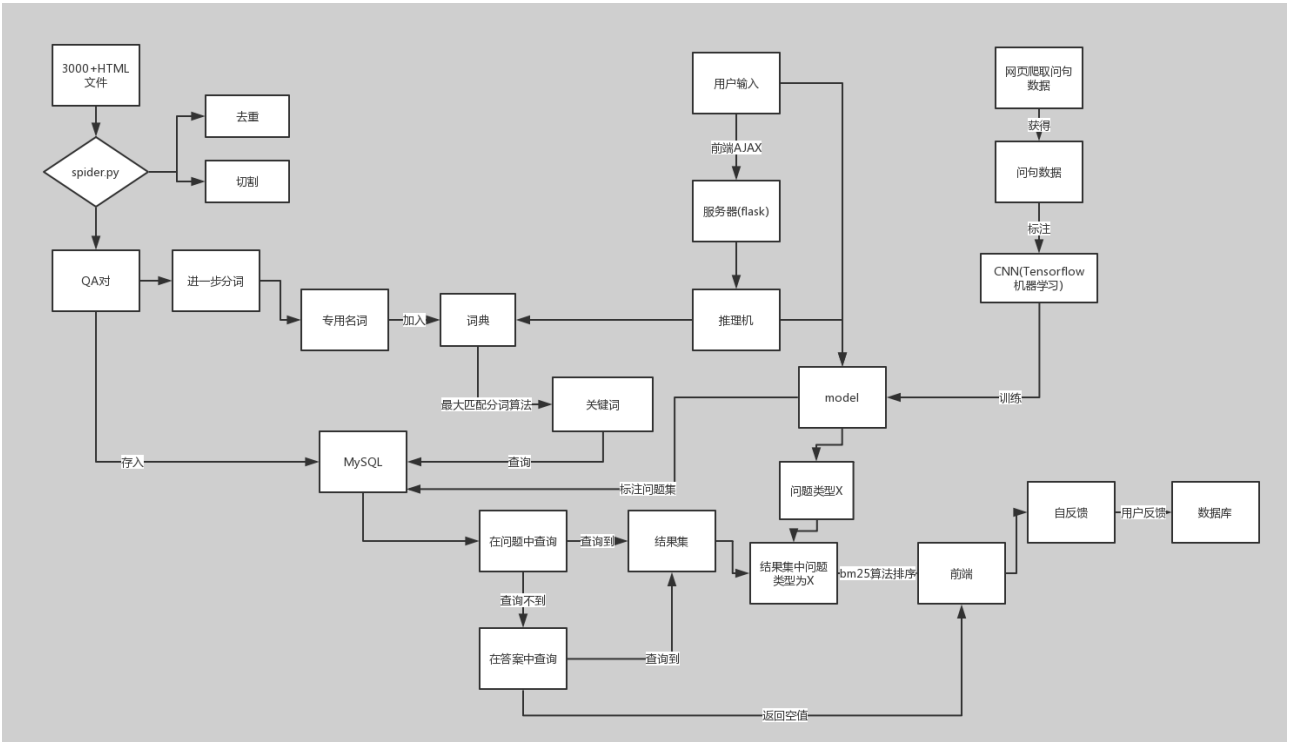
后台实现产品文档(HTML 文件)录入、知识库构建根据录入文档实现“QA 对”自动化生成，知识库实现基本 QA 对删除、增加、查询等操作功能。

前台实现 QA 对话界面，该界面可以基于用户提问，自动连接后台、并从知识库寻找答案，并呈现给用户，前台问题可以是由主题、关键词、短语构成。针对不同的问法的同种问题，可自动识别。

后台搜索引擎可根据问题的模糊程度返回多条答案或一条精确答案。在每条答案的基础

上，增加自反馈功能，增强用户体验。

2.系统结构



本文的系统设计结构如图所示。

3. 算法设计

本系统（后简称 WEBQA 系统）具有以下几个特征：

- （1）自动根据数据集生成知识库，通过文档结构树和自行设计的爬虫框架的方法提取 QA 对。无需从大量篇幅中做语句分析或是实体识别，从而提高系统效率和准确度。
- （2）自动收录专业名词训练词库。对用户提出的问题，更准确地提取关键信息，从而使命中问答对数增高。
- （3）利用 word2vec 模型处理中文语言特色引起的问题。
- （4）可以回答同种问题的不同问法的问句，并准确返回答案。
- （5）通过与用户的交互实现知识库自我更新，使其提供用户友好性。

3.1 系统模型

构建 WEBQA 系统模型，主要包括两个部分，一是自动生成知识库，建立问题与答案的

关系；二是创建推理机，知识推理需要知识库所存储的知识作为基础，不同的知识表达方式在一定程度上决定了特定的知识运用方式。

WEBQA 系统的基本设计思路分两个步骤，即知识库构建和推理机：

知识库构建是自动生成的,本通过自设计框架、借助标签从网页中提取 QA 对，并通过文本结构树的思想整合答案的类型。

推理机处理模块通过问题分析、问题分类取得用户的期望问题类型和关键词。用 word2vec 模型进行相似词辅助搜索,已解决中文语言特色带来的问题 。在此基础上，用关键词对知识库中总结的概述词和答案进行文本相似度计算，得到目标答案。

3.2 自动生成知识库

目前，自动生成知识库已作为实现实用的问答系统的一个基本组成部分。WEBQA 系统针对网页形式的数据集自动生成知识库，主要通过网页的标签进行答案的提取，并采用创新的文本结构树算法对答案和问题进行分割，如图 1 所示。

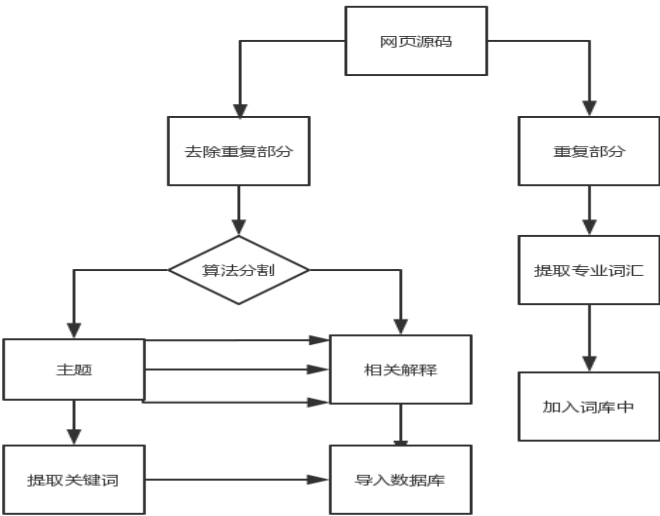


图 1 知识库自动生成流程

WEBQA 系统基于网站生成知识库，利用网页中已有的标签、结构进行 QA 对提取。以

[最新活动](#)
[产品](#)
[解决方案](#)
[EI企业智能](#)
[云市场](#)
[社区区](#)
[合作伙伴](#)
[支持与服务](#)
[Q](#)
[中文 \(简体\)](#)
[控制台](#)
[登录](#)

为您推荐：[弹性云服务器](#) [关系型数据库](#) [主机安全服务](#) [云硬盘](#) [虚拟私有云](#)

计算	存储	网络	安全	管理与部署
弹性云服务器	对象存储服务	虚拟私有云	Anti-DDoS流量清洗	云监控服务
GPU加速云服务器	云硬盘	弹性负载均衡	DDoS高防IP	云日志服务
FPGA加速云服务器	云硬盘备份	NAT网关 New!	Web应用防火墙	统一身份认证服务
裸金属服务器	云备份服务	弹性公网IP New!	漏扫扫描服务	云审计服务
专属云	专属存储服务	云专线	企业主机安全 New!	标签管理服务 New!
专属主机 New!	CDN	虚拟专用网络	密评管理服务	云报表服务
弹性伸缩	云文件服务	云解析服务	数据库安全服务 New!	云目录服务
镜像服务	数据快道服务	EI企业智能	安全体检服务	资源模板服务 New!
云容器引擎	专属企业存储服务	机器学习服务	态势感知 New!	企业应用
云容器实例 New!	数据库 New!	深度学习服务	软件开发服务 New!	云桌面
函数服务	云数据库 HWSQL New!	图引擎服务	项目管理	云通信
	云数据库 MySQL	实时流计算服务	代运维	即时通信
应用服务 New!	云数据库 PostgreSQL	MapReduce服务	流水线	会议
微服务云应用平台	云数据库 SQL Server	数据查询服务	代码检查	联络中心
应用编排服务 New!	文档数据库服务	表格存储服务	漏洞构建	语音通话
函数工作流	分布式缓存服务 Redis	数据仓库服务	部署	消息&短信 New!
容器镜像服务	分布式缓存服务 Memcached	云搜索服务	测试管理	隐私保护通话 New!
微服务引擎 New!	分布式数据库中间件	文字识别 New!	发布	物联网
消息通知服务	数据复制服务 New!	图像识别	移动应用测试	IoT平台
分布式消息服务		智能物流	CloudIDE	
应用性能管理	迁移	视频		
API网关 New!	对象存储迁移服务	媒体转码 New!		
云性能测试服务 New!	云数据迁移 New!	视频点播 New!		
区块链服务 New!	数据复制服务 New!	娱乐视频云服务 New!		
应用运维管理 New!				

近期用户关注度比较高的热门活动提醒：

词库应尽可能贴合知识库的方向。比如“云服务器是什么？”，用普通词库分出的关键词是“云”和“服务器”，而想要的是“云服务器”这个词，分词错误会影响后来的一系列操作，所以应采用相关的词库，将描述的词加进词库。把词频设为中等频率，因为过大的词频会影响推理机部分 TF-IDF 算法提取关键词，过小的词频又不能分出。

本系统的爬虫放弃对网页标签结构的研究, 即把一个网页当成纯文本形式, 将所有 h 标

为有效管理知识库和,我运用文档结构树的知识,将每条 QA 对一步步细分,以达到方便管理和查询的目的。如“帮助中心>机器学习服务>最佳实践>附录>修订记录”,这是一条 QA 对中问题的格式。是每层一步步细化下来的。在此基础上我设计数据库表用于存储 QA 对。在表设计上我考虑到了效率,将用于搜索引擎中的分词,去除停用词等文本预处理在生成知识库步骤中完成。在后期就不用再次操作了。答案是连同 html 的 tag 标签一起爬取下来的,

所以在前端返回答案时，可以看到表格，列表等网页样式的答案而不是单一的文字。将知识库存入数据库也满足了赛题中“实现 QA 对的增删查改”的要求。

数据库表格成分为：

- 标准问题：通过规则组合关键词生成（基于文档结构树技术）。
- 答案：带网页标签的答案
- url：数据来自的网页
- 描述词：此数据的描述，由网页中提取出（如 3.1 图）。
- 问题词组：通过分词（全模式分词算法）和去除停用词等步骤得到的词组。用于计算文本相似度。
- 答案词组：通过分词（全模式分词算法）和去除停用词等步骤得到的词组。用于计算

```
mysql> mysql> select descs from QA limit 30;
+-----+
| descs                                     |
+-----+
| -备案中心 &gt;-备案帮助-备案简介          |
| -备案中心 &gt;-备案帮助-政策法规          |
| -备案中心 &gt;-备案帮助-新手指引          |
| DevOps解决方案-DevOps解决方案            |
| DevOps解决方案-DevOps解决方案架构图      |
| DevOps解决方案-业务架构图                |
| DevOps解决方案-云容器引擎                |
| DevOps解决方案-产品推荐                  |
| DevOps解决方案-合作伙伴展示              |
| DevOps解决方案-图形化编排工具            |
| DevOps解决方案-基于华为云Docker容器开发  |
| DevOps解决方案-基于软件开发云的研发平台  |
| DevOps解决方案-容器镜像仓库              |
| DevOps解决方案-容器集群管理              |
| DevOps解决方案-应用场景                  |
| DevOps解决方案-开发测试生产环境          |
| DevOps解决方案-弹性伸缩                  |
| DevOps解决方案-架构优势                  |
| DevOps解决方案-流畅体验                  |
| DevOps解决方案-自动化运维管理            |
| DevOps解决方案-高效负载                  |
| null-上海                                |
| null-云南                                |
| null-内蒙古                              |
| null-北京                                |
| null-吉林                                |
| null-四川                                |
| null-天津                                |
| null-宁夏                                |
| null-安徽                                |
+-----+
30 rows in set (0.00 sec)

mysql>
```

图 3.1 数据库描述词样例

文本相似度。

```
mysql> desc QA;
```

Field	Type	Null	Key	Default	Extra
normal_question	varchar(255)	NO	PRI	NULL	
type	varchar(255)	YES		NULL	
url	varchar(1000)	YES		NULL	
answer	mediumtext	YES		NULL	
descs_words	mediumtext	YES		NULL	
answer_words	mediumtext	YES		NULL	

5 rows in set (0.00 sec)

图 3.2 数据库表信息

备案中心	>- 新手指引	>- 核验常见问题	特殊要求:	定义	support.huaweicloud.com/beian/new_hycjwt.html
list>江西、新疆地区的要求: 个人性质备案的客户, 需要在核验单网站负责人签字处加盖个人手印。</p><p class="l_list">广东个人核验单与其他省份不同 (广东个人样例), 广东备案主体为个人时, 核验单中需手写以下内容: "本人已履行网站备案信息当面核验手续, 承认网站备案信息和核验记录真实有效, 承诺本网站是个人网站, 未含企业、单位等非个人网站的信息, 承诺网站备案信息一旦发生变更, 将及时进行更新, 填报虚假备案信息、未履行备案变更手续、超出备案项\$提供服务的, 愿承担关闭网站并注销备案 (列入黑名单) 等相应处理。"</p><div class="moreImgBox l_img"><div class="img"><em class="expand"></div></div></div><div class="rec-item"><p class="caption">材料邮寄及现场核验地址<i class="foldIcon"></i></p><div class="content"><p class="info">收件地址: 广东省深圳市龙岗区坂田街道天安云谷2栋</p><p class="info">收件人: 备案专员 (收)</p><p class="info">联系电话: 4000-955-988转7</p></div></div><div class="rec-item"><p class="caption">办理拍照<i class="foldIcon"></i></p><div class="content"><p class="l_list">根据工信部要求, 备案时还需提供网站负责人当面核验照片。在系统初审通过后, 将有邮件通知您准备网站负责人的幕布照 (可到拍照点拍照或申请邮寄拍照专用幕布自行拍照), 将照片电子版发送到: hwclouds.ba@huawei.com 自行拍照指导--收到华为云的背景幕布后, 请按以下要求进行操作: </p><p class="l_list">拍照人必须与网站负责人为同一人。</p><p class="l_list">负责人需站在背景幕布中间位置, 只拍上半身即可。</p><p class="l_list">整个照片背景必\$都是华为云幕布, 且为蓝色, 显示效果: 幕布字迹清晰。</p><p class="l_list">请您避免身着红色或者蓝色上衣进行拍照</p></div></div><div class="rec-it\$特殊要求, 中心, 特殊, >, 常见, 要求, 核验, 新手, 指引, 备案, 常见问题					
黑名, 虚假, 可到, 只, 愿, 以下内容, 部, 红色, 发生, 非, 谷, 已, 进行, 人性, 背景, 生变, 工, 黑名单, 收, 求, 进, 手续, 专员, 单位, 拍上, 上半, 还, 将, 电子, 系统, 现场, 名单, 天安, 签字, com, 性质, 同位, 未, 准备, 岗区, 出备, 田, 服务, 相应, 都, 材料, 初审, 深圳, 信, 江西, 未, 含, 项目, 龙岗区, 单, 街道, 需要, 通知, 送到, 站, 龙岗, 新疆, 履行, 过后, 上半身, 地区, 云, 处, 需, 有, 负责人, 加盖, 要求, 字迹, 照片, 电话, 效果, 中, 提供, 整个, 承认, 注销, 华为, 避免, 列, 个人性, 责人, 实有, 人手, 拍照, 承诺, 请, 收件, 信息, 上半身, 地区, 云, 处, 需, 有, 负责人, 加盖, 要求, 字迹, 照片, 电话, 效果, 中, 提供, 整个, 承认, 注销, 华为, 避免, 列, 个人性, 责人, 实有, 人手, 拍照, 承诺, 请, 收件, 信息, 更新, 显示, 人为, 是, 以下, 变更, 有效, 收到, 超出, hwclouds, 地址, 申请, 填报, 位置, 幕布, 操作, 即可, 坂, 同一, 与其, huawei, 联系, 蓝色, 广东, 您, 指导, 必须, 邮寄, 收, 入, 联系电话, 新疆地区, 身着, 发送到, 转, 人样, 上衣, 且为, 深圳市, 办理, 手写, 照, 省份, 不同, 时, 主体, 负责, 关闭, 核验, 照人, 内容, 清晰, 例, ba, 网站, 专用, 一旦, 手, 印, 个人, 电子版, 备案, 记录, 人					

图 3.3 单条数据示例

在创建数据过程中，对爬取到的答案和问题做预处理为推理机部分做准备，进行去除网页标签，分词，去除停用词等操作。其中分词算法较初赛改良为全模式分词。这种分词的特点是把一句话中能分出的词全部分出，优点可以体现在：确保知识库中的词可以匹配上推理机中的词。例如“云服务器”和“服务器”这样的词在相似度算法中并没有匹配上，但实际上这表达了同一种意思。最后将这些分好的词写入每条数据库的 descs_words 和

answer_words 的属性中并作为推理机中用于匹配答案的信息。提前将答案预处理,可以节省后期在加载数据时浪费的时间和资源.

3.3 推理机处理模块

3.3.1 构建推理机

WEBQA 系统的答案提取由常用知识库的简单答案匹配和文本相似度匹配构成。

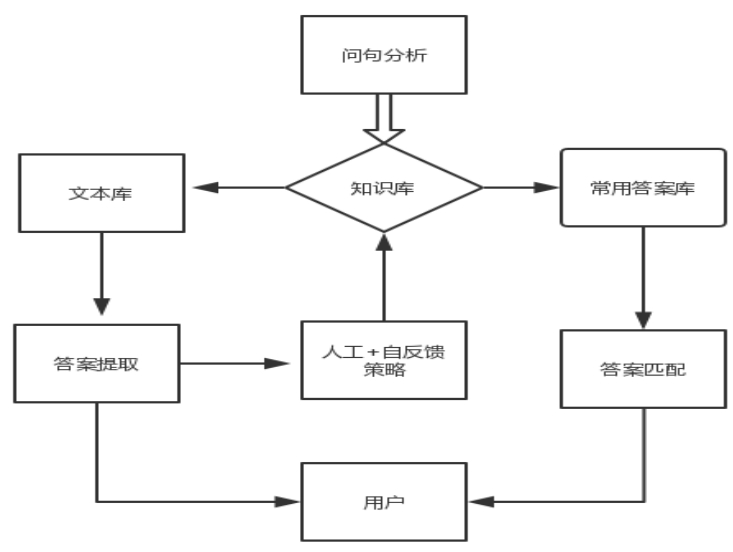


图 4 推理机构建流程图

常用知识库构成包括：文本概述和答案构成。文本概述是通过爬取网页建立的概述词，在一定层面上想比文本相似度更精确，更人性化。因为它本身就源自于人工的总结，所以在答案提取优先使用常用知识库，其实际查询规则是基于文本相似度。

在推理机构建中着重创新解决了以下问题：

判断网页内容是否与用户查询相关，这依赖于搜索引擎所采用的检索模型。检索模型是搜索引擎的理论基础，是对查询词和文档之间进行相似度计算的框架和方法。如图 5 所示，

检索模型所在搜索引擎系统架构位置。

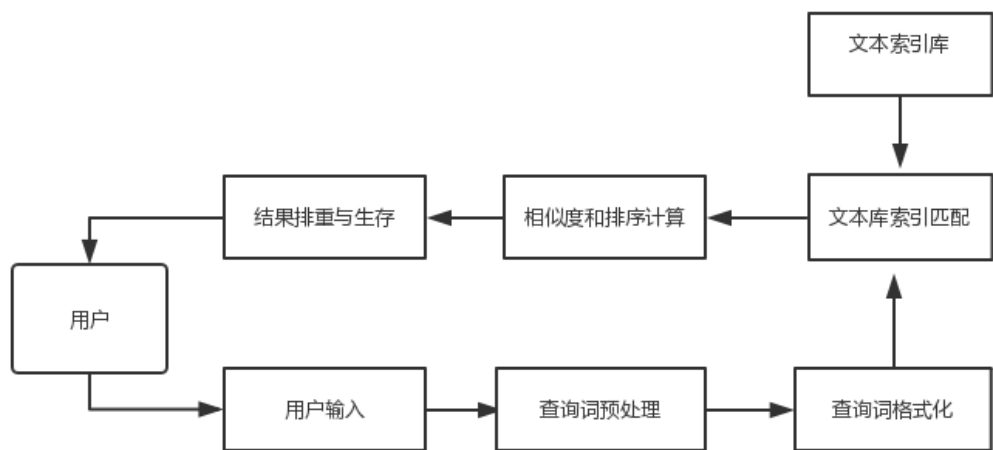


图 5 文本相似度的检索模型运行流程

本系统使用 **BM25** 模型用作计算文本相似度，使用改良的 **BM25** 算法对答案问题同时进行打分。通过分析数据，还提高了的文本长度影响分数的权重，因为概括性强的问题往往在比较短的答案中出现。

摒弃基于 **TF-IDF** 算法提取用户的问题中的关键词。这里我对数据进行分析，发现不适合用 **TF-IDF** 算法，因为每个词汇都十分重要，可替代性不强，所以删除了原 **BM25** 的 **TF-IDF** 部分。**TF-IDF** 算法的主要思想就是：如果某个词在一篇文档中出现的频率高，并且在语料库中其他文档中很少出现，则认为这个词代表问句的主题，即关键词。由于 **WEBQA** 词库中已有大量相关网站领域性的词，所以如果用这种方法去提取关键词会丢弃重要的词（同时出现频率也很高的词）。如“云专线”出现的频率很高。但是如果从关键词中排除它就很难定位答案。在用用户输入的关键词匹配答案时，采用同时在答案和问题打分的方式，同时赋予问题更高的权重，最终将两部分打分相加再排序返回前 3 条高分答案。因为答案和问题都有可能包含用户想要的到的信息。而赋予问题更高的权重，是因为问题是人工总结出来的，其概括性更强，更能代表这个 QA 对的主题。本系统还添加了自反馈系统，在返回问答系统回答的问题时，会产生与用户的交互，用户可以根据返回的答案满意与否，点击答案下的按钮，如满意，系统会自动将用户输入的关键词加入到那条数据的 `descs_words` 属性中。下次搜索时，更新过的关键词就会影响到得分。从人的角度增加了知识库的可用性。

通过分析中文语言特色，发现用户会使用相近的意思的词语代替知识库中的词搜索，因此做了词向量模型的训练，帮助解决这样的问题。词向量模型用于计算关键词的相似度，帮助比较每种近义词搜索的得分。但需要注意的是，专有名词并不能用近义词去替代搜索，因

为其可替代性很低，如果选择近义词，就不能代表原来意思。如“服务器”和“电脑”会被词向量模型标记为近义词。而在专业领域，这两个词词义差异很大。

本系统还尝试用 CNN 算法训练分类器模型，也是为了增加准确度。最终整合达到最高的准确度。在问句分析中运用问题分类。经统计，用户在提出问题时，习惯性用口语形式，问句形式不规范，进而无法准确分析作者的期望答案类型。因此，在知识库中，对一类问题用规范统一的问法保存。在推理机中，通过机器学习训练出的分类模型对用户输入的问句进行分类，用问句的类别去匹配数据库中这类别的答案，以更好的提供回答。用于训练问题分类模型的数据来自网络，通过爬虫技术获取，并自行标注，利用 CNN 算法训练。如表 1 所示，

用户输入“云服务器多少钱”转换为知识库中更为专业性的问句“云服务器的价格”。

专业性问句	口语性问句 1	口语性问句 2	问句分类
云服务器价格	云服务器多少钱？	云服务器多贵？	价格
云服务器功能	云服务器是什么？	云服务器有什么用？	功能

表 1 问句分类样例

4. 测试

测试使用每个网页中的“引导标签”，通过基于规则的方式生成测试数据，首先筛除部分对问题生成无用的标签，如：“帮助中心”、“概览”、“产品简介”、“快速入门”等等，接着通过一些规则自动生成问题，规则包括：

1. “引导标签”的最后一个标签以“什么”、“怎么”或“如何”开头，则不进行语法更改，直接使用最后一个标签作为生成的问题，如：

帮助中心>云服务器备份>产品介绍>什么是云服务器备份

->

“什么是云服务器备份？”

2. 最后一个“引导标签”以动词开头且这个动词不做形容词用，用“怎么”加最后一个标签生成问题，如：

帮助中心>弹性云服务器>快速入门>注册公有云

->

“怎么注册公有云？”

3. 最后一个“引导标签”以动词做形容词开头（这些词组包括“使用限制”、“计费方式”、“操作指南”、“开发指南”），则以主标签加“的”加最后一个标签加“有哪些”生成问题，如：

帮助中心 > MapReduce 服务 > 开发指南

->

“MapReduce 服务的开发指南有哪些？”

测试使用的规则不止这 3 条，但是这 3 条比较独特或常见，测试选取了随机 1000 个网页并生成了问题，这些问题在 [作者的 Github](#) 上可以看到。测试从中抽取了 100 条不重复的合理问题输入机器进行了黑盒测试，并人工判断输出的合理性。结果比较令人满意，100 个问题中超过 90 个能给出合理的答案。