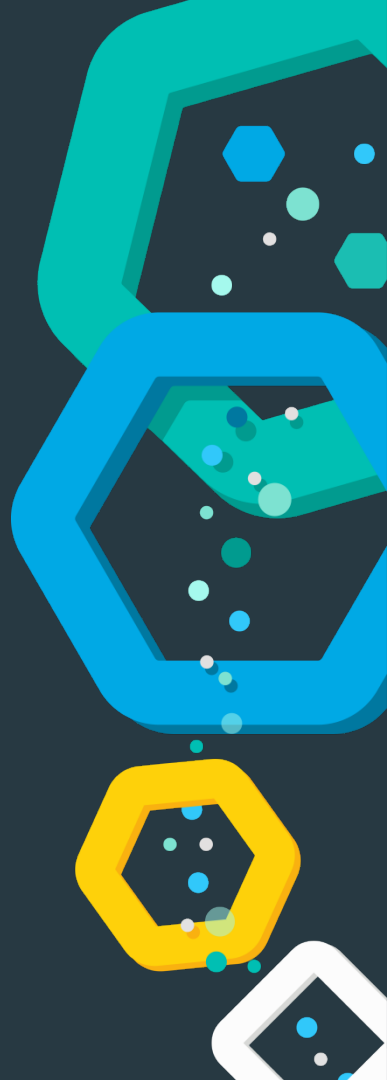


又一个爬虫

曾勇 (Medcl)



什么是爬虫

- 又叫 Robot、Bot 或 Crawler
- 简单来说
 - 自动探索网站
 - 帮你访问整个站点
 - 自动为你收集网站信息
 - 自动更新
 - 抽取处理网页内容
 - 存储索引和快照
 - 。 。 。

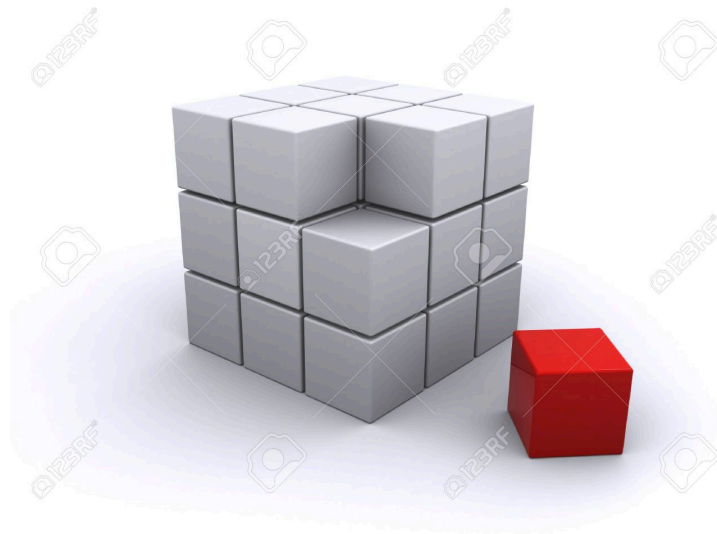


为什么造这个轮子？

现在已经有很多开源的爬虫了: *Scrapy*, *Nutch*, *Heritrix* 等等 . [1,2,3]

但是！

- 大多仅仅是一个爬虫框架！ *ES* 与 *Lucene*
- 需要熟悉各种与爬取任务本身无关的知识！
- 开发和部署的环境复杂，痛苦！
- 太重了！
- 分布式、管理、监控、扩展复杂！
- 结果就是一大堆凌乱无法维护的脚本，或是一大堆技术拼凑的大杂烩。



1.<http://bigdata-madesimple.com/top-50-open-source-web-crawlers-for-data-mining/>

2.<https://github.com/BruceDone/awesome-crawler>

3.<https://gitee.com/explore/starred/spider>

所以

- 狗爬, Gopa

Golang + pá chóng (爬虫)

<https://github.com/infini7byte/gopa>

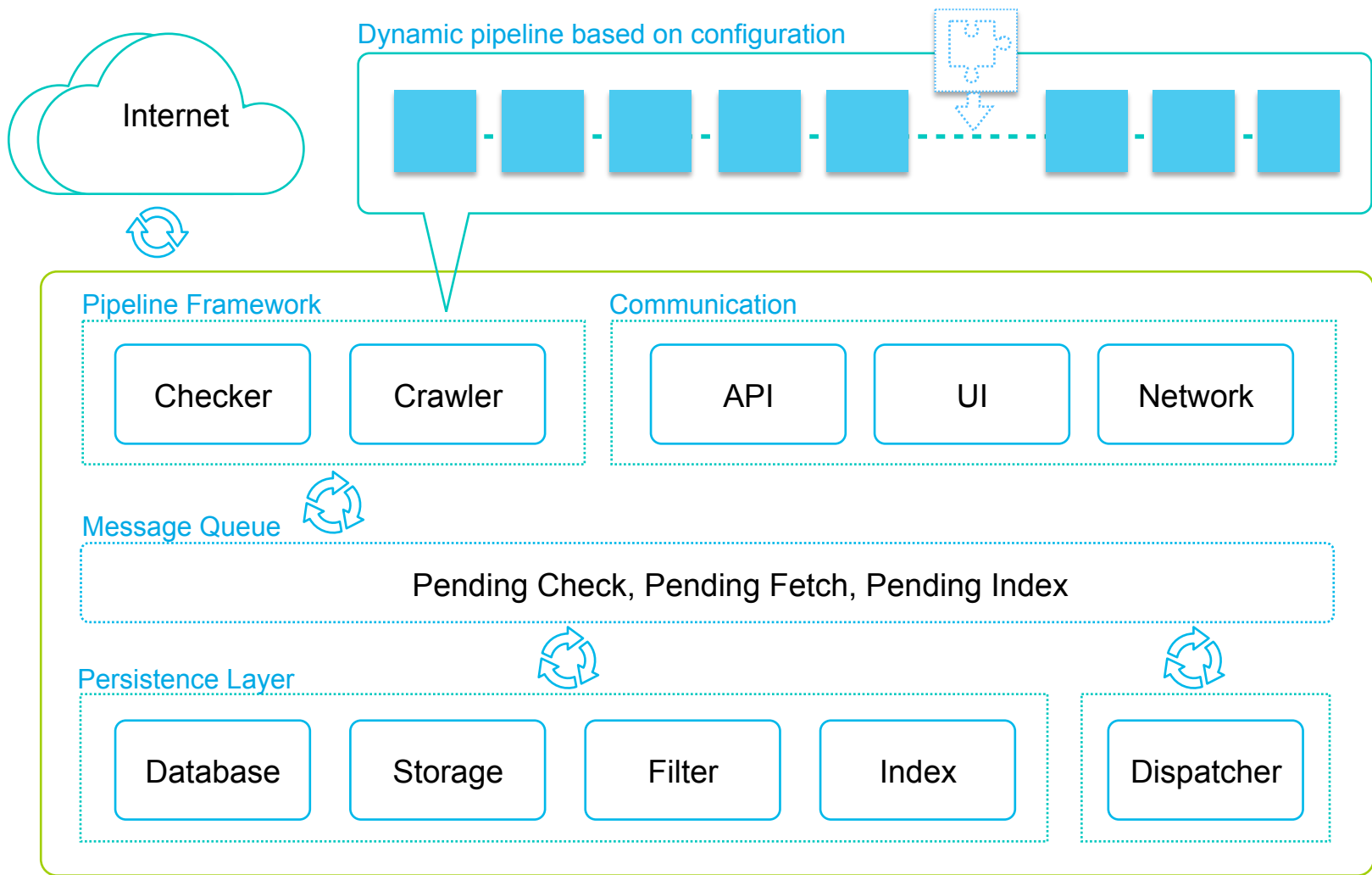


目标

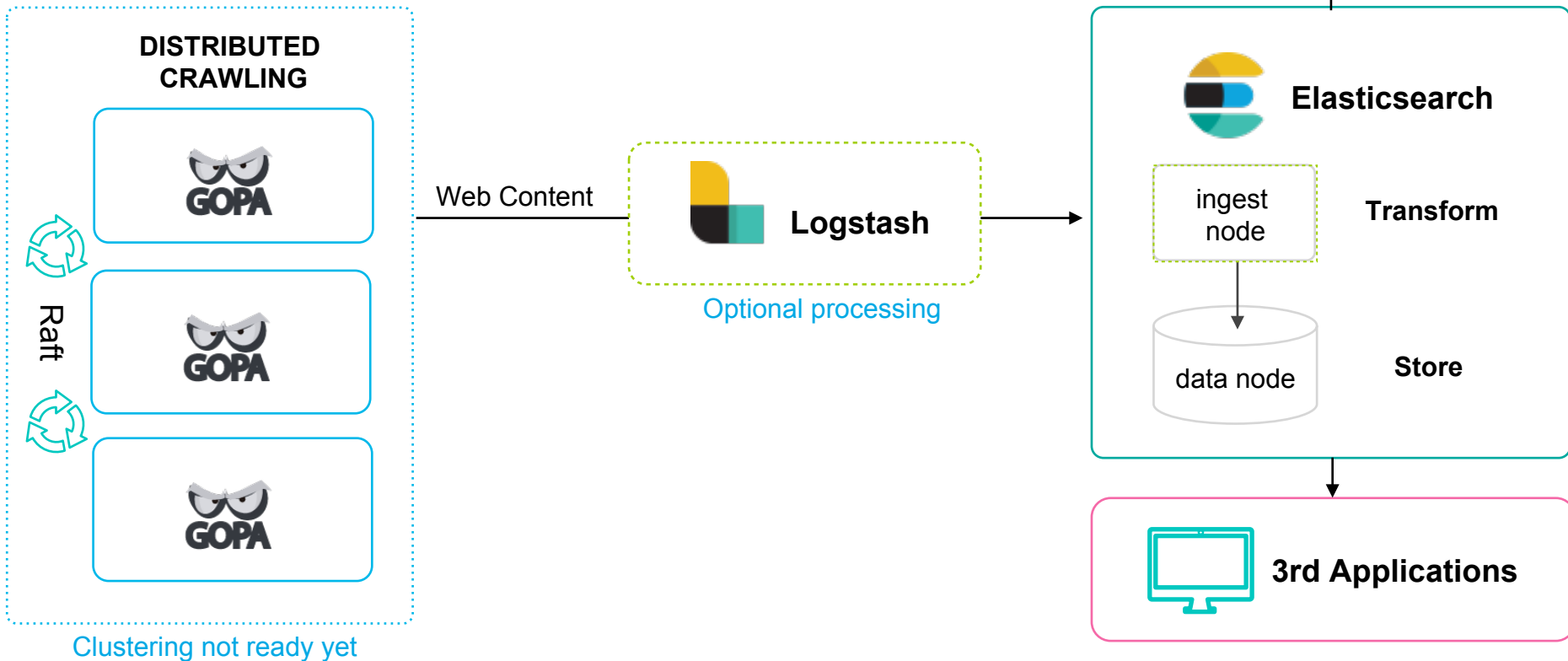
- 轻量级，内存占用 $< 100\text{MB}$
- 易于部署，无运行时和环境依赖
- 方便使用，无需编程和脚本技能
- 开箱即用
- 提供 RESTful API 和 UI
- 简单，可伸缩，可扩展

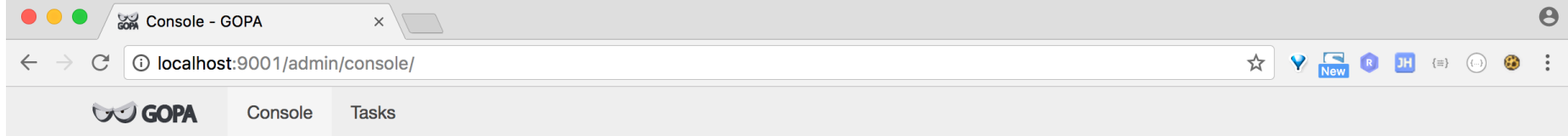


DEMO

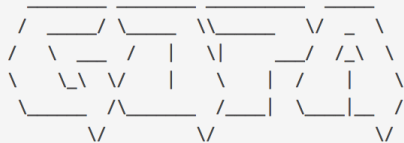


Elastic Integration Overview





Connection established.



[gopa] 0.9.0_SNAPSHOT

///last commit: 08232b4, Sat Sep 23 19:11:37 2017 +0800, medcl, normalize content-type ///

Press <s> to quick input; Type HELP for more details.

Realtime logging

DEBUG ▾

FilePattern, eg: crawler*.go

FuncPattern, eg: runPipeline*

MessagePattern, eg: *timeout

▶ START

■ STOP

```
[08:51:48] [DEBUG] [url_normalization.go:321] [Process] finished normalization,, https://elasticsearch.cn/topic/%E7%A4%BE%E5%8C%BA%E6%B4%BB%E5%8A%A8, /topic, /社区活动.htm
[08:51:48] [DEBUG] [url_normalization.go:188] [Process] domain mismatch,elasticsearch.cn vs www.elastic.co
[08:51:48] [DEBUG] [checker.go:181] [execute] ignored url, https://elasticsearch.cn/article/189
[08:51:48] [DEBUG] [url_normalization.go:321] [Process] finished normalization,, https://elasticsearch.cn/topic/Elastic%7BON%7D17, /topic, /Elastic{ON}17.html
[08:51:48] [DEBUG] [index.go:270] [Search] search: http://dev:9200/gopa-task/_search
[08:51:48] [DEBUG] [webhunter.go:255] [ExecuteRequest] let's: POST, http://dev:9200/gopa-task/_search
[08:51:48] [DEBUG] [save_snapshot.go:72] [Process] save snapshot to db, url:https://elasticsearch.cn/article/254,domain:elasticsearch.cn,path:/article,file:/254.html,save
```

Total 5

[www.elasticsearch.cn\(10\) es-guide-preview.elasticsearch.cn\(1\) elasticsearch.cn\(2497\) conf.elasticsearch.cn\(2\) grok.elasticsearch.cn\(1\)](#)

Total 1

URL	LastUpdate	NextCheck	Status
http://grok.elasticsearch.cn	~1 min ago	~24 hours later	success

1

[GitHub](#) | [Issues](#) | [Releases](#) | [Changelog](#)

Licensed under [Apache License, Version 2.0](#).

Dashboard / Editing Gopa Dashboard (unsaved)

Save Cancel Add Options Share Reporting < Last 15 minutes >

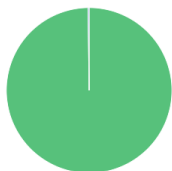
Search... (e.g. status:200 AND extension:PHP)

Add a filter +

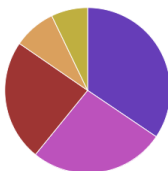
Task Updated



Task Breath

0
1

Task Depth

4
5
3
6
2

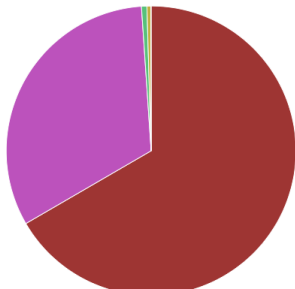
Domain status

Domain

Links

elasticsearch.cn	7,949
www.elasticsearch.cn	10
conf.elasticsearch.cn	2
grok.elasticsearch.cn	2
es-guide-preview.elasticsearch.cn	1

Task Status

3
5
0
2
4

Snapshot Status

File Type

File Count

text/html; charset=utf-8	5,304
text/html	3

Path Tags

The container is too small to display the entire cloud. Tags might be cropped or omitted.

category-12_sort_type-hot_day-7 /category-13_is_recommend-1 /category-15_is_recommend-1

/category-10_sort_type-hot_day-7 /category-10_is_recommend-1 /category-13_sort_type-unresponsive

/category-12 /_question/1 /topic /_question/68 /category-11_sort_type-unresponsive

/category-17 /category-1 /_question/12 /event /category-11 /category-16

/category-13 /_question/7 /category-14 /_question/7 /category-15

/category-11_is_recommend-1 /book/elasticsearch_definitive_guide_2x /account/register /category-18

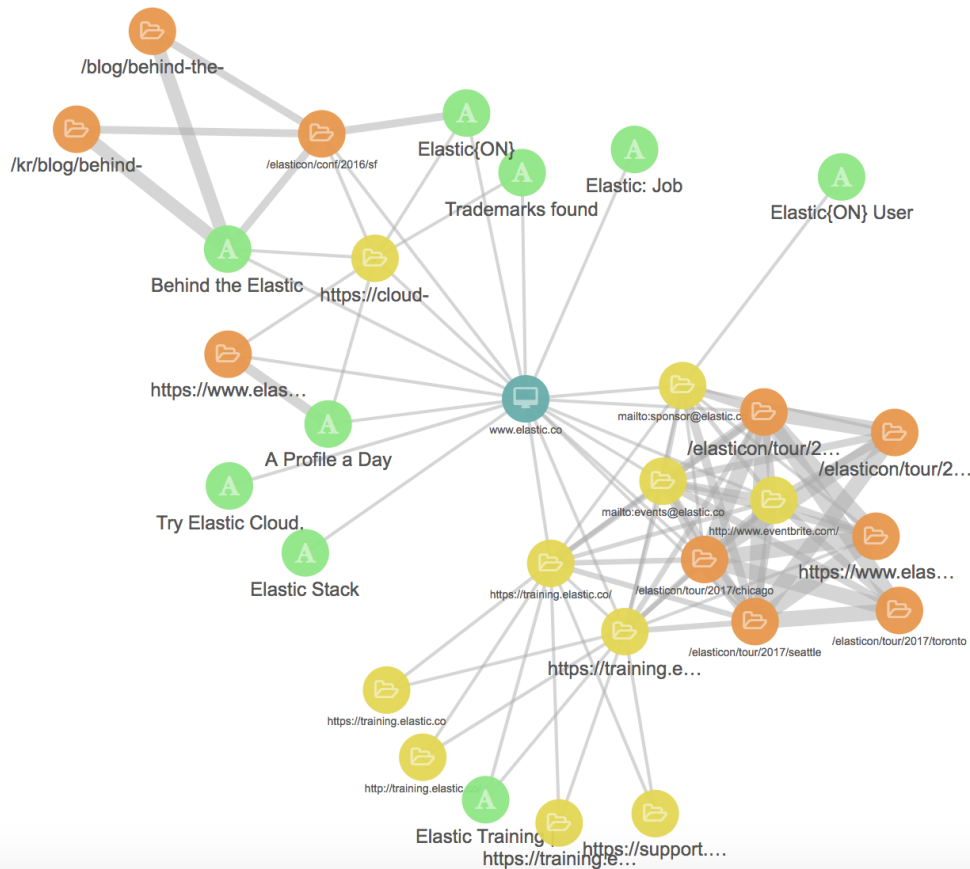
/category-12_sort_type-hot_day-7 /account/find_password /category-10_sort_type-unresponsive

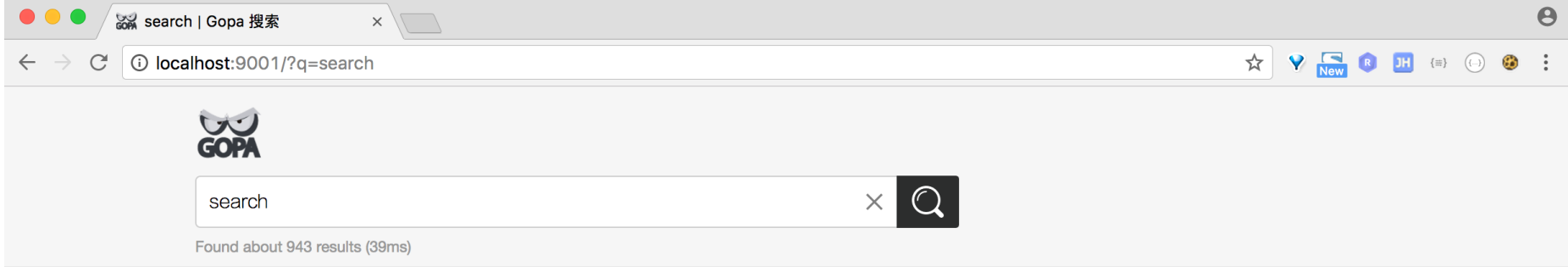
/category-11_sort_type-hot_day-7 /category-13_sort_type-hot_day-7

GOPA url relations with title

gopa-index

foo AND bar NOT baz





Open Source Search & Analytics

Open Source **Search** & Analytics · Elasticsearch | Elastic Questions?

<https://www.elastic.co/cn/> 1 weeks ago

Powering Data Search, Log Analysis

Powering Data **Search**, Log Analysis, Analytics | Elastic Questions?

<https://www.elastic.co/fr/products> 1 weeks ago

Powering Real-Time Search at Microsoft

Powering Real-Time **Search** at Microsoft | Elastic Questions?

<https://www.elastic.co/fr/elasticsearch/2015/sf/powering-real-time-search-at-microsoft> 1 weeks ago

Powering Real-Time Search at Microsoft

Powering Real-Time **Search** at Microsoft | Elastic Questions?

<https://www.elastic.co/jp/elasticsearch/2015/sf/powering-real-time-search-at-microsoft> 1 weeks ago

Elasticsearch: RESTful, Distributed Search

Elasticsearch: RESTful, Distributed **Search** & Analytics | Elastic Questions?

<https://www.elastic.co/jp/products/elasticsearch> 1 weeks ago

Search Template | Elasticsearch Reference [

Search Template | Elasticsearch Reference [5.6] | Elastic Questions?

<https://www.elastic.co/guide/en/elasticsearch/reference/current/search-template.html> 1 weeks ago

Search After | Elasticsearch Reference [5.0

Language

- en(869)
- zh(26)
- fr(15)
- ja(15)
- de(9)
- ko(9)

Domain

- www.elastic.co(672)
- discuss.elastic.co(225)
- training.elastic.co(24)
- localhost:8081(19)
- info.elastic.co(3)

Content Type

- text/html; charset=utf-8(895)
- text/html; charset=UTF-8(24)
- text/html(19)
- text/xml; charset=utf-8(5)

Protocol

- https(924)
- http(19)

Task Phrase

谢谢!

<https://github.com/infinittbyte/gopa>

