

Uknow InfoHub

Design Document

BIXLRSMB Team

Contents

1	Introduction	2
1.1	Purpose	2
1.2	Background	2
1.3	Definition	2
1.4	Reference	2
2	Architecture	2
2.1	Fetcher	2
2.1.1	Function	3
2.1.2	Performance	3
2.1.3	Input	3
2.1.4	Output	3
2.2	Prefilter	3
2.2.1	Function	4
2.2.2	Performance	4
2.2.3	Input	4
2.2.4	Output	4
2.3	Postfilter	4
2.3.1	Function	5
2.3.2	Performance	5
2.3.3	Input	5
2.3.4	Output	5
2.4	API Website	5
2.4.1	Function	5
2.4.2	Performance	5
2.4.3	Input	5
2.4.4	Output	5

3	Data Storage	6
3.1	Session	6
3.2	Item	6
4	Interface	6

1 Introduction

1.1 Purpose

For that there are many developers that got involved into this project, this document was wrote, in order to provide a mannual that developers can refer to when they are confused or forgot some protocol. Besides, this is also an introduction for new developers, who may join after this project was done. It will act as a overview of this project, to which new hands can refer.

1.2 Background

Uknow InfoHub is a collector that aimed to gather information. It's a project which was assigned as team homework of 20132014 Fall, Softerware Engineering class of Tsinghua University. And it was bootstraped by blxlrmb group, which has blahgeek, jiakai66, ppwwyyxx, vera, vuryleo and zxytim as members. It would be run on a VPS located in U.S., provided by Linode Ltd. co.

1.3 Definition

1.4 Reference

2 Archtecture

2.1 Fetcher

As the core component, fetcher is designed to be extremely flexible. It will be inherited into various kinds of fetchers to deal with different fetch target. All the fetchers will be executed in a cluster as a time based work, just like web spiders, keeps fetching all the time.

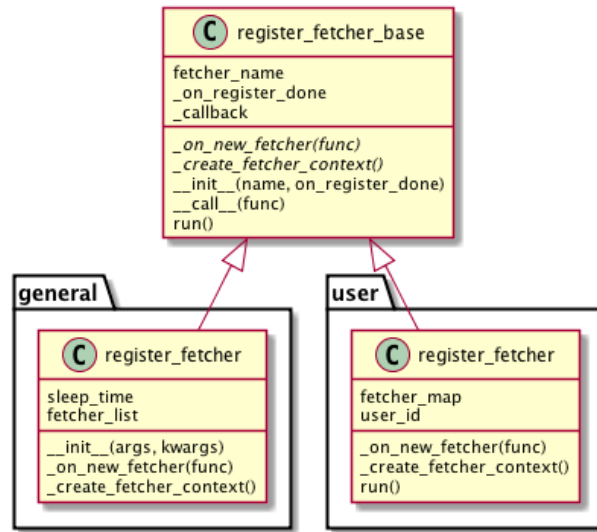


Figure 1: Fetcher

2.1.1 Function

The fetcher will be assigned to a target website or other kind of source. Each time it was called, it will return a list of informations, almost in raw text but with some basic information, such as source, origin url, fetched time and so on.

2.1.2 Performance

Fetcher shall be called at least one time per day. For hot resources such as SNS, it's fetcher must have high speed, so each time it shall return result in minute level.

2.1.3 Input

Fetcher don't need input after it's well configured. However, for augmentbilty reason, a callback function is acceptable. The callback funtion will deal with the result by default.

2.1.4 Output

Fetcher will return a FetcherContext which wraps the raw information. For default situation, it will be dealed by Prefilter.

2.2 Prefilter

What users want definitely is not raw information. Thus a prefilter will be applied on each piece of information after it's fetched. Prefilter will be called each time a fetcher returned a piece of information.

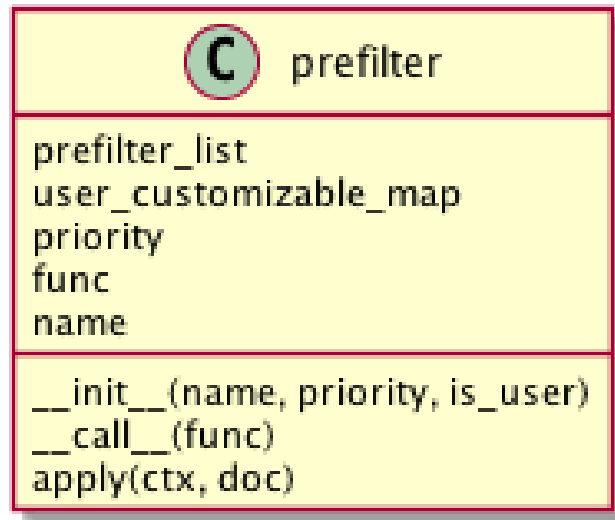


Figure 2: Prefilter

2.2.1 Function

Prefilter is the first classifier applied on the information. It can be system defined or user customized. Itself has priority levels, thus different prefilter will be executed in order. It will put different tags on piece of information, so the next step it can be easily selected.

2.2.2 Performance

It shall return result almost immediately after it was called. So it only has several seconds of time.

2.2.3 Input

A FetcherContext item which was just fetched by a fetcher.

2.2.4 Output

A FetcherContext item which has multiple tags base on its content read by prefilter. Notice that a item may be abandoned by a prefilter for it may be duplicate, redundant or illegal.

2.3 Postfilter

Users can't read all of information that we fetched. Thus a postfilter is designed to search for information that users want to read. It will be called each time user make a action that need get new information.

2.3.1 Function

Postfilter is designed aimed to select the best information that users want. It will pick out information that have the same or similiar tags with what user request. More, it will sort the information list in better order so user can get what they are most desired to get at first skim.

2.3.2 Performance

For that users must wait for it when ask for more information, it shall have extremely high Performance, as fast as possible.

2.3.3 Input

A user defined query request, has various properties and limits.

2.3.4 Output

A list of information, which fits users request well.

2.4 API Website

All things talking above is though as the backend and core items of this project. What exposed to client software is the api website. Which will run in Linode server so that every client can find it easily.

2.4.1 Function

API website will check the requester's permission and validate the request format. Then it will call the components described above to get the response, then provide them to the requester.

2.4.2 Performance

The same, users must wait for it when posting request, it shall have extremely high performance.

2.4.3 Input

Various kinds of request.

2.4.4 Output

Responses to them.

3 Data Storage

3.1 Session

Session that determining which user is currently posting request is stored in redis which is a light key-value pair database. Though only one value can be stored with one key, its high performance makes it's the best to store sessions data.

The sessions is stored automatically without any design, for a third party package will handle all of this.

3.2 Item

Item mainly is the only things that we need to store other than user's information which needed for all apps. An item will be mapped into a FetcherContext object in program, and it will be stored into database simply with all attributes by key-value pair.

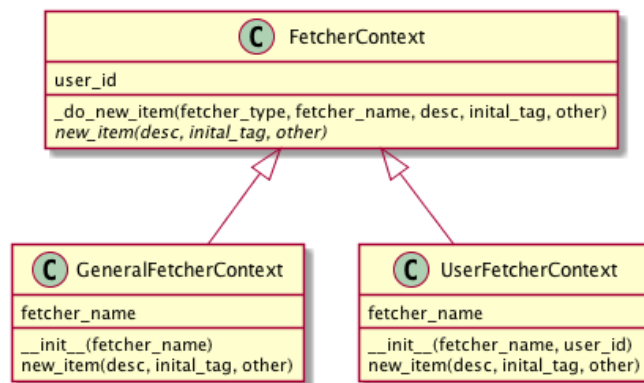


Figure 3: FetcherContext

4 Interface

This part will be discussed with caution dreadfully. For that interfaces may changed a lot when developing. Part that we can be sure of is that JSON is used when communicating between client and API server. Then the API server will call various method provided by different classes that fit the requests.

For that this project is still on the bootstrap step, and it may change everyday, we decided to mantain the interface information with wiki and API doc, other than write it done on paper. So, requests' format is still under discussion and it would change a lot. We don't think providing them now is a good idea.

More, the prototype system does not containing a client program, so there's no need for stable interface design up to now.