

Old School Cralwer

Problem ID: crawler

When building a search engine, one of the important components is the *web crawler*, meant to crawl through the web of web pages on the Internet by following links on the pages. In this problem, you are to implement a web crawler that talks to a web server in order to crawl a site.

Task

Write a program that communicates with a web server in order to find the links of all the HTML pages that can be accessed by clicking links starting at the HTTP URL `http://www.fake.webserver`.

All entities served by the server will be either:

- A valid **HTML 4.0** web page
- A valid **CSS 1.0** style sheet
- A valid **ECMAScript** script

More specifically, you should **print all the links to HTML pages a user could access** starting at the `http://www.fake.webserver` page as they were presented to a user in Internet Explorer 4 running on Windows 98.

You can safely disregard any HTTP URLs that point to hosts other than `www.fake.webserver` when solving the problem.

Server communication

Normally, you'd implement such a crawler by using DNS to lookup the IP address for the host and then connecting to it. However, since Kattis does not support networking, you will instead talk to the server by writing your requests one at a time to standard out, and reading the responses one at a time from standard in. To separate responses containing entity bodies from each other, each entity body will be terminated by with the bytes `0x0D 0x0A` – it is guaranteed that they are never appear in an entity body.

The server you will be talking to only implements **HTTP/1.0** GET requests.

If you ever print an invalid request, your program will be judged *Wrong Answer*.

Output

Once your program has crawled the web site, print a single line containing the canonical HTTP URLs of all HTML pages.