



SCHOOL OF
COMPUTER SCIENCE

Information
Management Group

HCW— HuCEL Technical Report 2, June 2009

HuCEL: Keywords Experiment II Manual

Paul Waring

Human Centred Web Lab
School of Computer Science
University of Manchester
UK

The World Wide Web contains a vast corpus of information describing a variety of events, but this information is poorly interconnected. The aim of the HuCEL project is to provide a solution to this problem by automatically generating associative links between related events. This manual describes how to re-run an investigation into how users select keywords from Web pages and how these keywords rank in relation to those generated by an algorithm.

HCW

Human Centred Web

HuCEL

The aim of the HuCEL project is to investigate how related events can be connected on the Web, in order to improve navigation of the hypertext space and enable users to serendipitously discover new information. The HuCEL Web pages may be found at: <http://hcw.cs.manchester.ac.uk/research/hucel/>.

HuCEL Reports

This report is in the series of HCW HuCEL technical reports. Other reports in this series may be found in our data repository, at <http://hcw-eprints.cs.man.ac.uk/view/subjects/hucel.html>. Reports from other Human Centred Web projects are also available at <http://hcw-eprints.cs.manchester.ac.uk/>.

Contents

1	Introduction	1
2	System Requirements	1
3	Data Storage	1
4	Installation	4
5	Flow Control	4
6	Common Errors	4
7	Summary	5
8	Associated Files	6

Human Centred Web Lab
School of Computer Science
University of Manchester
Kilburn Building
Oxford Road
Manchester
M13 9PL
UK

Corresponding author:
Paul Waring
tel: +44 (161) 275 6239
pwaring@cs.man.ac.uk

tel: +44 161 275 7821
<http://hcw.cs.manchester.ac.uk/>

1 Introduction

This manual explains how to re-run the second keywords experiment of the HUCCEL project, which is an investigation into how users select keywords from Web pages and how these keywords rank in relation to those generated by an algorithm.

The experiment can be re-run without any changes to the code supplied with this manual. However, some work will be required on the part of the investigator, both to set up the experiment and to recruit participants. It is also possible to re-analyse the data collected during the first run of the experiment, conducted between 29 September 2008 and 5 October 2008.

This report assumes that the reader is familiar with the PHP programming language and basic SQL queries, and is capable of uploading files to a publicly accessible Web server using a mechanism such as FTP.

2 System Requirements

In order to run the experiment code, certain minimum software requirements must be met. The experiment may run under older versions of software, but no guarantees are made as to whether this will be the case, and users are strongly recommended to use the latest versions where possible. The minimum requirements include:

- Web server capable of supporting PHP 5 – Apache 2.2 is recommended.
- PHP 5.2.0¹
- SQLite 2.1.0²

The experiment code should run on any operating system and platform which is capable of meeting the above system requirements. However, all instructions in this report assume a Unix-like operating system such as Linux, FreeBSD or Mac OS X.

3 Data Storage

All of the data collected in the course of the experiment is stored in an SQLite³ database. To ensure maximum backwards compatibility, SQLite version 2 is used.

The majority of configuration options are also stored in the database, so it is possible to re-run the experiment with different Web pages and keywords without changing the underlying PHP code. Furthermore, the HTML output is generated from templates, so the look and feel of the experiment can also be altered without changing the PHP code. It should only be necessary to edit `experiment.php` if changes to the logic of the experiment are required – for example, if users were to be shown the keywords suggested by a different participant, rather than their own.

¹Support for accessing SQLite databases is built in to PHP 5 by default, so no action should be required to activate this functionality unless it has been explicitly disabled.

²The database used for this experiment is created in SQLite 2. Attempting to use a database created in SQLite 3 will result in error messages suggesting a ‘corrupt database’.

³<http://www.sqlite.org/>

The schemas for the various database tables are described in the following sections.

Table 1: Schema for `participants` table

Field Name	Data Type	Additional	Notes
<code>id</code>	INTEGER	PRIMARY KEY	Unique number for identifying individual participants
<code>starttime</code>	INTEGER	NOT NULL	Timestamp for beginning of experiment
<code>endtime</code>	INTEGER		Timestamp for end of experiment (0 if user does not complete)
<code>gender</code>	CHAR(1)	NOT NULL	Set to 'm' or 'f'
<code>age</code>	VARCHAR(10)	NOT NULL	Age range of the user
<code>time_web</code>	VARCHAR(10)	NOT NULL	Time user spends on Web each week
<code>english_native</code>	CHAR(1)	NOT NULL	'y' if user is a native English speaker, 'n' if not

Table 2: Schema for `webpages_config` table

Field Name	Data Type	Additional	Notes
<code>id</code>	INTEGER	PRIMARY KEY	Unique number for identifying individual Web pages
<code>original_url</code>	VARCHAR(255)	NOT NULL	Original URL of the Web page
<code>local_path</code>	VARCHAR(255)		Path to the local copy of the Web page
<code>keywords</code>	VARCHAR(255)	NOT NULL	The keywords generated by our parser

Table 3: Schema for `keyword_ratings` table

Field Name	Data Type	Additional	Notes
<code>participant</code>	INTEGER	NOT NULL	Foreign key on <code>participants.id</code>

starttime	INTEGER	NOT NULL	Timestamp for when user began looking at this page
endtime	INTEGER	NOT NULL	Timestamp for when user finished looking at this page (0 if user quit experiment at this point)
webpage	INTEGER	NOT NULL	Foreign key on <code>webpages_config.id</code>
webpage_keywords_rating	INTEGER	NOT NULL	Rating for the keywords generated by our parser (1-5)
user_keywords_rating	INTEGER	NOT NULL	Rating for the keywords suggested by the user (1-5)
preference	INTEGER	NOT NULL	Keywords preference value (1-5, where 1 is strongly prefer user keywords and 5 is strongly prefer parser keywords)

Table 4: Schema for `webpages_user_keywords` table

Field Name	Data Type	Additional	Notes
participant	INTEGER	NOT NULL	Foreign key on <code>participants.id</code>
starttime	INTEGER	NOT NULL	Timestamp for when user began looking at this page
endtime	INTEGER	NOT NULL	Timestamp for when user finished looking at this page (0 if user quit experiment at this point)
webpage	INTEGER	NOT NULL	Foreign key on <code>webpages_config.id</code>
keywords	VARCHAR(100)	NOT NULL	Keyword list suggested by user for this page

4 Installation

Once all of the system requirements have been fulfilled, installing the code is a simple case of uploading the files to a publicly accessible directory (e.g. `public.html`) via a mechanism such as FTP or the `scp` command. After uploading the files, the correct permissions must be set in order for the experiment to run correctly.

The majority of files within the experiment will only require read access for all users (`chmod 644`) and directories will require read and execute access (`chmod 755`). However, there are some exceptions where files or directories require different permissions:

- `templates.c` – all users should have read, write and execute access (`chmod 777`).
- `data.db` – all users should have read, write and execute access (`chmod 777`).

In addition to the above requirements, the directory which contains all of these files (e.g. `public.html`) needs to allow all users read, write and execute access.

5 Flow Control

The flow of the experiment code follows a simple process:

1. Find out which page the user has just submitted.
2. Process the form data from this page and save it to the database.
3. Decide what the next page to display should be – at most there will be two options to choose from.
4. Fetch any information from the database which is required to display the next page (e.g. Web page configuration details).
5. Display the next page.

There are two ways for the participant to break this process. Before the experiment begins, the user is shown an example page and asked if he understands what is required during the experiment. If the user selects ‘no’ at this point, he will be redirected to the final page and no useful data will be recorded. Alternatively, the user can finish the experiment by closing the browser window. In this case, any information entered on the current page will not be saved, but all data entered on previous pages will be retained. Finally, if the participant ends the experiment before the final page, no end time will be recorded.

6 Common Errors

Some common error messages which may occur when installing the experiment code include:

- ‘Unable to open database file’: Usually indicates that the permissions are not set correctly on the directory which contains the SQLite database.
- ‘Corrupt database’: Indicates that the database is in SQLite 3 format, as opposed to SQLite 2.
- Blank page displaying: Suggests that the permissions on the `templates_c` directory are not set correctly.

7 Summary

Using the instructions contained in this manual, readers should be able to re-run the HuCEL keywords experiment or extend the experiment to capture additional data. The additional files included with this manual also allow readers to re-analyse the data collected in the original run of the experiment.

8 Associated Files

In addition to this manual, the repository also includes the following files:

- `experiment-code.zip`: A compressed file containing all of the code necessary to re-run the experiment, including an SQL file for recreating the database.
- `data.db`: An SQLite database containing all of the data obtained from the experiment.

Within the `experiment-code.zip` archive are the following folders and files:

- `experiment.php`: The PHP code which provides an interface between the user and the database, responsible for controlling the flow of the experiment as described in Section 5.
- `webpages`: Directory containing static copies of all the Web pages used in our run of the experiment, so that all participants will see the same page.
- `templates`: Smarty⁴ templates for all the pages which will be shown to participants in the course of the experiment.
- `index.html`: The first page of the experiment, which participants should be directed to once the code is up and running.
- `style.css`: Stylesheet for the experiment.
- `validation.js`: Client-side validation routines (JavaScript).
- `create.sql`: SQL commands to recreate the database structure and populate it with initial configuration.
- `empty.db`: An SQLite database set up with all the required tables and configuration, but no user data.

⁴Smarty is a template engine for PHP. It can be downloaded from <http://www.smarty.net> and is available as a package for numerous platforms.