

*I*Robot System User's Manual

IRobotSoft.com

Sep 14, 2011

Catalog

CATALOG	1
1. INTRODUCTION TO THE IROBOT SYSTEM	4
1.1. About the IRobot Manual.....	4
1.2. What is IRobot	4
1.3. Why Choosing IRobot.....	5
1.4. Definitions.....	5
1.5. IRobot Installation	6
1.6. The IRobot System Interface	6
1.7. Robot Execution.....	7
2. CREATE NEW ROBOTS	9
2.1. Record Actions Automatically.....	10
2.2. Design Actions Manually.....	14
A) Go to URL	15
B) A Click	15
C) A List of Links	16
D) Take Data	16
E) Take Table	16
F) Submit a Form	16
G) Logon Form	17
H) Open a Frame	17
I) Sent Emails	17
J) A Schedule	17
2.3. Repeat on Next Pages	18
2.4. Open Links in New Windows.....	20
2.5. Create Variables.....	21
A) Create Variables for Extraction Data.....	21
B) Create Variables for Table Data	23
2.6. Save Variables into Databases.....	25
2.7. Draw Data from Databases	28
A) Define New Data Sources.....	29
B) Connect Robot Input with Data Sources.....	30
2.8. Scheduled Run Tasks.....	32

3.	BUG FIXING.....	32
3.1.	<i>Debugging Robot Errors</i>	34
A)	Test a Single Action	34
B)	Test Run Selected Actions.....	34
C)	View the Log File.....	34
3.2.	<i>Repairing Robot Actions</i>	35
A)	Change the Extraction Query by Wizard	35
B)	Check the Action Property	37
4.	MANAGING MULTIPLE ROBOT TASKS	37
4.1.	<i>Calling Robot Tasks</i>	37
4.2.	<i>Organizing Robot Tasks</i>	38
5.	IROBOT REFERENCES	39
5.1.	<i>Event Definitions</i>	39
5.2.	<i>Internal Variables</i>	41
5.3.	<i>Dataset Variables and Functions</i>	42
A)	Dataset Variables.....	42
B)	Dataset Functions.....	43
C)	Use Dataset Variables for Form Submission.....	45
5.4.	<i>Internal Functions</i>	46
D)	String Functions	46
E)	Date Time Functions	46
F)	Floating Number Functions	47
G)	Data Type Functions	47
H)	File Functions.....	48
I)	Interface Functions.....	49
J)	Browser Functions	50
K)	Automation Functions.....	51
L)	Crawling Functions.....	52
M)	Setting Functions	52
5.5.	<i>Detailed Action Properties</i>	53
1)	Go to URL	53
2)	A Click	54
3)	A List of Links	55
4)	Take Data	57
5)	Take Table	57
6)	Submit a Form	58
7)	Logon Form	59
8)	Open a Frame.....	59
9)	Sent Email.....	60
10)	A Schedule	60
6.	FREQUENTLY ASKED QUESTIONS	61
6.1.	<i>Can I create a dummy action as a placeholder?</i>	61
6.2.	<i>The robot skips certain actions during navigation!</i>	61
6.3.	<i>How can I scrape data from pop-up windows?</i>	61

6.4.	<i>Save Variables does not work with MS-Access database, but it works with XML</i>	62
6.5.	<i>How to run robots from command line?</i>	62
7.	ADDITIONAL INFORMATION.....	62
COPYRIGHT		62

1. Introduction to the IRobot System

1.1. About the IRobot Manual

This manual introduces the use of IRobot system for visual Web scraping and Web automation. The manual is intended for both novice user and experienced users. New users are encouraged to watch **video demos** at: <http://irobotsoft.com/help/> before reading this manual.

Readers without any programming experience can still read the first two sections and be able to record robots to automate simple websites. You can also draw data from databases and save data to text files using the visual tools.

After getting more familiar with the IRobot system, you may start encountering problems when designing more complex robots. Please refer Sections 3 and 4 for tools to locate robot debugs and fix errors. A frequent error our users reported is that the robot does not follow Next pages. You can check out Section 2.3 (Page 18, Repeat on Next Pages) for various options to fix this. After reading these sections, readers would be able to design more robust robots for complex Web sites.

The power of IRobot system comes with a full range of advanced functions and utilities that allow you to design even the most complex robot. Section 5 and 6 introduce advanced features like events, datasets, internal functions, etc. for advanced robot design. They are written in a reference style, so readers can quickly locate relevant parts.

1.2. What is IRobot

IRobot (named for *Internet Robot*) is a visual automation tool to create robot agents, or *irobots*, for Web data collection. An irobot agent is able to navigate Web sites, fill in Web forms, extract data, compute and transform data on the fly, and integrate directly with local databases. Using the user-friendly interfaces, you don't need to have programming skills to create irobots; but with some programming skills, you can create more powerful irobots. IRobot is the ultimate Web automation tool you would need to analyze and aggregate data from the Web.

IRobot runs on MicroSoft Windows NT, XP, Vista, and Win 7, and requires MicroSoft Internet Explorer (IE) 6.0 and above for Web automation. Currently IRobot only support IE for Web automation and does not support other Web browsers.

1.3. Why Choosing IRobot

You choose IRobot because:

- You need a visual tool to automate Web form submission and Web data extraction;
- You want to repeatedly collect data from multiple Web sites;
- You want an easy way to save collected data into databases;
- You want to test your Web sites automatically and repeatedly;
- You want to integrate Web data with your in-house databases;
- You want to process Web data across multiple sources.

IRobot is developed entirely in C++ and visual C++, and works extremely efficient. Internally, IRobot embeds an IE Web browser. Running IRobot for Web navigation works exactly like a normal user navigating on the Web via IE. In fact, if you change any IE security/option settings, it will affect IRobot as well.

1.4. Definitions

Robot (or irobot): An encoded script file containing Web navigation and data computation rules. A robot file usually has a suffix “.irb”. Robot files are password protected so that personal information for Web submission cannot be read directly by others.

Robot action: An action completes certain Web function, such as Web navigation, form submission and data extraction. Variables and event rules can be defined with a robot action for data computation.

Internal variable: Internal variables maintained by the IRobot system during execution.

User-defined variable: Variables defined by users in the robot. A user-defined variable has a global scope and exists during the entire robot execution. **Variable name** should start with an English letter, and followed by a number of letters, digits, or underscore (_) without any space or special characters. So `Abc_1` is a valid variable name, but `Abc-1` is not.

Expressions: IRobot follows most standard SQL syntax for **conditions** and **expressions**. A string is enclosed by two single quotation marks ('), such as 'test string'. IRobot does not recognize double quotation marks ("), so do not use it. If the string itself contains single quotation marks, use another single quot to escape each. So 'It's a good day' is a valid string. Within a string, `\r`, `\n`, `\t`, and `\\` stand for line-feed, new-line, TAB, and backslash-char, respectively. More supported functions are defined in section 5.4.

Events: Events are defined with robot actions. Events are usually used for **data computation**.

Tuple: Data extracted from a Web page are always treated as a **table**. (Even a single item is considered a table with one row and one column) Each row of the table data is termed a **tuple**. IRobot accesses each tuple (or row) of the table at a time, and do computations in between via **events**.

Robot task: A sequences of robot actions comprises a robot task. A robot task works like a **function** or a procedure and can be called from other tasks from events.

IRobot system: A visual interface where users can design **irobots**.

IRobot Engine: A core library for robot execution written in C++. The IRobot engine can be embedded in higher-layer applications using ActiveX controls. Refer the [IRobot ActiveX Programmer's Guide](#) for more help.

HTQL: A Web query language used in IRobot for Web data extraction. Refer [The Hyper-Text Query Language \(HTQL\) – A Web Programmer's Guide](#) from <http://htql.net/> for detailed syntax of HTQL. IRobot supports all HTQL syntax in the Guide. IRobot users can practice HTQL following instructions in page 37.

Data source: Any **external database** (text file, XML file, Access DB, MSSQL, MySQL etc.) defined in a robot for batched data processing.

Dataset: An **in-memory database** to store data read from data sources. A dataset is a named object with multiple fields. Each field can be accessed by the dot (‘.’) operation following the dataset name.

1.5. IRobot Installation

IRobot is published in free evaluation versions in a RAR or ZIP package ([irobot-eval.rar](#) or [irobot-eval.zip](#)). For installation, download the latest package, create an empty directory, for example, “C:\Program Files\irobot\”, herein called IROBOT directory, and exact the “irobot.exe” from the RAR or ZIP package into that directory. The package may contain some “.irb” files for demonstration purpose. Extract those “.irb” files to the same IROBOT directory.

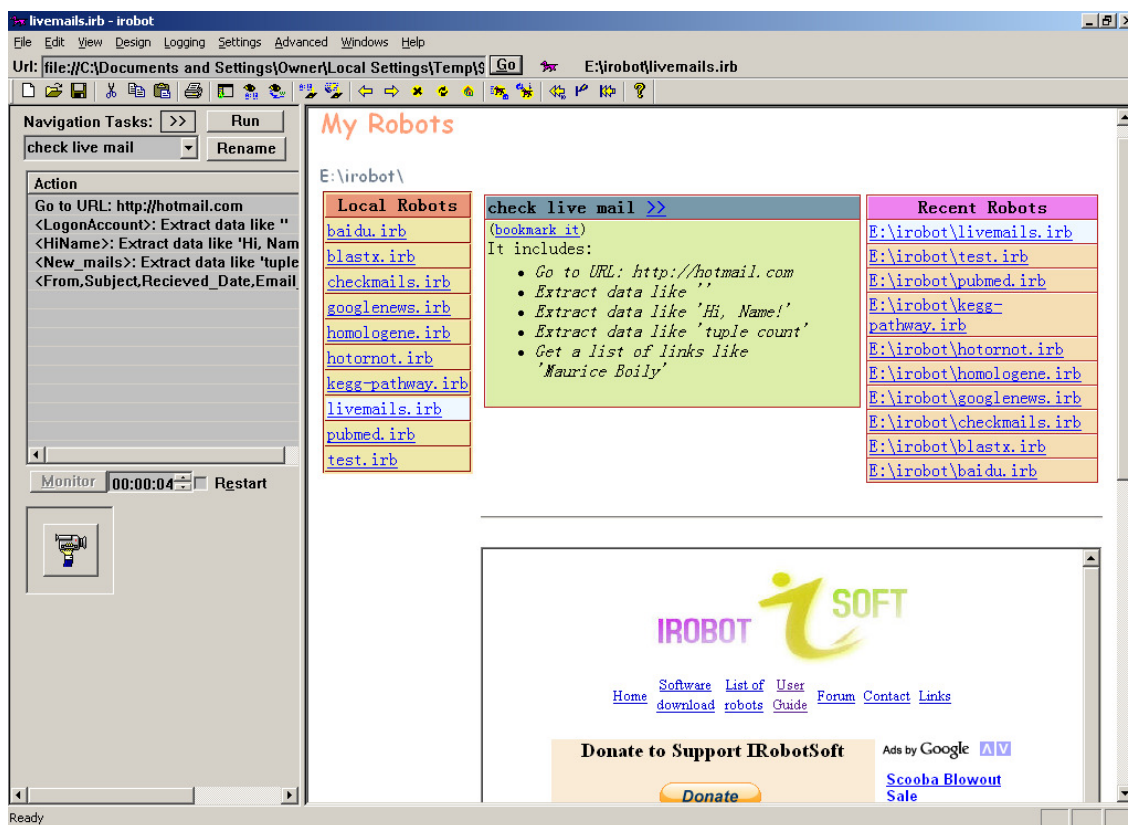
A first run of IRobot will create two additional directories under the IROBOT directory, namely:

IROBOT\data: The default directory for IRobot to save downloaded data.

IROBOT\system: The system directory for IRobot to save temporary data and settings.



1.6. The IRobot System Interface


The IRobot system interface is shown like:



On the left of the interface, it is a control panel. The control panel displays the content of a robot, and includes buttons for running, recording and renaming robot tasks.

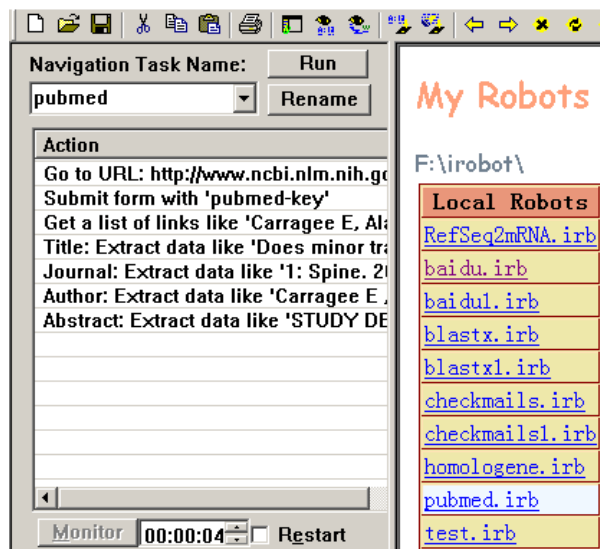
On the right of the interface, there is a browser window showing lists of robots in the IROBOT directory, in the recent execution history, and on the IRobotSoft Web site. You can

always show the list of robots by clicking the  icon on the toolbar: 

There is an icon in the tray-bar shown like . If you minimize the IRobot system, IRobot system will be hidden in the tray-bar. Double click the tray-bar to show the IRobot again.

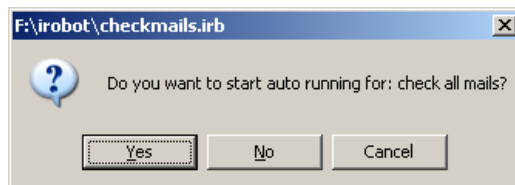
1.7. Robot Execution

Click on a robot file from the right of the IRobot interface to open the robot file, for example, click on the “pubmed.irb” under the “My Robots”. The robot tasks and actions will be shown in the left panel. The following figure shows the list of actions in the “pubmed.irb”.



Choose a correct robot task from the task drop-down list, such as the “pubmed” task, and click the [Run] button to execute the robot.

Some robots are customized to run automatically once it is open. In this case, IRobot system will show a dialog box to ask whether to run a default robot task, for example, if you open the “checkmails.irb”, it will ask you whether to run a default task:



You can click [Yes] to run it or [No] or [Cancel] to stop it.

When a robot is running, the tray-bar icon is animated like

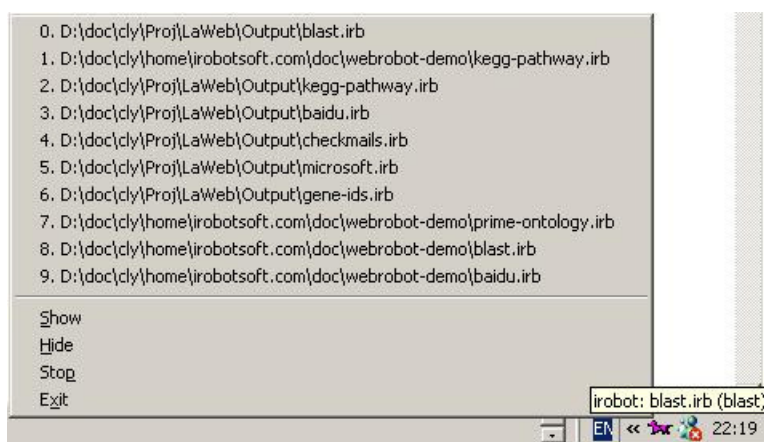


A running robot can be in an active status or in a sleeping status. When a robot is navigating the Web or processing data, it is in the active status. When a robot is running but is idling, it is in a sleeping status. A sleeping robot is still running but is waiting for the next scheduled time to act on certain thing.

A running robot can be stopped in two ways: one by clicking the [Stop] button on the control bar like



and the other by right-clicking on the irobot tray-bar icon and clicking on the stop menu item, like:



2. Create New Robots

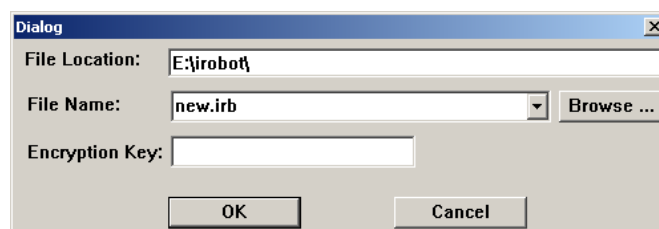
This section will demonstrate how to create a new robot to automatically download abstracts from the [Pubmed literature search](#) and save the abstracts in a local XML file. This is a very common Web scraping scenery and it is recommended for every new starter to read through. After reading this section and recreating the pubmed robot, you will be able to create similar robots for Web scraping or Web automation.

What is the pubmed robot?

The robot will go to the [Pubmed literature search](#) website, ask the user for a keyword. It will then submit the keyword, open a list of links in the page, extract multiple attributes including the abstract from the result page, and follow the “Next” pages continuously. The robot will save the extracted data in an XML file. The final robot can be found from the “[pubmed.irb](#)” on the irobotsoft website: <http://irobotsoft.com/robots/pubmed.irb>. A video demo can be found at: <http://irobotsoft.com/help/>.

Getting start:

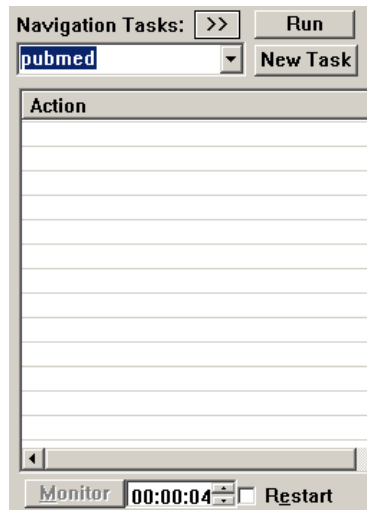
First, we need to create a new robot. Click on the menu **File**→ **New**, and create a new robot file. You may want to name the new robot in a suffix of ‘.irb’, for example, ‘new.irb’, as in the following interface.



You can set an encryption key for the new robot file. If set, the encryption key will be asked whenever the robot file is open. We will leave it empty for now, and click [OK].

The newly created robot has no any Web navigation task. To create a new navigation task, put a task name in the dropdown box under the Navigation Task Name from the left panel,

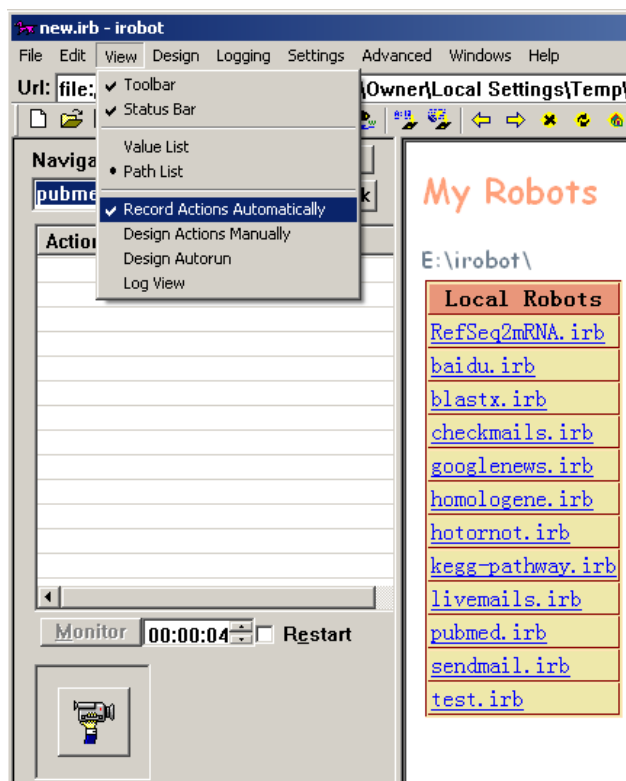
such as “pubmed”, as shown below, and click the [New Task] button to create the new task.



Web actions can be added to the robot task in either of the two ways: 1) by recording; and 2) by step-by-step customization. We will first show how to record a basic “skeleton” of the new robot.

2.1. Record Actions Automatically

To record Web navigation actions, open the robot recorder from menu: View→ Record Actions Automatically and the robot recorder will be shown in the left panel, as in the following figure:



We will press on the recorder button and start recording, like



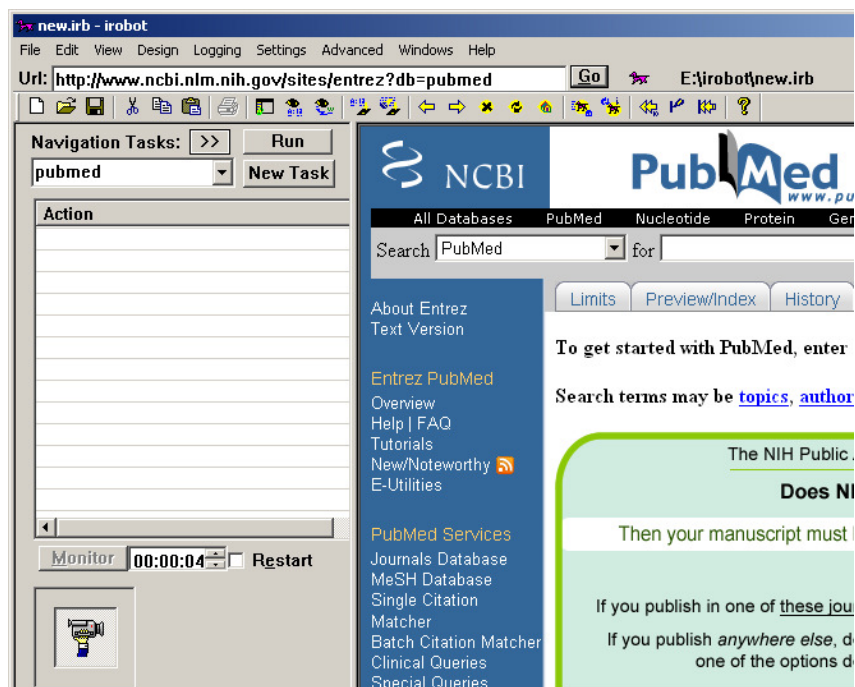
then navigate from the browser as usual:

1) Go to the Pubmed web site: Enter

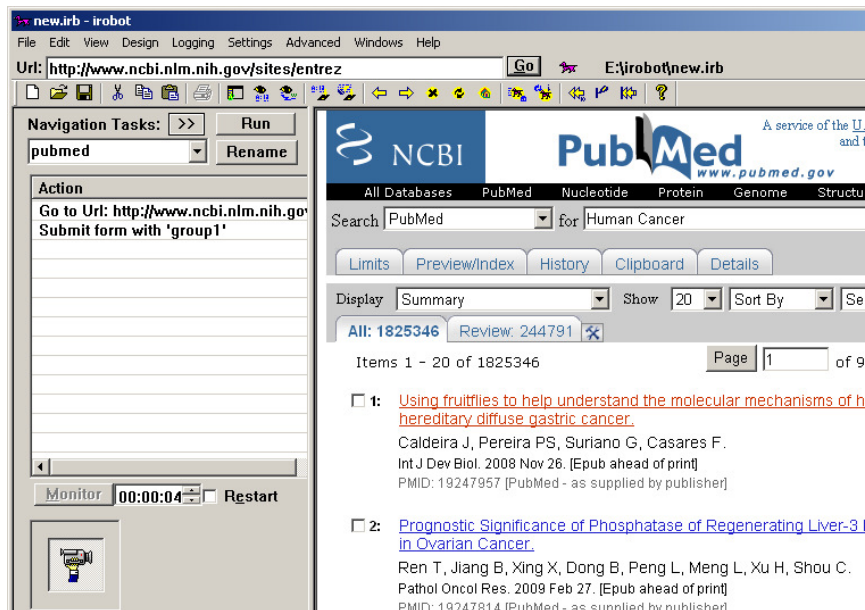
“<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=pubmed>” in the URL bar and click [Go] as in the following figure:



The browser will bring up the Pubmed Web site like the following (the Go to URL action will show up in the next step).

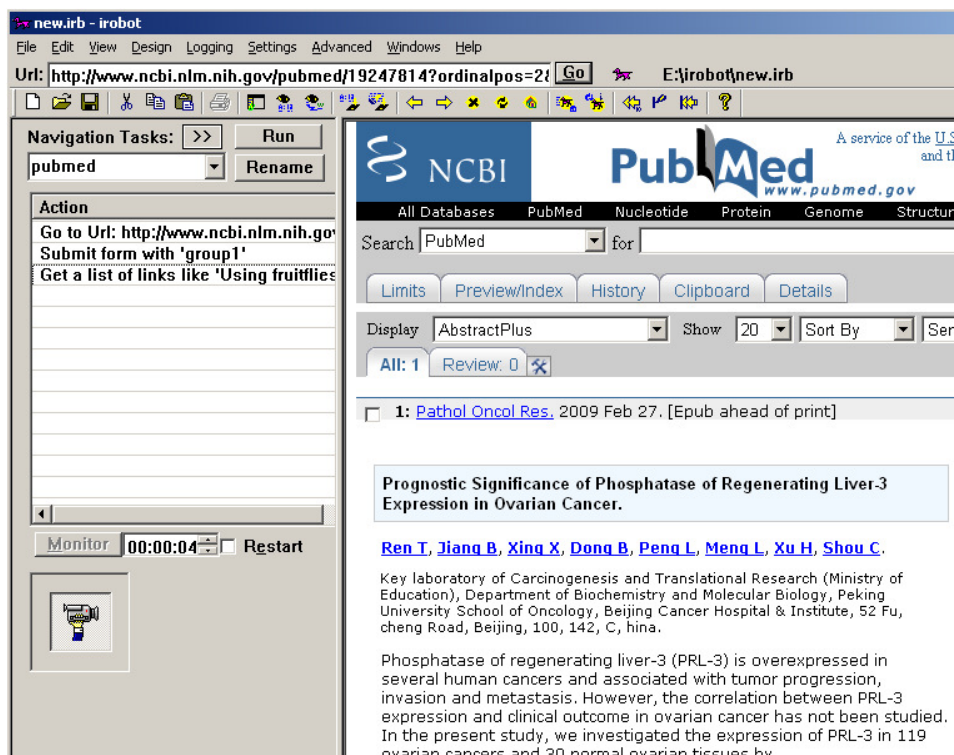


2) Enter a keyword in the Web page, like “Human cancer” and click [Go] in the Web page. It shows:



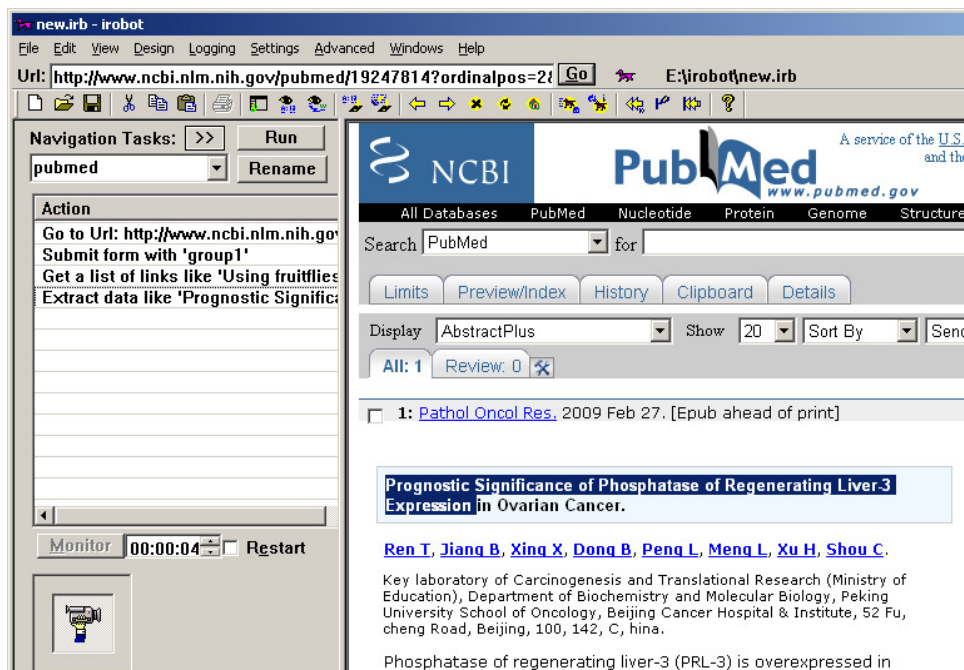
You will notice there is a new task action “Go to Url: ***” followed by another action “Submit form with ‘group2’” shown under the action list. (The Go to URL action always appears when the next action is recognized).

3) Click on a link on the result page, like the second link: “Prognostic Significance of ***”, and it will show a page like:



A new action will also appear under the action list: “Get a list of links like ***” (it may take a while before the recorder figures this action out ;). So, be patient, don’t try to click too fast when recording robots).

4) Mark the title of the abstract until a new action “Extract data like ***” appears under the action list, like:

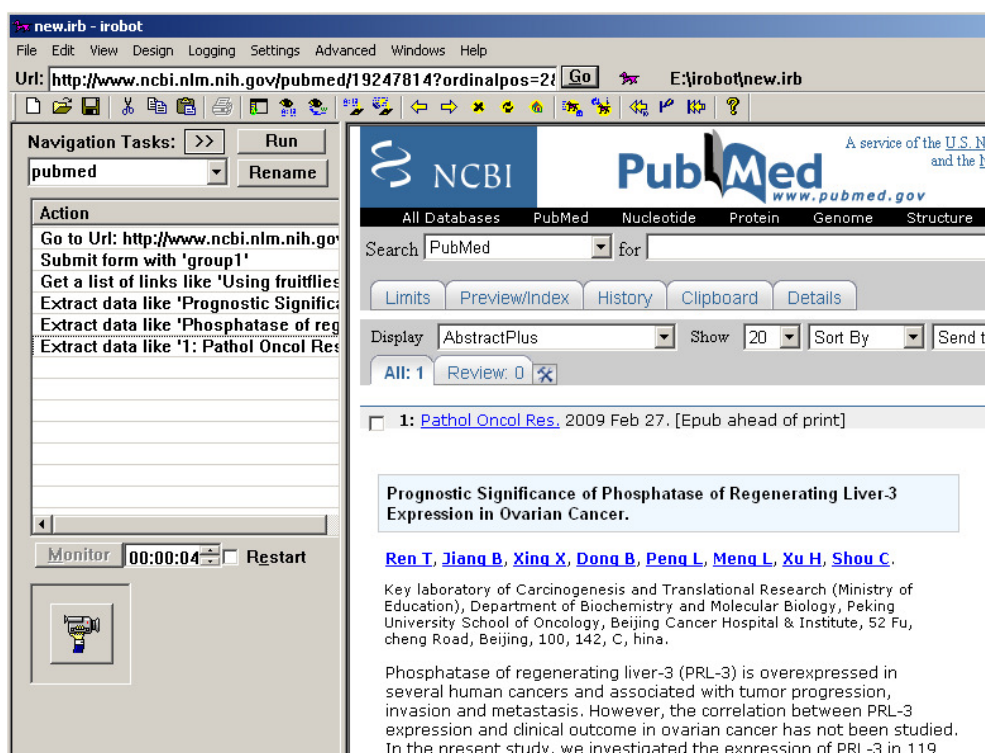


5) Mark the abstract until a new action “Extract data like ***” appears under the action list.

6) Mark the journal information until a new action “Extract data like ***” appears under the action list.

7) Release the recorder button on the left panel to stop recording.

Now we have the skeleton of the “pubmed” robot available as in the following:



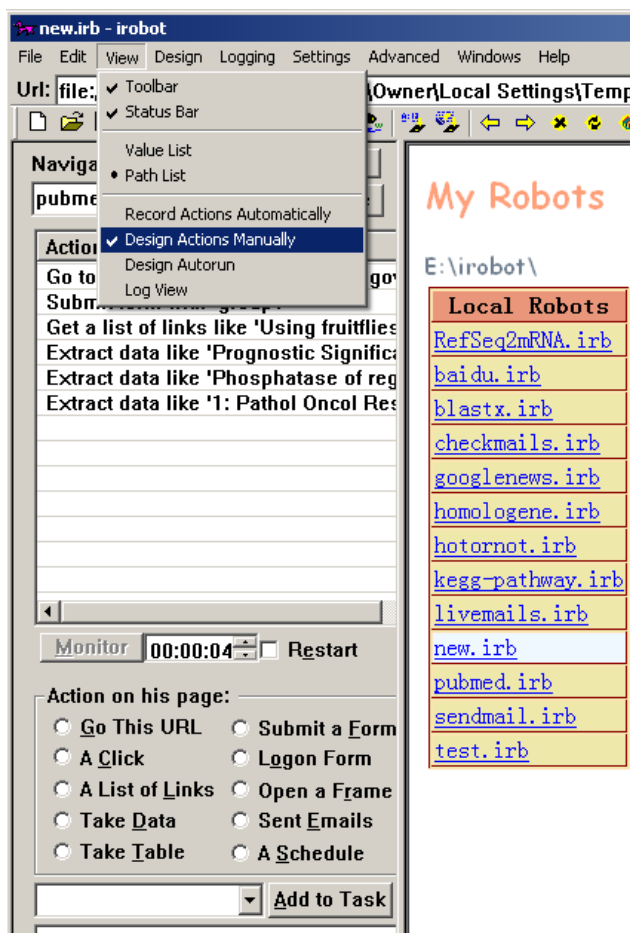
Press the [Run] button on the left panel to test-run this robot. You will see it does everything correctly. You may want to skip the following section “Create Actions Step-by-Step” if the recording can fit your need.

2.2. Design Actions Manually

An alternative way for robot creation is to create actions step by step. This allows advanced users to have a finer control over the creation of robot actions and to better cope with bugs appearing in the robot recording process. (The step-by-step creation is optimized for data recognition instead of recording speed by the recorder. This may solve some errors originated from ill-structured Web pages. Also note that the “Change Extraction ...” menu item is a further option if such error occurs.) We will describe instructions for each type of the action without much detailed explanations. Readers may follow the instructions and

re-create the example demonstrated in the previous robot recording section.

To create actions step-by-step, open the robot designer from menu: View → Design Actions Manually, as:



You can design ten types of robot actions using robot designer: Go to URL, A Click, A List of Links, Take Data, Take Table, Submit a Form, Logon Form, Open a Frame, Sent Emails and A Schedule, as listed in the designer panel. Detailed instruction follows:

A) Go to URL

- i) Navigate to an URL in the browser window;
- ii) Click on the 'Go This URL' radio button in the designer window;
- iii) Click the [Add to Task] button.

B) A Click

- i) Show an example page in the browser window;
- ii) Mark the link to be clicked in the browser window (move your mouse above the

- link, press the left button, drag the link slightly, and release the left button; or simply right-click on the link);
- iii) Click on the 'A Click' radio button in the designer window;
 - iv) Click the [Add to Task] button.

C) A List of Links

- i) Show an example page in the browser window;
- ii) Mark one of the links to be followed in the browser window (move your mouse to above the link, press the left button, drag the link slightly, release the left button; or simply right-click on the link);
- iii) Click on the 'A List of Links' radio button in the designer window;
- iv) Click the [Add to Task] button.

D) Take Data

- i) Show an example page in the browser window;
- ii) Mark the data to extract in the browser window;
- iii) Click on the 'Take Data' radio button in the designer window;
- iv) Click the [Add to Task] button.

E) Take Table

- i) Show an example page in the browser window;
- ii) Mark a row of data to be extracted in the browser window;
- iii) Click on the 'Take Table' radio button in the designer window;
- iv) Click the [Add to Task] button.

F) Submit a Form

- i) Show an example page with forms in the browser window;
- ii) Mark the button that can submit the form in the browser window (move your mouse to above the submit button, press the left button, drag the button slightly, release the left button; or simply right-click on the button);
- iii) Click on the 'Submit a Form' radio button in the designer window;
- iv) The designer window is switched to a form value window;
- v) Enter a new name in the drop-down box under the 'Form Values:' (or you may select an existing name from the drop-down box);
- vi) Click on the [Replace Form] button in the form value window.
- vii) Press [OK] in the form value window;
- viii) The form value window is switched to the designer window;

- ix) Click the [Add to Task] button.

G) Logon Form

- i) Show an example page with logon forms in the browser window;
- ii) Mark the button that can submit the form in the browser window (move your mouse to above the submit button, press the left button, drag the button slightly, release the left button; or simply right-click on the button);
- iii) Click on the 'Logon Form' radio button in the designer window;
- iv) The designer window is switched to a form value window;
- v) Enter a new name in the drop-down box under the 'Form Values:' (or you may select an existing name from the drop-down box);
- vi) Click on the [Replace Form] button in the form value window, as in the above figure.
- vii) Press [OK] in the form value window;
- viii) The form value window is switched to the designer window;
- ix) Click the [Add to Task] button.

H) Open a Frame

- i) Show an example page with multiple frames in the browser window;
- ii) Mark any text in the target frame in the browser window;
- iii) Click on the 'Open a Frame' radio button in the designer window;
- iv) Click [Add to Task] button.

I) Sent Emails

- v) To be discussed.

J) A Schedule

- i) Click on the 'A Schedule' radio button in the designer window;
- ii) Click the [Add to Task] button;
- iii) Double-click the new 'Schedule tasks' action in the action list box;
- iv) Press the [Insert] button on the page in the browser window;
- v) Add new schedule tasks.

Examples are given below:

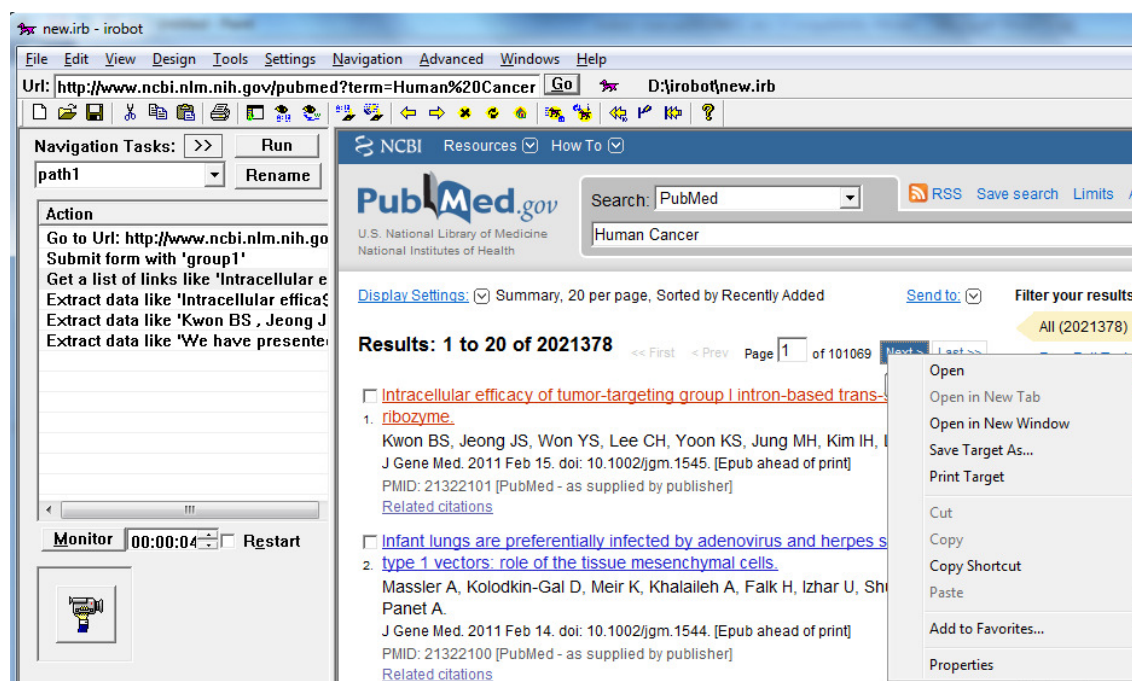
Scheduled Tasks

	Type	Interval	Base	Action	ActionType	Session
<input type="checkbox"/>	DoEvery	1.5	Minute	check yahoo mails	CallPath	yahoomail1
<input type="checkbox"/>	DoEvery	1.5	Minute	check gmails	CallPath	gmail1
<input type="checkbox"/>	DoEvery	1.5	Minute	check hotmails	CallPath	hotmail1
				Modify	Return	Insert
				Delete		

2.3. Repeat on Next Pages

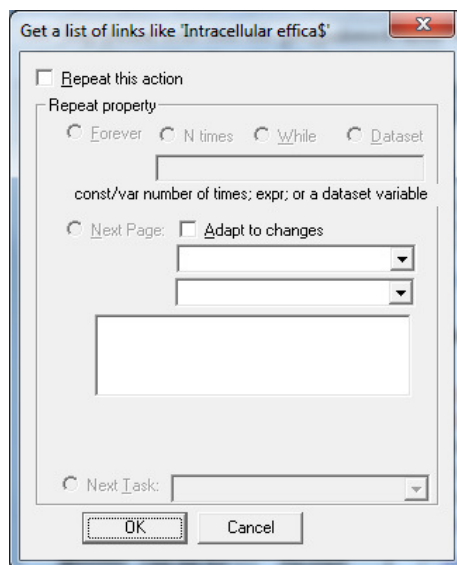
The “new.ibr” robot we have created can open every entry from a Pubmed result page automatically. How if we want to let it follow the Next pages continuously, and from each next page, repeat the actions it does on the first page? We will describe how to do it in the following.

1) Bring up a page including a list of pubmed-search results in the browser panel like:

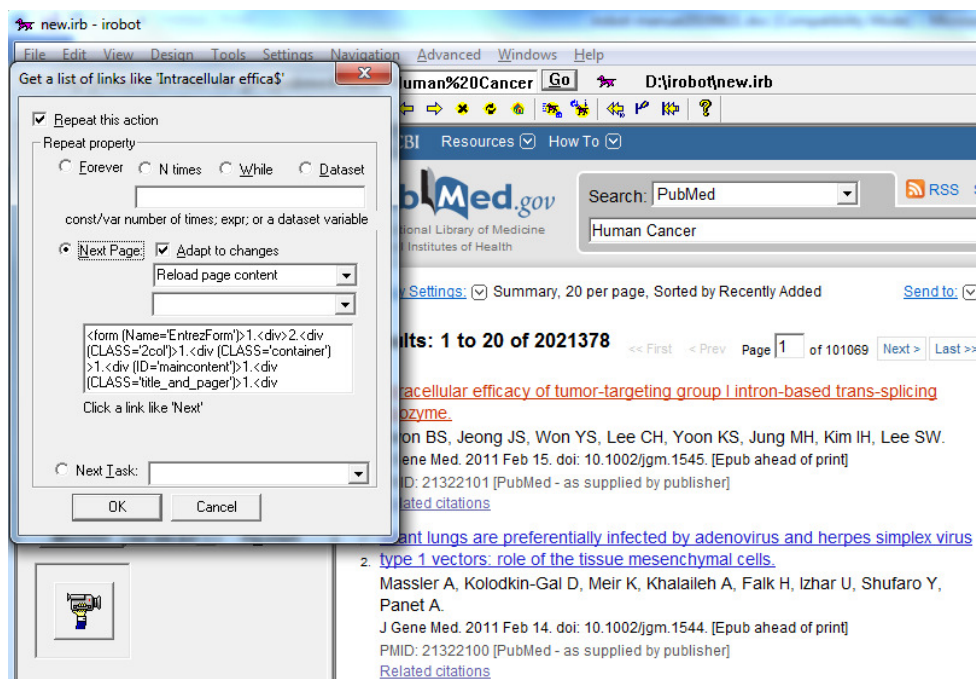


2) Mark the Next link from the Web page (right-click on the link but do not select any option from the popup menu, or move the mouse above the link, press the left button, drag the link slightly, and release the left button) (The highlighted Next link is shown in the above figure).

3) Right click on the action “Get a list of links like ***” from the left panel, and select the “Repeat Property ...” from the popup manual. It shows a dialog:



4) Click on the [Repeat this action] check box, and click on the [Next Page] radio button from popup dialog, as in the following figure. There will be a query shown in the text box, and a message “Click a link like ‘Next’” shown under the text box. Please make sure that the Next link in the Web page has been marked before you clicking on the [Next Page] radio button.



5) We are almost there. For most websites, we are already done and the robot can now click the next pages properly. However, for some JavaScript heavy websites, such as the Pubmed, the robot still refuses to go to the next page properly. We then need some trials and errors with the following options:

Option 1: Below the “Adapt to change,” select “Reload page from URL” instead of the

“Reload page content”. This works if the URL can bring us back to the list-of-links page.

Option 2: select “Do NOT reload page & use URL”. This works if the URL of the Next link can bring us to the next page directly.

Option 3: select “Do NOT reload page & click”. Further, close the repeat property, and open the property page of the “Get a List of Links” action: Check the “To open in new browser”. This usually works if the above two options failed.

Option 4 (rarely used): If option 3 failed and the pages are not opened in the new window, then try this option. Uncheck the “To open in new browser” in the property page, and keep the “Do NOT reload page & click” in the repeat property. Then right-click the “Get a List of Links” action and select “Events...”. Add a new event, Event: *Completed each tuple; Value: GoBack(). This option will tend to fail if any link is not opened correctly. So try NOT to use this option. It is more likely that you need to change the query of the repeat property to make it work; read further.

For Pubmed, the option currently working is Option 3.

Tip 1: To save time on testing the repeat property, you can temporarily add an event to the “Get a List of Links” action, with Event: After each tuple, Condition: Tuple>1, Return: End this page. This will allow the action to open only two links in each page before clicking the next.

Tip 2: Sometimes, the robot goes only to the next 2nd page, and comes back to the 1st page after the 2nd page is completed. This is because the HTQL query in the Repeat property is not robust enough. You then need to manually find a better query that can consistently find the Next link. A typical query that works is:

<a (tx like 'Next%')>

So try this query in the Repeat property if there is a Next link on the page.

5) Click the [OK] button to close the Repeat property dialog. Now we can run the robot, and it will open each link page and follow Next pages automatically.

2.4. Open Links in New Windows

The “new.irb” robot we have created can open every entry from a Pubmed result page automatically. Sometimes, the robot will try to reload the table page before going to the next pages. Then, you can let it run more efficiently with a little tweak—you can open links in a new window, and click the next link without reloading the table page (also refer option 3 in Section 2.3). You need to do the following changes, and also refer the online version of the “pubmed.irb”.

- 1) In the property of the “Get a list of links” action, check the checkbox after the “To open in new browser”;
- 2) Change the “Before each click:” to “do NOT reload & click it”;

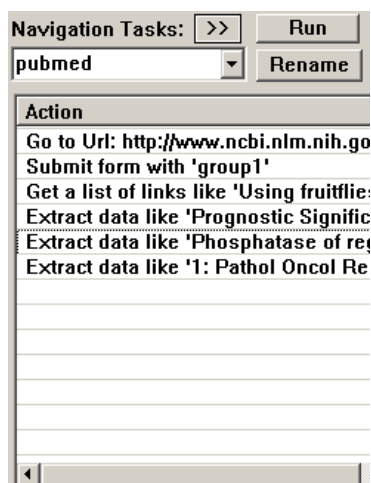
3) In its repeat property, select “do NOT reload page & click” from the dropdown list.

Now you will see that the robot opens each link in a new window, and go to the next pages without reloading the content and become more efficient.

2.5. Create Variables

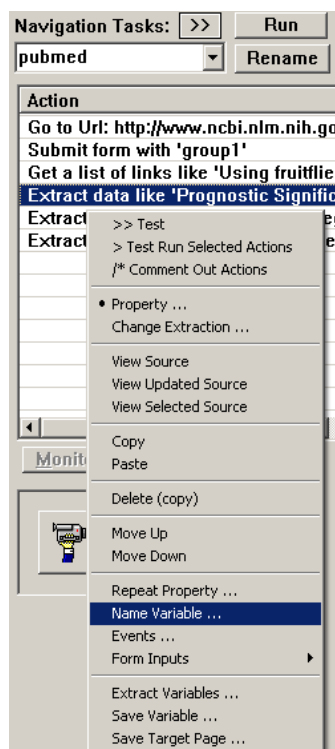
Now we have created the “new.irb” robot to do literature search and Web navigations. Our goal, however, is to extract the literature information and save them into a database. To achieve this goal, we still need to: 1) extract target Web information during robot navigation, 2) define variables to store the extracted data, and 3) save the variable data into a local database. We will illustrate how to define variables to store the extracted abstract author, title and text in this section.

From the “pubmed” task created previously (also shown below) we will define four variables: *AbstractURL*, *Title*, *Abstract*, and *Journal* for the 3rd, 4th, 5th and 6th action, respectively.

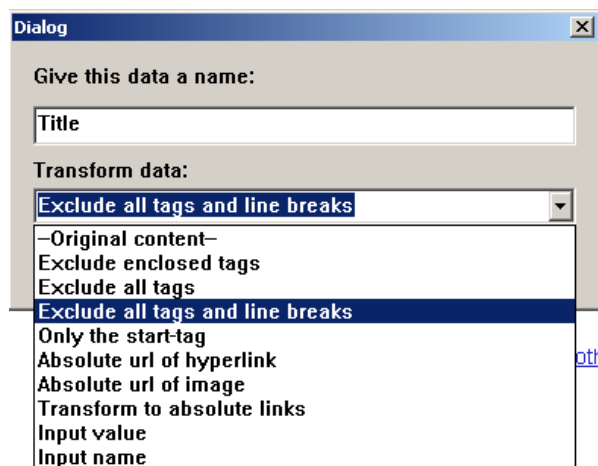


A) Create Variables for Extraction Data

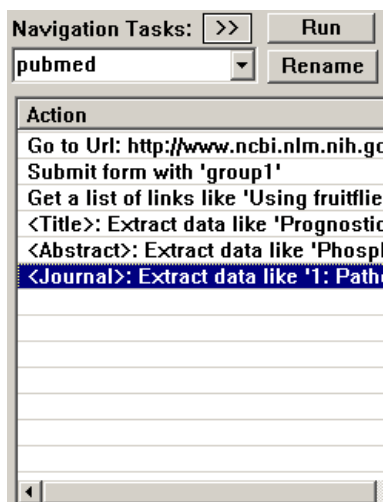
Right-click on the forth action: “Extract data like ‘Analysis of ***’”, and select the “Name Variable ...” menu item, as below.



It shows a dialog box as in the following. Give it a name “Title”, select a data transformation option: “Exclude all tags and line breaks”, and press [OK].

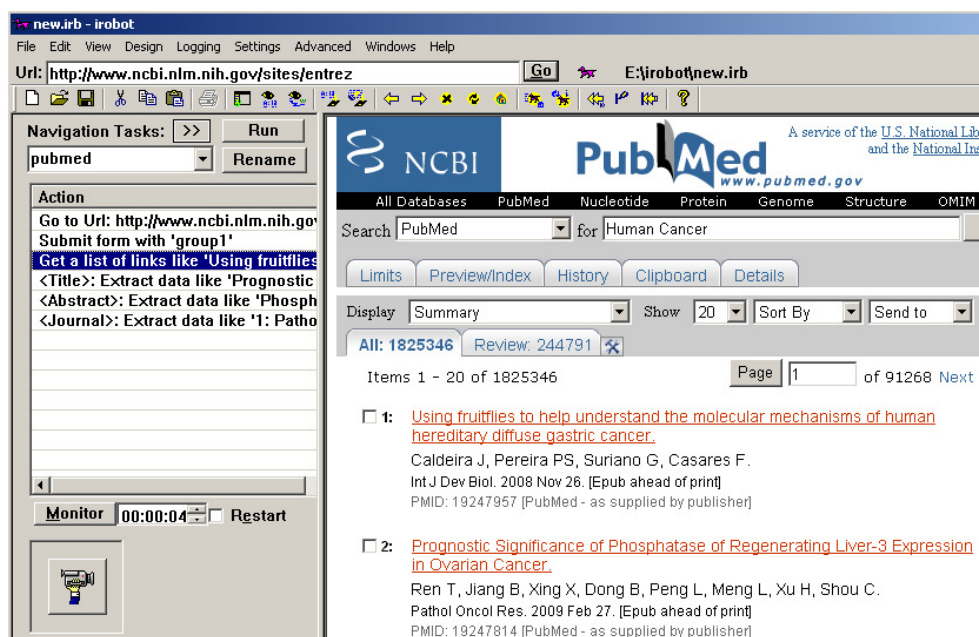


Similarly, define Abstract and Journal for the 5 and 6th action. The “pubmed” task becomes:



B) Create Variables for Table Data

Finally, we are going to define the AbstractURL from the 3rd action, which is more complex. Bring up a result page including a list of entries from the pubmed search in the browser panel like:



Then right-click on the 3rd action “Get a list of links like ***”, and select the “Name Variables ...” menu item. It shows:

Name Table Fields

Column	Name	Transformation
COLUMN1		-Original content-
COLUMN2		-Original content-
COLUMN3		-Original content-
COLUMN4		-Original content-
COLUMN5		-Original content-
COLUMN6		-Original content-
COLUMN7		-Original content-

COLUMN1	COLUMN2	COLUMN3	COLUMN4	COLUMN5	COLUMN6
1:	Using fruitflies to help understand the \$	Caldeira J, Pereira PS, Suriano G, Casar\$	Int J Dev Biol	. 2008 Nov 26. [Epub ahead of print]	PMID: 19247957 [PubMed - as supplied by \$
2:	Prognostic Significance of Phosphatase o\$	Ren T, Jiang B, Xing X, Dong B, Peng L, \$	Pathol Oncol Res	. 2009 Feb 27. [Epub ahead of print]	PMID: 19247814 [PubMed - as supplied by \$

As you can see, the search result page serves as an example page for data columns. Put a name “AbstractURL” after the “COLUMN2”, select a transformation “Absolute url of hyperlink” after it, and click the [Update]. It shows:

Name Table Fields

Column	Name	Transformation
COLUMN1		-Original content-
COLUMN2	AbstractURL	Absolute url of hyperlink
COLUMN3	Authors	Exclude all tags and line breaks
COLUMN4		-Original content-
COLUMN5		-Original content-
COLUMN6		-Original content-
COLUMN7		-Original content-

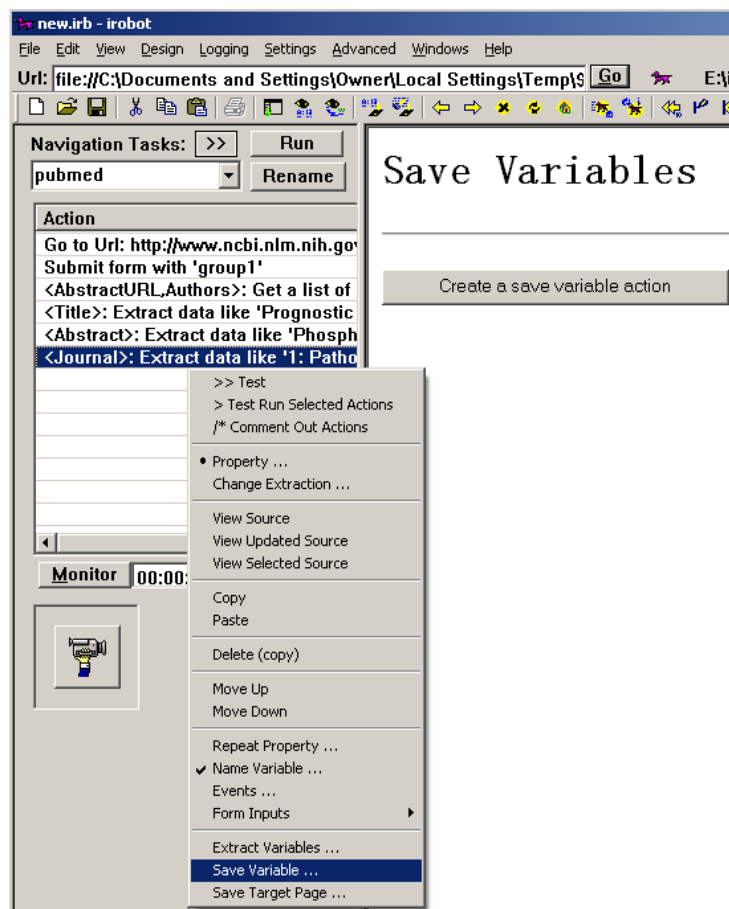
COLUMN1	AbstractURL	Authors	COLUMN4	COLUMN5
1:	http://E:/pubmed/19247957?ordinalpos=1&i\$	Caldeira J, Pereira PS, Suriano G, Casar\$	Int J Dev Biol	. 2008 Nov 26. [Epub ahead of print]
2:	http://E:/pubmed/19247814?ordinalpos=2&i\$	Ren T, Jiang B, Xing X, Dong B, Peng L, \$	Pathol Oncol Res	. 2009 Feb 27. [Epub ahead of print]

Notice that URLs are shown in the sample table in the lower half of the page. Similarly, we may define the *Authors* for “COLUMN3”, with the option “Exclude all tags and line breaks”.

Now click [Return]. We have finished defining the five variables.

2.6. Save Variables into Databases

We are going to save the literature information into an XML file database including the AbstractURL, Title, Journal, and Abstract variables. Right-click on the 6th action “Abstract: ***”, and select the “Save Variable ...” menu item, it shows:



Click on the [Create a save variable action], it shows:

Save Variables

Delete All

When:

Under condition: (Can be a boolean expression.)

Define new variables: [Add](#) [Update](#)

#	Name	Value	Value Type	Transformation of Value	Delete
---	------	-------	------------	-------------------------	--------

[Delete](#)

Save on condition: target type:

Target Datasource: [New](#)

Save Variables: [Add](#) [Update](#) [Move up](#) [Move down](#)

#	TargetField	Expression	Expr Type	Delete
---	-------------	------------	-----------	--------

Sorting by fields:

(example: field1, field2 DESC, field3 NUMBER DESC)

[Add new save](#)

Select the “XML file” from the “target type: [Database]” drop-down list, enter “saved-abstracts.txt” in the “Save to file: []” input box (it will appear once the “XML file” is selected), and click the “Add” following the “Save Variables:”. It shows:

Save Variables

Delete All

When:

Under condition: (Can be a boolean expression.)

Define new variables: [Add](#) [Update](#)

#	Name	Value	Value Type	Transformation of Value	Delete
---	------	-------	------------	-------------------------	--------

[Delete](#)

Save on condition: target type:

Save to file: in name:

Save Variables: [Add](#) [Update](#) [Move up](#) [Move down](#)

#	TargetField	Expression	Expr Type	Delete
<input checked="" type="checkbox"/>	<input type="text"/>	<input type="text"/>	<input type="text" value="Expression"/>	<input type="button" value="Delete"/>

Sorting by fields:

(example: field1, field2 DESC, field3 NUMBER DESC)

[Add new save](#)

Enter “AbstractURL” in both the “TargetField” and the “Expression” boxes, click “Add”;
 Enter “Title” in both the “TargetField” and the “Expression” boxes, click “Add”;
 Enter “Journal” in both the “TargetField” and the “Expression” boxes, click “Add”;
 Enter “Abstract” in both the “TargetField” and the “Expression” boxes, click [Update] button;

It shows:

Save Variables

[Delete All](#)

When: >> After each tuple

Under condition: (Can be a boolean expression.)

Define new variables here: [Add](#) [Modify](#) [Comment](#) [Uncomment](#)

#	Name	Value	Value Type	Transformation of Value	Delete
Delete Save on condition: target type: XML file Save to file: save-abstract.txt Browse... in name: Record tag: Variables TargetField as tag? <input checked="" type="checkbox"/> Save Variables: Add Modify Move up Move down Comment Uncomment					
#	TargetField	Expression	Expr Type	Delete	
<input type="checkbox"/>	AbstractURL	AbstractURL	Expr	Delete	
<input type="checkbox"/>	Title	Title	Expr	Delete	
<input type="checkbox"/>	Journal	Journal	Expr	Delete	
<input type="checkbox"/>	Abstract	Abstract	Expr	Delete	

Sorting by fields:
Not Unique & Sort (example: field1, field2 DESC, field3 NUMBER DESC)

[Add new save](#)

Modify
Move up
Move down
Return

We would like to sort the output by Title and Journal without duplication. So we set the “Title, Journal” after the “Sorting by fields:” and select the [Unique & Keep Old Data & Append File] from the drop down list. It is like this:

Save Variables

Delete All

When: >> After each tuple
Under condition: (Can be a boolean expression.)

Define new variables here: [Add](#) [Modify](#) [Comment](#) [Uncomment](#)

#	Name	Value	Value Type	Transformation of Value	Delete
Delete Save on condition: target type: XML file Save to file: save-abstract.txt Browse... in name: Record tag: Variables TargetField as tag? <input checked="" type="checkbox"/> Save Variables: Add Modify Move up Move down Comment Uncomment					
#	TargetField	Expression	Expr Type	Expr Type	Delete
<input type="checkbox"/>	AbstractURL	AbstractURL	Expr		Delete
<input type="checkbox"/>	Title	Title	Expr		Delete
<input type="checkbox"/>	Journal	Journal	Expr		Delete
<input type="checkbox"/>	Abstract	Abstract	Expr		Delete

Sorting by fields: Title, Journal
[Unique & Keep Old Data & Append File](#) (example: field1, field2 DESC, field3 NUMBER DESC)

[Add new save](#)

Modify Move up Move down Return

Click [Return](#) and we are done.

Now run the robot, and open the “saved-abstracts.txt” file from the IROBOT installation directory. You can see the abstract information is saved correctly and there are no duplicated results.

2.7. Draw Data from Databases

We are going to automate the literature search with multiple keywords. For example, we have three candidate keywords and we want to submit them for Pubmed search in batch:

Human cancer

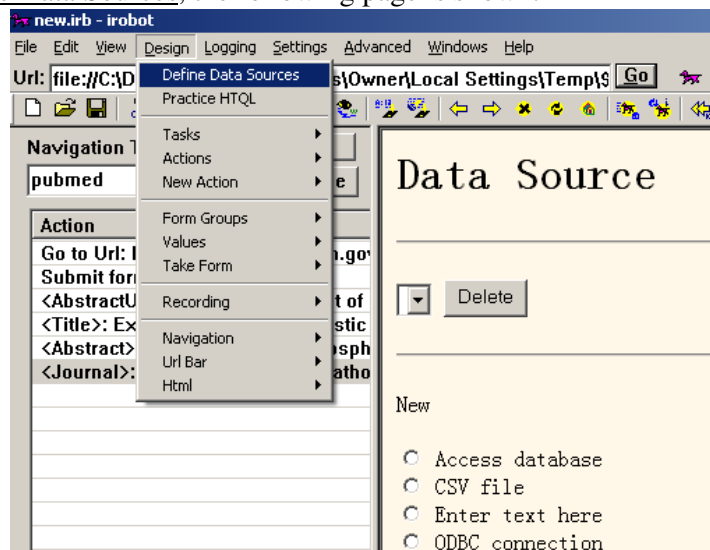
Tumor

Breast cancer

We can do this in two methods: 1) setup a connection between the data source and the form inputs, and let iRobot automatically draw input from the data source; and 2) load the data into a dataset variable, manually repeat on the dataset, and use dataset fields as variables for form inputs. The second method requires some programming skill and will be introduced later in Section 5.3 Dataset Variables and Functions. This section introduces the first method with two steps: **Error! Reference source not found. Error! Reference source not found.**, and **Error! Reference source not found. Error! Reference source not found.** This method is relatively simple and does not require any programming skill.

A) Define New Data Sources

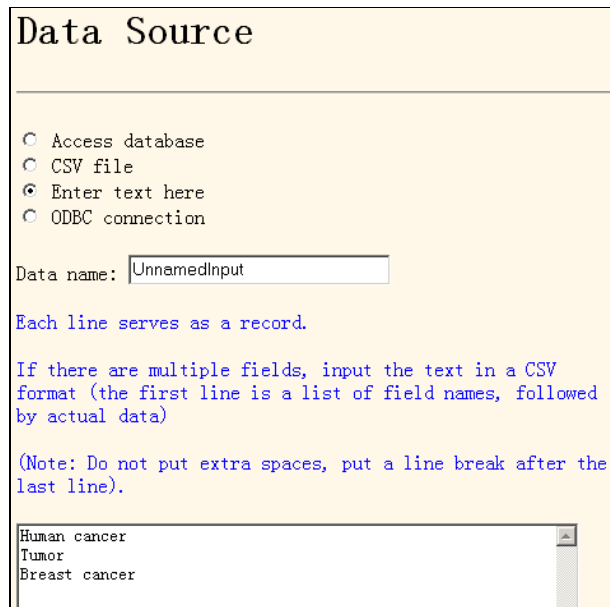
First, we will create a *data source* to include the above three entries. Select the main menu Design → Define Data Sources, the following page is shown.



Select the “Enter text here” radio button, which means that the data will be drawn from some text inputs. If you have keywords in an Access database, or in a CSV file, choose the corresponding radio button instead. After clicking the “Enter text here” radio button, it shows:

 The 'Data Source' dialog box is shown with the 'Enter text here' radio button selected. Below the radio buttons, there is a text input field labeled 'Data name:' containing the text 'UnnamedInput'. Below this, there is a text area with the following instructions: 'Each line serves as a record.', 'If there are multiple fields, input the text in a CSV format (the first line is a list of field names, followed by actual data)', and '(Note: Do not put extra spaces, put a line break after the last line)'. At the bottom of the dialog box is a large text input area for entering the data.

Give the data source a name, default as “UnnamedInput”, and fill the keywords in the text box, each keyword in a line:

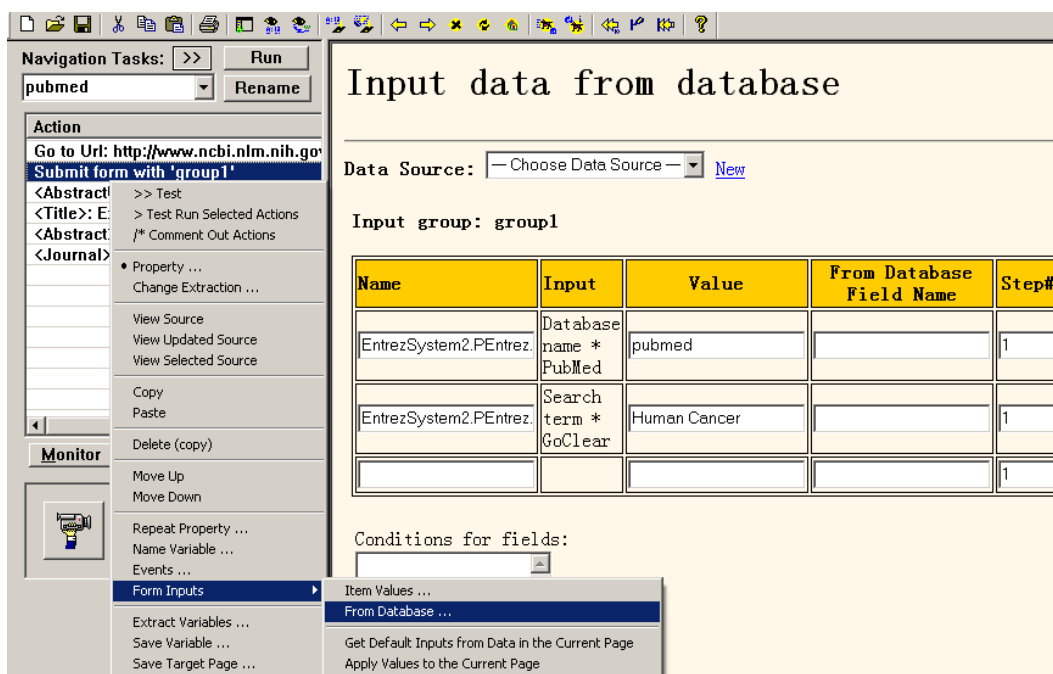


The screenshot shows a window titled "Data Source" with a yellow background. It contains four radio buttons for selecting a data source: "Access database", "CSV file", "Enter text here" (which is selected), and "ODBC connection". Below the radio buttons is a text input field labeled "Data name:" containing the text "UnnamedInput". Underneath the input field, there is instructional text in blue: "Each line serves as a record." followed by "If there are multiple fields, input the text in a CSV format (the first line is a list of field names, followed by actual data)" and a note in blue: "(Note: Do not put extra spaces, put a line break after the last line)." At the bottom of the window is a text area containing the text: "Human cancer", "Tumor", and "Breast cancer" on separate lines.

Then press [Next], a confirmation page is shown, click on the 'Done' link, and we return to the original data source page.

B) Connect Robot Input with Data Sources

Next, we will show how to modify the previously created Pubmed robot to submit searches using values in the "UnnamedInput" data source we have just created. Right-click on the second action "Submit form with group1" and from the popup menu, select "Form Inputs → From Database ...", as in the following:



Input data from database

Data Source: [New](#)

Input group:

Name	Input	Value	From Database Field Name	Step#
EntrezSystem2.PEntrez.name *	PubMed	pubmed		1
EntrezSystem2.PEntrez.Search term *	GoClear	Human Cancer		1
				1

Conditions for fields:

Draw data from data source:

Table: UnnamedInput

#	Field	Type
1	Line	memo

Variables

#	Field	In Task
1	Title	pubmed
2	Abstract	pubmed
3	Journal	pubmed
4	AbstractURL	pubmed
5	Authors	pubmed

Now select the “UnnamedInput” from the dropdown menu following the “Data Source”, as above, and the “Table: UnnamedInput” is shown on the right. We see the only field in the “UnnamedInput” table is the “1 | Line | memo”.

We are going to submit the “Line” field from the database to the “EntrezSystem2.* | Search term * GoClear | Human cancer” input of the form. So simply put the field name “Line” in

the corresponding input box as shown below. (Or you can first click the input box, then click the “Line” field under the Table: UnnamedInput, and the “Line” field will automatically set in the input box):

Input data from database

Data Source: UnnamedInput (Html, TEXT) [New](#)

Input group: group1

Name	Input	Value	From Database Field Name	Step#
EntrezSystem2.PEntrez	Database name *	pubmed		1
EntrezSystem2.PEntrez	Search term *	Human Cancer	Line	1
	GoClear			1

Conditions for fields:

Draw data from data source:

Iterate Through All Tuples

Table: UnnamedInput

#	Field	Type
1	Line	memo

Variables

#	Field	In Task
1	Title	pubmed
2	Abstract	pubmed
3	Journal	pubmed
4	AbstractURL	pubmed
5	Authors	pubmed

Finally, click the [Confirm] button. We are done with the batch definition.

Now re-run the Pubmed robot, and we will see it continuously draw input from our list of keywords and do the search with each keyword.

An alternative method to connect robot input with data sources is to use dataset variables. Refer section 5.3 Dataset Variables and Functions for this.

2.8. Scheduled Run Tasks

We can easily configure the Pubmed robot to run automatically every few hours. First, type in a new task name after the “Task” and click the “New Task” button. Then, open the design view from menu: View → Design Actions Manually. From the design panel, add an action “A Schedule” and double click the action. From the web page, you can add a schedule to call the “pubmed” task every few hours.

3. Bug Fixing

It is well understood that HTML pages are frequently changing and any rules-based system may fail when substantial changes have occurred. Fortunately, IRobot allows you to easily divide and conquer complex navigation tasks into smaller pieces. Here are some tips that

may help you fix any bugs:

- 1) Keep in mind that each action or sequence of actions in IRobot can be tested independently. So first test each action thoroughly and make sure it is performing as expected. Subsection **Error! Reference source not found.** “**Error! Reference source not found.**” shows some useful tools for testing robot actions and pinpointing error actions.
- 2) A typical cause of error is that robot actions are not robust enough, and they no longer work with new Web pages. Subsection **Error! Reference source not found.** “**Error! Reference source not found.**” demonstrates some simple tools to fixing them.
- 3) However, if you have pinpointed the error, and the above tools cannot solve the problem, you may need to manually recompose the Query in the Property of each action, and make them robust across all pages. A typical method is to change their extraction query from multiple pages, compare the different queries it generated automatically, and come up with your best query.

The remaining section guides you through steps for checking robot errors, and tools for fixing the robots.

You may need to have a basic understanding of how IRobot works before fixing bugs. The basic logic of a robot task is a recursive invocation of a set of actions. In other words, actions are not working sequentially, but *recursively*. For example, if there are three actions A1, A2, A3 in a robot task T, then, you may think of task T as a function like the following:

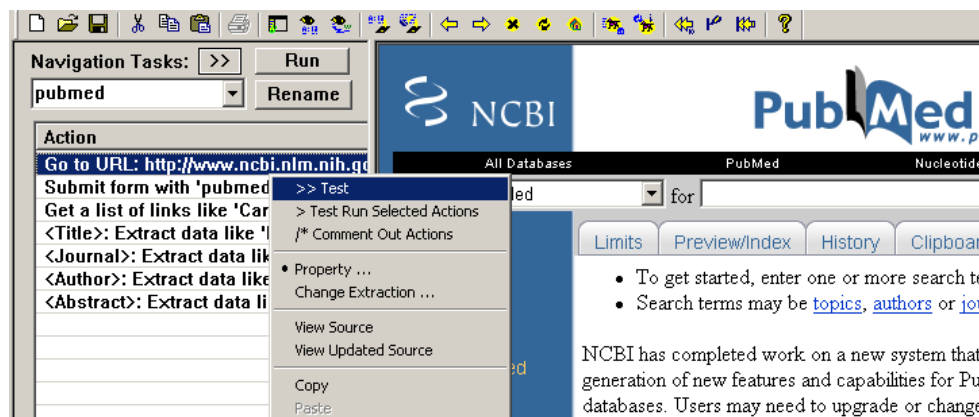
```
Function T()
  For each tuple a1 in A1
    Do some computation on a1
    For each tuple a2 in A2
      Do some computation on a2
      For each tuple a3 in A3
        Do some computation on a3
      End
    End
  End
End
```

Without understanding this recursive logic, you may find some weird behavior from your robots. Note that if no tuple is found from the current action, the rest of the actions will NOT be executed. An exception is for the “Extract Data” action as the current action, where no matter the data is found from the page or not, the next action will be invoked.

3.1. Debugging Robot Errors

A) Test a Single Action

The simplest way to test a robot action is to just right click on the robot action and select the “Test” menu item.

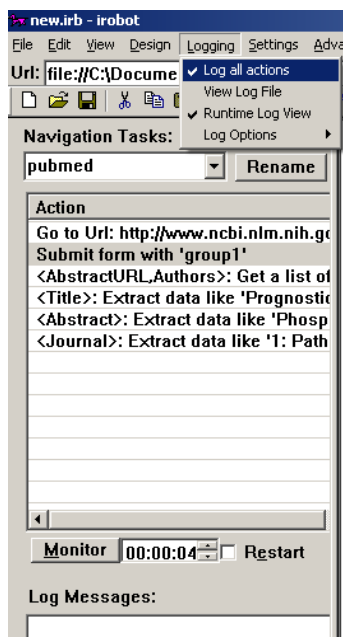


B) Test Run Selected Actions

The “Test” option simply checks if the action is working as expected. It does not evaluate any variables, events, or task calls associated with the action. In order to check these associated actions, you can use the other menu item: “Test Run Selected Actions”. This can be combined with the following action log method to check if variables are assigned correctly, events are executed in the right order, or task calls have been invoked as wanted.

C) View the Log File

The best way to debug a robot is to open its action log. You can log robot actions by opening the main menu → Logging → Log all actions. Then run the robot, and open the “IROBOT\system\action.log” file.



It may give you something like this:

```

0|2007/08/24 19:47:43|===== Start: =====
16|2007/08/24 19:47:43|[pubmed]*1(0)* : Url : http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=PubMed&itool=toolbar
16|2007/08/24 19:47:46|[pubmed]*2(0)* : Form : pubmed-key
16|2007/08/24 19:47:51|[pubmed]*3(0)* : Table : <FORM (Name='EntrezForm')>1.<TABLE (ID='resultview')>1.<TABLE
16|2007/08/24 19:47:54|[pubmed]*4(0)*Title : Item : <FORM (Name='EntrezForm')>1.<TABLE (ID='resultview')>1.<TAB
2|2007/08/24 19:47:54|Title:="<H2>The utility of focused assessment with sonography for trauma as a triage tool i
2|2007/08/24 19:47:54|Title="The utility of focused assessment with sonography for trauma as a triage tool in mul
16|2007/08/24 19:47:54|[pubmed]*5(0)*Journal : Item : <FORM (Name='EntrezForm')>1.<TABLE (ID='resultview')>1.<T
2|2007/08/24 19:47:54|Journal:="<SPAN class=ti><INPUT id=UidCheckBox type=checkbox value=17715308 na
2|2007/08/24 19:47:54|Journal:="1: J Ultrasound Med. 2007 Sep;26(9):1149-56."
  
```

You can then check if actions are taking the right data, and variables are correctly assigned.

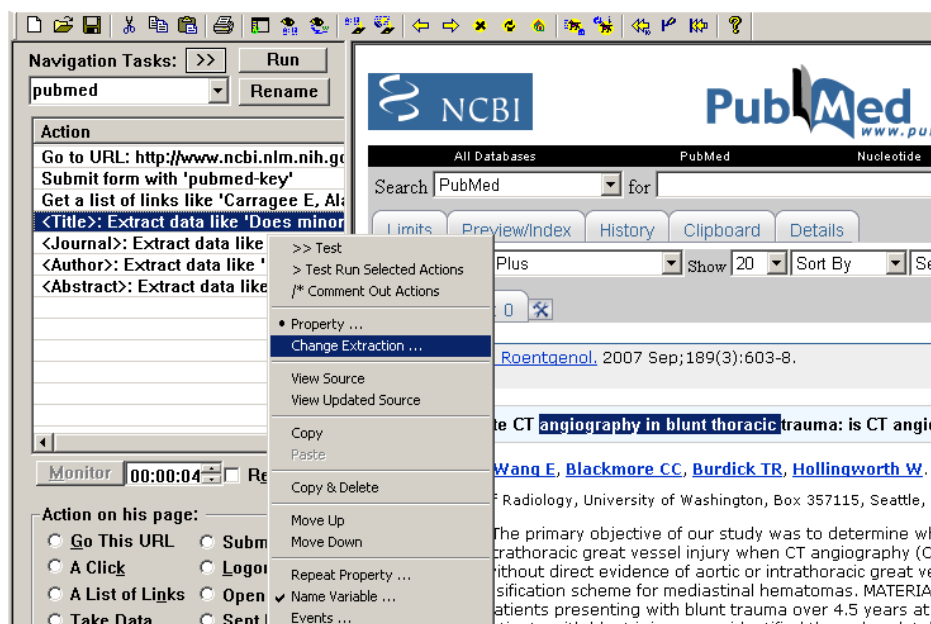
3.2. Repairing Robot Actions

If errors are detected, you may want to repair the robot actions or recreate the action.

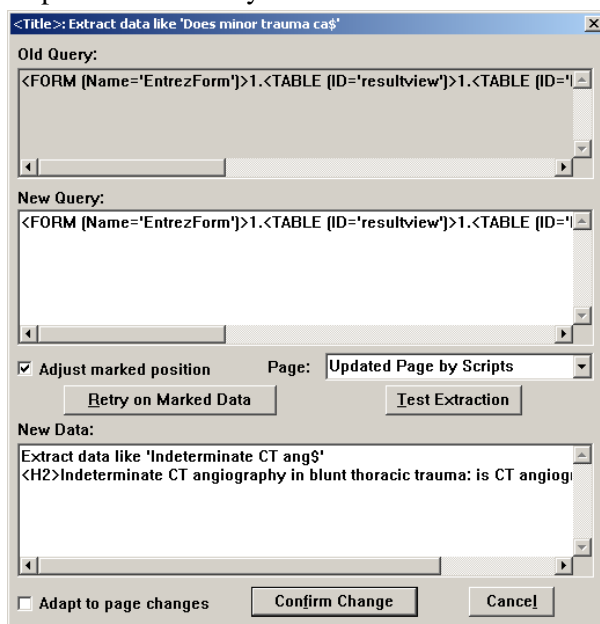
Repairing is desirable if you want to preserve the variables or events that are associated with an action instead of recreating everything from the scratch.

A) Change the Extraction Query by Wizard

The most common error is because the Web page has changed substantially and the action can no longer take the right data from the page. In this case, the first choice is to use the “Change Extraction ...” wizard from the action menu. But before doing so, you first *mark* the correct data in the Web page, and the “Change Extraction ...” wizard will try to relocate data based on your marked position and regenerate a data extraction query in HTQL.

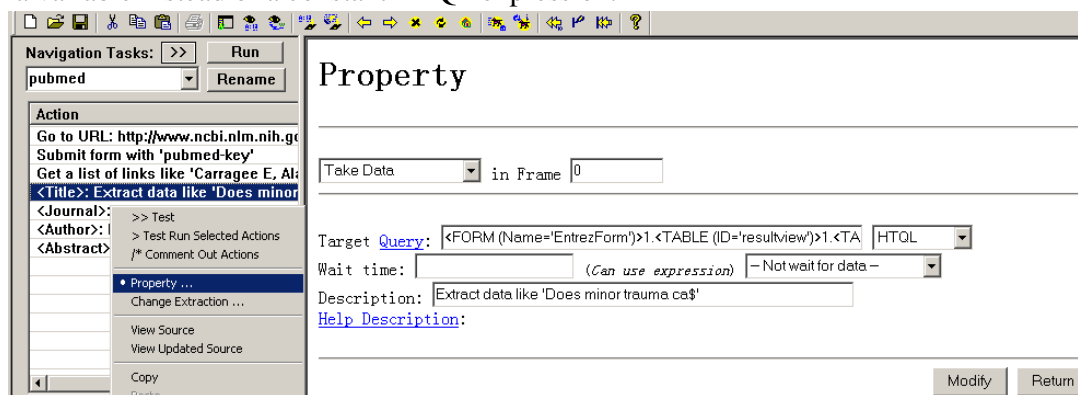


There are multiple options in the “Change Extraction ...” wizard as shown in the following figure. Both old and new queries will be shown in the wizard interface, as well as the data being extracted. If the data is extracted not as expected, you can 1) change slightly the marked positions in the page and retry on the marked data; 2) change some options in the wizard panel and regenerate the query; and 3) if you are familiar with HTQL syntax, you may change the new query expression manually and test it.



B) Check the Action Property

The wizard simplest changes the data query in the action property. You may change the query directly without invoking the wizard. The action property also includes other parameters that are not present in the wizard. Right click on the action and select the “Property ...”, you will able to see the target query in the property page. In this property page, you may also customize other options for the action, for example, you may use queries in a variable instead of a constant HTQL expression.

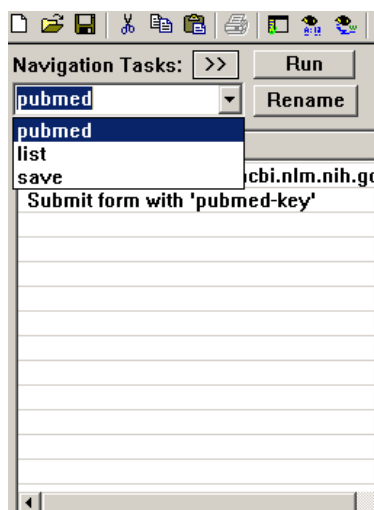


For advanced users, you may click on the “Query” link in the property page to practice the use of different HTQL queries. Remember, however, that you need to navigate to the target Web page before opening the property page, and all the queries will be targeted to the target page.

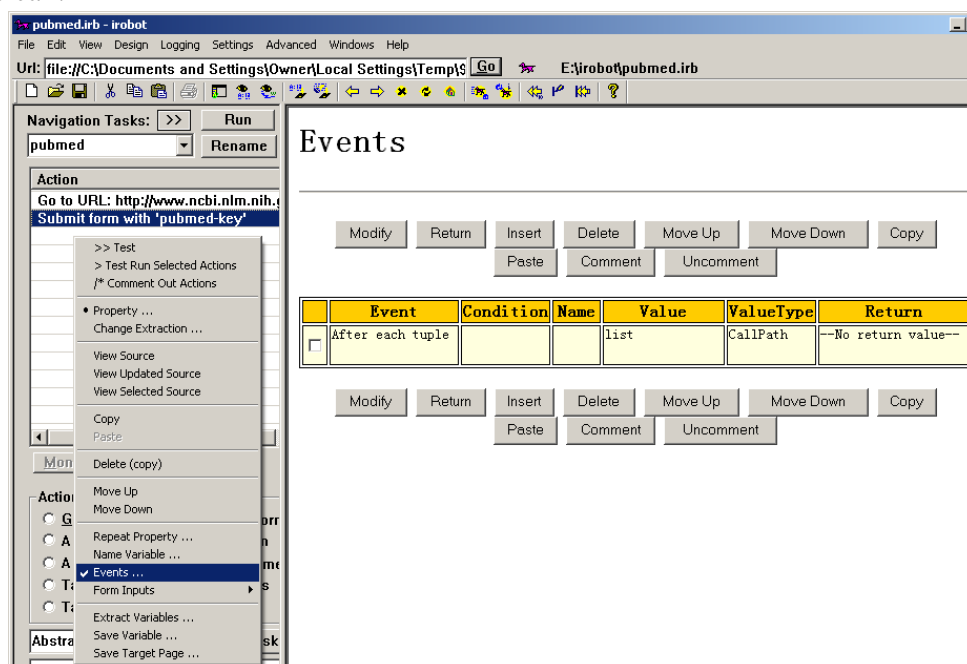
4. Managing Multiple Robot Tasks

4.1. Calling Robot Tasks

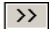
You may break down a complex robot into multiple robot tasks to modularize the robot program. These multiple tasks can then call each other recursively just like making functional calls in a regular program. The following figure shows that the sample robot “pubmed.irb” is divided into three tasks: “pubmed”, “list”, and “save”.

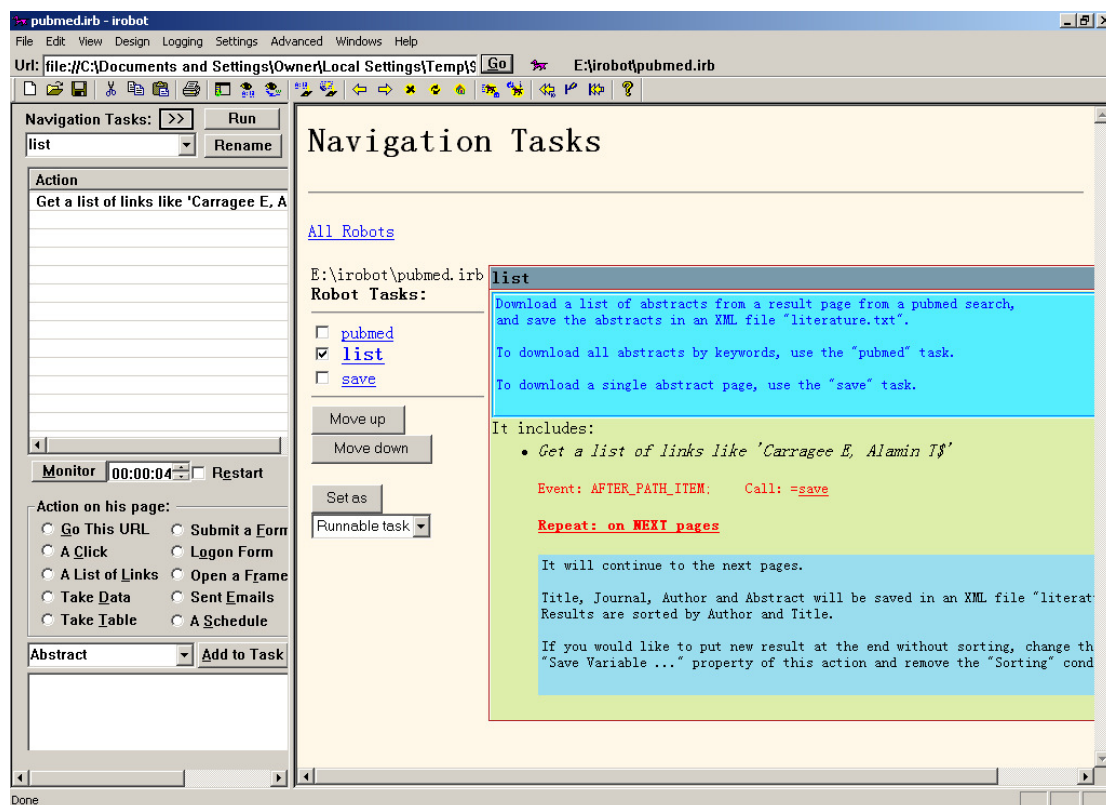


Robot tasks can be called from events. You may create an event by right-clicking on an action, selecting the “Events ...” menu item, and inserting a new event. The following figure illustrate how the “pubmed” task calls the “list” task from an even of the “Submit form with 'pubmed-key'” action. The “CallPath” under the Value Type column specifies that the Value “list” is a task to call.



4.2. Organizing Robot Tasks

You may click on the  after the “Navigation Tasks” to have an overview of the robot tasks. You can select a set of tasks to show a summary of them and move the relative position of robot tasks.



You may also set a task as “Runnable task”, “Disabled task” or “Normal task” here. A runnable task will be listed when you left-click on the tray-bar icon, where you can run the task without showing the visual interface. Runnable tasks will also be shown in the main browser window as quick launches.

5. IRobot References

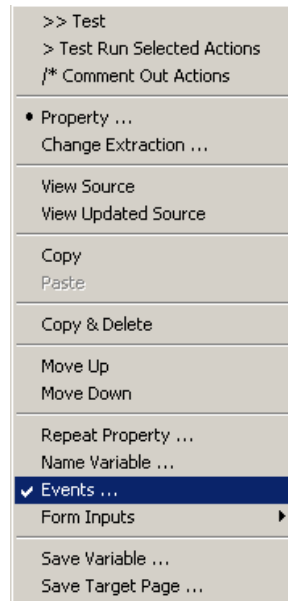
5.1. Event Definitions

Each action fires a number of events. You can use event rules to create variables, evaluate expressions, and test conditions using these events. You are able to create event rules when right-click on any action, and select the “Events...” The list of available events, in the order they are fired from the action, is:

- Before ANY page: Before any page is open;
- Before each page: Before the current page is open;
- Read page: The page is just opened and HTML source is available;
- Before each tuple: Before processing each tuple of the page;
- Read tuple: Tuple variables are read from the page;
- Before action: An action will be taken on the current tuple;
- After action: Action on the current tuple is completed;
- After each tuple: After each tuple is processed;

- i) Completed each tuple: After the subsequent robot actions associated with the current tuple have been completed;
- j) After each page: All tuple on the page has been processed;
- k) After ALL pages: After all repeated actions are completed.

Action events are shown when right-clicking a chosen action from the action list box in the control panel, and choosing the “Events ...” from the pop-up menu, as shown in the following figure.



The event page is shown in the browser window like:

Events

	Event	Condition	Name	Value	ValueType	Return
<input type="checkbox"/>	Before action	ErrorCode>=0	a	ChangeFormInputs ('Please put login information in the browser and press OK.')	Expression	
<input type="checkbox"/>	Before action	ErrorCode>=0 and a=-1			Expression	End this page
<input type="checkbox"/>	After each tuple			Sleep(3)	Expression	--No return value--

Modify Return Insert Delete Move Up Move Down

Each event is defined by:

- a) Event: The event when the associated expressions will be evaluated;
- b) Condition: Evaluate the associated expressions only when the condition is satisfied;
- c) Name: A variable will be created to contain the result of the Value expression;
- d) Value: An expression for evaluation;

- e) **ValueType**: Specify how the Value expression will be evaluated; it can be an expression, a string, a task call, or an HTQL expression against the current Web page (UpdatedPage);
- f) **Return**: Some control flag or error code, for example, to end the execution on the current Web page.

5.2. Internal Variables

In addition to variables defined explicitly by users, there are a number of internally defined variables you can use directly.

TaskName: The current task the robot is running.

Tuple: The index of the current row of data in a page that the robot is acting on. When the robot completed all action, Tuple is reset to 0.

CurrentPage: For *repeat* actions, CurrentPage is the count of repeats on a action.

LastError: The status of the current action

- a) *LastError* <0: the action is in error;
- b) *LastError* =0: success;
- c) *LastError* =1: ignore the current tuple and try the next tuple;
- d) *LastError* =2: end the current action (and its subsequent actions) regardless of the rest tuples;
- e) *LastError* =3: pause navigation on the current tuple;
- f) *LastError* =4: retry on the current tuple.

ActionName: Take the current action name, each name corresponding to an action type.

As a list:

- a) Go to URL: *ActionName*='URL';
- b) A Click: *ActionName*='Click';
- c) A List of Links: *ActionName*='Table';
- d) Take Data: *ActionName*='Item';
- e) Take Table: *ActionName*='Table';
- f) Submit a Form: *ActionName*='Form';
- g) Logon Form: *ActionName*='Logon';
- h) Open a Frame: *ActionName*='Frame';
- i) Sent Email: *ActionName*='Email';
- j) A Schedule: *ActionName*='Schedule'.

SourcePage: The original HTML source of the current page.

SourceUrl: The original URL of the current page.

UpdatedPage: The updated HTML source (by JavaScript) of the current page.

UpdatedUrl: The updated URL (by URL forwarding) of the current page.

TargetPage: The updated HTML source of the page currently in the browser.

TargetUrl: The updated URL of the page currently in the browser.

5.3. Dataset Variables and Functions

A) Dataset Variables

A dataset variable is a set of data organized like a table in a database. You may load a dataset from a database or any string. The dataset can be set in the repeat property to loop over all tuples, or you can manually loop through all tuples. Dataset fields can be accessed like a regular variable.

We will demonstrate this functionality by showing how to read URLs from a text file and repeatedly navigate through all URLs. Suppose we have a text file "test.txt" with data (the separator can also be a TAB, ':', ';' or ','):

```
aaal111|http://msn.com
bbbl222|http://yahoo.com
cccl333|http://google.com
```

We are going to repeat a 'Go to URL' action for all the URLs in this text file. The steps are:

- 1) Create a 'Go This URL' action using any target page;
- 2) Insert a new event for this action: Event=Before any pages; Name=tuple_count; Value=newdataset.loadData('test.txt', '|', 'txt');

This loads data from file 'test.txt' to a new dataset variable 'newdataset', and return the total number of tuples in the tuple_count.

- 3) In the "Repeat Property ..." of this action, check the "Repeat this action" checkbox, select the "Dataset" radio button, and enter "newdataset" below the ratio button;

This allows the robot to repeat on every tuple of the "newdataset".

- 4) In the "Property..." of this action, change the Url: to "newdataset.column3", and select the "Variable" after it.

This allows the robot to use URLs in the dataset for navigation. Dataset fields like "newdataset.column3" can be used anywhere like a regular variable. Default field names are column1, column2, for text files. For CSV, XML, and Access files, actual field names in the database will be used.

Dataset functions will be introduced in the following subsection.

B) Dataset Functions

Member functions for a dataset variable include:

1. loadData(dataset_location, [sql_query], [dataset_type], [htql_query])
2. addData(dataset_location, [sql_query], [dataset_type], [htql_query])

Load/add data from a data source. The dataset will be emptied first in the loadData() function.

The dataset_location can be any text file, CSV file, XML file, Access file, or a predefined data source.

The sql_query is an expression for querying the database.

The dataset_type can be 'txt', 'csv', 'xml', 'access', 'str', 'htql', 'var', or empty:

- *For dataset_types 'txt', 'csv', 'xml', and 'access', the dataset_location is the file location.*
 - Example: a.loadData('test.txt', '', 'txt');
- *For dataset_type 'str', the dataset_location is the actual string, and the htql_query is the field separator, such as a comma ','.*
 - Example: a.loadData('aaa,bbb,ccc', ',', 'str', ',')
 - Dataset a will include: 'aaa', 'bbb', and 'ccc'.
- *For dataset_type 'htql', the dataset_location is the actual string, and the htql_query is an HTQL query.*
 - Example: a.loadData('aaa,bbb,ccc', ',', 'htql', 'I', 'I')
 - Dataset a will include 'aaa', 'bbb', and 'ccc'.
- *For dataset_type 'var', the dataset_location is the name of another dataset variable.*
 - Example: a.loadData('b', '', 'var')
 - Dataset a will copy all tuples from dataset b.
- *For empty dataset_type, the dataset_location is a data-source name, and the sql_query is the SQL query that can be executed directly in the data source. In this case, the SQL query can be delete or update statements.*
 - Example: a.loadData('newdatasource', 'select * from newdatasource')
 - Load all tuples in the newdatasource to dataset a.
 - Example: a.loadData('newdatasource', 'delete from

newdatasource where ID=1')

Delete a tuple in the newdatasource.

The `htql_query` is an expression for querying the source file or source database. If `htql_query` is given, the `sql_query` will be the SQL query to further filter tuples from the dataset (where the "FROM TABLE" clause can be omitted).

2. `nameFields(fields)`

Rename dataset fields.
The 'fields' is a list of comma separated field names.

3. `next()`, `prev()`

Move the dataset to the next/previous tuple. Note that in the above example, the `next()` is called implicitly by the repeat property.

4. `first()`, `last()`

Move the dataset to the first/last tuple.

5. `query(sql_query, tuplejoint, fieldjoint)`

6. `query(sql_query, tuplejoint, fieldjoint)`

Query the dataset with `sql_query`. If `tuplejoint` is not set, only the first result tuple will be returned. If `fieldjoint` is not set, only the first field will be returned. Results are concatenated together as a string using `tuplejoint` to joint tuples, and `fieldjoint` to joint fields.

7. `sql(query)`

Filter the dataset with a query. Tuples not satisfying the condition are deleted from the dataset.

8. `formatHtml()`

Format data in the dataset as an HTML table.

9. `EOF()`

Test if the dataset is ended.

10. `MatchPatterns(string, [separator=',',] [jointor=',',] [casesensitive='false',] [ds_is_pattern='false',] [string_is_pattern='false',] [from_tuple='str/tuple',]`

```
[max_matches=0]);
```

Match patterns in the tuples or in the *string*. Only the first field of the tuple is used. The function searches tuples that match any pattern in the *string*. Multiple patterns in *string* are separated by *separator*. Matching results are jointed by *jointor*. Matching results are taken from tuple if *from_tuple* = 'tuple', or from *string* if *from_tuple* \neq 'tuple'.

C) Use Dataset Variables for Form Submission

Once a dataset variable is defined, its fields can be used like a regular variable. A typical use is for form submission. For example, in the above example, if we would like to submit the `newdataset.column1`, i.e., values of "aaa", "bbb", and "ccc" to google search repeatedly, then we need to design the robot as follows:

1) Create a 'Go This URL' action to go to `http://google.com/`;

2) Insert a new event for this action: Event=Before any pages; Name=tuple_count; Value=newdataset.loadData('test.txt', 'txt');

This loads data from file 'test.txt' to a new dataset variable 'newdataset', and return the total number of tuples in the `tuple_count`.

3) In the "Repeat Property ..." of this action, check the "Repeat this action" checkbox, select the "Repeat" radio button, and enter "newdataset" after the ratio button;

This allows the robot to repeat on every tuple of the "newdataset".

4) Add a submit form action to submit the google form.

5) Double click the form action, you will be in the "Form Values" panel. Find the input corresponding to the search term, which can be identified by the sample text you input to the search box, right-click on it, and select "From Variable...". Then input the "`newdataset.column1`" (without quotation marks) in the popup box.

5.4. Internal Functions

D) String Functions

NewString=substr(String, StartPosition0, [Len])

Get the substring of a string from StartPosition0 of length Len.

Pos0=strfind(String, Substring, [StartPos0])

Find the first occurrence of Substring in String, starting from StartPos0. If no occurrence is found, return -1. If *StartPos0* ≤ 0, search from the first position.

NewString=strcat(String1, String2, ...)

Concatenate multiple strings into a single string.

Len=strlen(String)

Return the length of a string.

NewString=replace(String, Source, Target)

Replace any occurrence of Source in String to Target.

NewString=ltrim(String, TrimmedChars)

Delete from the left of String any occurrence of chars in TrimmedChars.

NewString=rtrim(String, TrimmedChars)

Delete from the right of String any occurrence of chars in TrimmedChars.

FormattedString = format(fmt, ...)

Format a string. Refer *printf* function in C.

Result = htql(TextSource, HtqlQuery)

Evaluate an HTQL expression against data in TextSource. Only the first result is returned.

Result = get_emails(Text, fmt)

Find emails from Text. If *fmt* is set, Results are formatted in *fmt*; otherwise, Results are formatted as: “name” <email>:. In *fmt*, %m will be replaced by email, %n will be replaced by name, multiple emails will be concatenated after formatted.

E) Date Time Functions

TimeInt=time()

Return the current time in integer.

DateVar=date(), sysdate()

Return the current date.

NewNumber=to_number(String)

Convert a string into a number.

StrDate=to_char(IntDate, DateFormat)

Return a string representation of a date in a certain format.

IntDate=to_date(StrDate, DateFormat)

Return an integer date from a date string in a certain format.

F) Floating Number Functions

RandNum=rand([RandMax])

Return a random number less than RandMax; If RandMax>1, RandNum is integer, if RandMax=1, RandNum is a floating number.

IntNum=round(FloatNum)

Return the rounded integer number of a floating number.

IntNum=floor(FloatNum)

Return the truncated integer number of a floating number.

IntNum=ceil(FloatNum)

Return the next integer number of a floating number.

FloatingNum=get_number(text, index)

Find a *index*'th number from the *text*.

G) Data Type Functions

Test=ischar(var)

Test if *var* starts with an alphabet character, except any leading spaces.

Test=isnumber(var)

Test if *var* starts with a digit of '0'-'9', except any leading spaces.

Test=isblank(var)

Test if *var* is a composed of blank characters.

Test= isphonenumnumber(var)

Test if *var* is a valid phone number, i.e., as a combination of '0'-'9', '+', '-', '(', ')', and '.', and has a length of 6-20 characters.

H) File Functions

SaveUrlFile(URL, FileName)

Download a file at *URL* and save to a file *FileName*.

Data=UrlData(URL)

Get data from *URL*. URL can be a local or http file.

Data=ReadFile(FileName)

Load data from local *FileName*.

SaveFile(FileName, DataString)

Save *DataString* to a file of *FileName*.

AppendFile(FileName, DataString)

Append *DataString* to file *FileName*.

DeleteFile(FileName)

Delete file *FileName*.

CapturePage(ImageFilename)

Capture the current web page and save it as an image. Only BMP file format is supported now.

Name = GetFileName(FileName)

Get the file name of *FileName*.

Path = GetFilePath(FileName)

Get the path of *FileName*.

MkDir(DirName)

Create a directory named *DirName*.

UniqueFileName = GetUniqueFileName(Dir, FileName, [SubDir])

Get a unique file name in *Dir*, or *Dir/SubDir/*, with the base name given in *FileName*. If *Dir* or *Dir/SubDir/* does not exist, the directory will be created by this function. The returned *UniqueFileName* does not include the *Dir* or *SubDir* prefix.

List = DirFiles(Match, 'FILE | DIR | RECUR | PATH | NAME | PATHNAME | WRTIME | CRTIME | SIZE')

Get file information. *Match* is a filename, may include wildcards * and ?.
 'FILE|DIR': search filename or dir name or both.
 'RECUR': search subfolders.
 'NAME|PATH|PATHNAME|SIZE|WRTIME|CRTIME': return the list of information for each file, separated by the same separator (e.g., '|' here) each line.

SendEmail(email_data, to_email)
 Send *email_data* to a recipient email *to_email*.

I) Interface Functions

NewValue=InputText(Description, DefaultValue)
 Prompt user for input.

Count=InputVariables(Title, Name1, Value1, Desc1, Size1, Name2, Value2, Desc2, Size2 ...)
 Prompt user for setting of multiple variables at the same time.
Title: title of the input dialog;
Name: variable name;
Value: default variable value;
Desc: description of the variable;
Size: size of the input box if it is a number, or formatted selection options, e.g.:
 |SELECT|option1:value1|option2:value2
 |RADIO|option1:value1|option2:value2
 |CHECKBOX|option1:value1|option2:value2

Sleep(Seconds1, Seconds2)
 Sleep for a certain amount of time in seconds; If *Seconds1*<0, then sleep for a random amount of time less than the absolute value of *Seconds1*; If both *Seconds1* and *Seconds2* <0, then sleep for a random amount of time in between the absolute values of *Seconds1* and *Seconds2*.

CallPath(TaskName, IrbFile, Passwd, Session)
 Call *TaskName* task in robot file *IrbFile* with password *Passwd*. If *Session* is specified, the new task will run in a separate window named *Session*.

CallParallel(TaskName)
 Call *TaskName* task in a parallel thread. The parallel thread only uses HTTP socket browser, and will not execute JavaScripts.

ShowMessage (Message, WaitSeconds)
 Show *Message*. If *WaitSeconds* is given, close the message after *WaitSeconds* time.

PathAlertMessage (Message, LinkTask)

Add an alert *Message* in the taskbar icon, and play a sound. If *LinkTask* is a task name, the alert *Message* is linked to the *LinkTask* when clicked, otherwise, it works as a message ID. Typical use: *PathAlertMessage (Message, time(0))*.

J) Browser Functions

GetCookie(CookieName)

Get the cookie associated with a *CookieName*.

SetCookie(CookieName, CookieValue)

Set a cookie in the current page.

ClearCookies([Url])

Delete cookies associated with *Url*. If *Url* is not given, delete the cookies associated with the current URL and URLs from the same domain. If *Url*=='*', then delete all cookies in the system.

SetProxyServer([IP], [Port], [SessionAgentName])

Set the proxy server with *IP* and *Port*. Remove proxy setting if no parameters are given. If *SessionAgentName* is not provided, the proxy setting will affect system-wise user settings. Otherwise, it will serve as the agent name of the current session, and the proxy setting will affect this session only. Use, for example, "Microsoft Internet Explorer" for the *SessionAgentName*.

SetAttribute(Html, AttributeName, AttributeValue)

Modify HTML attribute in the current page, *Html* specifies the tag for modification. *AttributeName* can be any attribute in the tag, or 'innerText', 'innerHTML', 'outerText', and 'outerHtml'.

SetInput(InputName, InputValue)

Set input data in the current page.

MouseMove(ItemHtml, [offset_x, offset_y])

Move mouse over the item specified by *ItemHtml*, with an offset of *offset_x*, *offset_y*.

MouseClick(ItemHtml)

Simulate a mouse click at *ItemHtml*, or at the current mouse position.

MouseRightClick (ItemHtml)

Simulate a mouse right click at *ItemHtml*, or at the current mouse position.

Refresh()

Refresh the current page.

GoForward()

Navigate forward.

GoBack()

Navigate backward.

ReloadPage('URL/LOCAL/');

Reload page.

ClickItem(item_htql, MaxWaitSec)

Click a page item specified by *item_htql*.

K) Automation Functions

Titles = GetWindows([Title])

Get application window titles. If *Title* is given, it finds the best matched title from all application titles. Otherwise, it gets all application window titles and separates them with a new line.

RunApp(CommandLine) or

StartApp(CommandLine)

Run application using *CommandLine*.

SendAppText(ApplicationTitle, TextString)

Find the application with title *ApplicationTitle*, and send text *TextString* to it.

SendAppKeys(ApplicationTitle, KeyString)

Find the application with title *ApplicationTitle*, and send keys *KeyString* to it. Refer *SendKeys* in .NET.

SendAppKeysWait(TimeSec, KeyString, WindowTitle, ChildWinTitle, Delay)

Find the application with window title *WindowTitle*, and send keys *KeyString* to its child window whose title is *ChildWinTitle*. Wait for *TimeSec* in the background if the application or the child window cannot be found. If *TimeSec* < 0, continue robot actions while waiting; otherwise, wait until the window is found. After the application is found, wait for *Delay* in seconds before sending the *KeyString*.

CloseApp(ApplicationTitle)

Close an application.

Copy(text)

Copy *text* to the clipboard.

Text = Paste()

Get *Text* from the clipboard.

L) Crawling Functions

SetDownloadParameter(threads, max_queue)

Set the number of parallel threads and the size of waiting queue.

SetDownloadCompletionTask (Id, CallTask)

Call *CallTask* after crawler *Id* completed.

Id = SearchEmail (StartUrl, EmailFile_CallTask, MatchPatterns, ExcludePatterns, MaxDepth)

Search emails from *StartUrl* and its subdirectories. If *EmailFile_CallTask* is a task name, call the task on each crawled Web page, otherwise, save emails to the file named *EmailFile_CallTask*. If *MatchPatterns* is given, only URLs including the pattern will be crawled; if *ExcludePatterns* is given, any URLs including the pattern will be excluded; if *MaxDepth* is given, the crawler will only go *MaxDepth* steps from *StartUrl*. This function only works in the paid version.

Id = SearchSiteLinks (SourceSites, TargetSites, EmailFile_CallTask, MatchPatterns, ExcludePatterns, MaxDepth)

Search links from *SourceSites* to *TargetSites*. Urls in *SourceSites* to *TargetSites* are separated by new lines (\n). See *SearchEmail()* function for other parameters.

Id = CrawlWebsite (SourceSites, CallTask, MatchPatterns, ExcludePatterns, MaxDepth)

Search links in *SourceSites*, and call *CallTask* on each crawled page. Urls in *SourceSites* are separated by new lines (\n). See *SearchEmail()* function for other parameters.

Id = DownloadPage (SourceSite, SaveDir, MatchPatterns, ExcludePatterns, MaxDepth)

Crawl *SourceSite*, and save all files in *SaveDir*. See *SearchEmail()* function for other parameters.

M) Setting Functions

SetProgram ('Window', 'Show/Hide/Maximize/Minimize')

Show IRobot program.

SetProgram ('Program', 'Exit')

Quit IRobot program.

SetProgram ('ToolPanel', 'Show/Hide/Toggle')

Toggle the tool panel.

SetSettings ('Speed', 'Fast/SuperFast/Slow/VerySlow/NonStop')

Set IRobot navigation speed.

SetSettings ('Popup', 'block/nonblock/InBrowser')

Set IRobot popup options.

SetSettings ('PopupSize', 'Large/Middle/Small')

Set IRobot popup window options.

SetSettings ('Silent', 'Yes/No/True/False')

Set IRobot silent browsing options.

5.5. Detailed Action Properties

1) Go to URL

The property page looks like:

Property

Go to Url in Frame (in the updated page)

Url:

Description:

[Help Description:](#)

* [Events](#) *

The Url can be a String like 'http://mail.yahoo.com/' in the example. The Url can also be taken from a Variable or an Expression, such as 'http://mail' + '.yahoo.com'.

Make sure to click the [Modify] button after any change.

The "Description" will be shown in the action list when showing this action. You can use the Help Description to add memo or notes for this action.

2) A Click

The property page looks like:

Property

A Click in Frame (in the updated page)

Target **Query**:

To open in new browser: ☐

Wait navigation? Before each click:

Description:

[Help Description:](#)

* [Events](#) *

“A Click” is similar to “A Link”, which can be chosen from the drop-down list. “A Click” is more general because it can also click on a JavaScript link. The property page of “A Link” looks like:

Property

A Link in Frame (in the updated page)

Item **Query**:

To open in new browser: ☐

Link tag:

Description:

[Help Description:](#)

* [Events](#) *

The “in Frame” specifies which child frame the link is in. For the first child frame, it is 1, the second child 2, and so on. Frame 0 is the current web page including any content created by dynamic scripts. Frame -1 is to use the source code of the current web page without any script interpretation. If there are several layers of frames, use something like “2.1.3” to access “the second child frame, and the first child frame, and the third child

frame.”

The “Query” link tests if the target query locates the correct link. In order to use the test correctly, navigate to the target page before opening the property page.

“Target query” is an HTQL expression to locate the link to click. The HTQL expression can also be taken from a Variable or an Expression.

If you check the “To open in new browser”, the target page will be opened in a new browser.

“Wait navigation?” specifies whether to wait for navigation after the click or not.

For the “Before each click”, you can decide if you want to reload the web page for clicking or not. It is most useful when you use this in the “Repeat property” for another action, so that you need to bring up the page before clicking on the Next link. Refer “Repeat property” in section “Repeat on Next Pages” for more help.

The “Link tag” specifies the HTML tag of the target link.

Make sure to click the [Modify] button after any change.

3) A List of Links

The property page looks like:

Property

A List of Links in Frame (in the source page)

Target [Query](#):

```
<form (Name='EntrezForm')>1.<div (ID='RightCol')>1.<div (ID='RightColContent')>1.<div (ID='ViewPanel')>1.<div (CLASS='DocSumTpl')>1.<div (CLASS='contentbox-left')>1.<div>3.<div>1-0{ COLUMN1=<div (CLASS='rprtNum')>1.<b>1;
```

Field Index for Links: (0: No links, 1~.: Field index)

To open in new browser: ☐

Tuple order: Before each click:

Do the next action:

Description:

[Help Description](#):

* [Events](#) *

The “in Frame” specifies which child frame the link is in. For the first child frame, it is 1, the second child 2, and so on. Frame 0 is the current web page including any content created by dynamic scripts. Frame -1 is to use the source code of the current web page without any script interpretation. If there are several layers of frames, use something like “2.1.3” to access “the second child frame, and the first child frame, and the third child frame.”

The “[Query](#)” link test if the target query locates the correct list of links. In order to use the test correctly, navigate to the target page before opening the property page.

“[Target query](#)” is an HTQL expression to locate the list of links to click. The HTQL expression can also taken from a [Variable](#) or an [Expression](#).

“[Field Index for Links](#)” specifies the table column where the links are in. If the field index is 0, the action is turned to ‘Take Table’, without following any links.

If you check the “[To open in new browser](#)”, the target pages will be opened in a new browser.

“[Tuple order](#)” specifies if to open the link sequentially or reversely. The “[Before each click](#)” specifies whether to reload the web page for link clicking or not. If the “[Do the next action](#)” is “[after completed this page](#)”, then it will click all links before doing the next action. Otherwise, it will repeatedly do the next action for each link.

Make sure to click the [\[Modify\]](#) button after any changes.

4) Take Data

The property page looks like:

The screenshot shows a web form titled "Property". At the top left is a dropdown menu labeled "Take Data". Below it is a "Target Query:" label followed by a text input field containing the HTQL expression: `<form (Name='propform')>1.<p>1`. To the right of this field is a dropdown menu currently set to "HTQL". Below the "Target Query" field is a "Wait time:" label followed by a text input field containing the placeholder text *(Can use expression)*. To the right of this field is another dropdown menu set to "-- Not wait for data --". Below the "Wait time" field is a "Description:" label followed by a text input field containing the text: `Extract data like 'Target Query : <pre>$'`. At the bottom of the form are two buttons: "Modify" and "Return".

The “[Query](#)” link tests if the target query extracts the correct data. In order to use the test correctly, navigate to the target page before opening the property page.

“[Target query](#)” is an HTQL expression to extract the data. The HTQL expression can also be taken from a [Variable](#) or an [Expression](#).

You can let the robot wait until the target data is shown or until the data disappears (work only when the page is refreshing by itself). “[Wait time](#)” specifies the wait interval in seconds.

Make sure to click the [\[Modify\]](#) button after any change.

5) Take Table

Take table is the same as “A List of Links”, when the “[Field Index for Links](#)” is set to 0 or left empty. Example page:

Property

A List of Links ▾

Target [Query](#):

```
<p>1-0{
COLUMN1=<a>1;
COLUMN2=<a>1:xx;
COLUMN3=<input (Name='Value')>1:fx;
COLUMN4=<select>1.<option>1;
```

HTQL ▾

Field Index for Links: (0: No links; 1~.: Field index)

Description:

The “[Query](#)” link tests if the target query locates the correct list of links. In order to use the test correctly, navigate to the target page before opening the property page.

“[Target query](#)” is an HTQL expression to locate the list of links to click. The HTQL expression can also be taken from a [Variable](#) or an [Expression](#).

“[Field Index for Links](#)” specifies the table columns where the list of links is defined in the target query. If the field index is 0, the action is turned to ‘Take Table’.

Make sure to click the [\[Modify\]](#) button after any change.

6) Submit a Form

The property page looks like:

Property

Submit Form ▾

Form [values](#): String ▾ [From database](#)

[Form](#) location HTQL:

[Submit button](#) HTQL:

To match form action:

Description:

The “values” link brings up the form value window from the control panel (Click on the [OK] button in the form value window to close it).

“Form values” specifies the form values to be filled in the target form. The form values can be specified as a “Variable” or an “Expression”. Form values can be drawn from databases, which will be explained later.

The “Form” link tests if the target form is located correctly. In order to use the test correctly, navigate to the target page before opening the property page.

“Form location HTQL” is an HTQL expression to locate the form to be submitted.

“Submit button HTQL” is an HTQL expression to locate the submit-button of the form. A special ‘-none-’ expression tell the robot not to submit the form (only fill form values).

“To match form action” tells the robot submit the form only the action matches the specified URL.

Make sure to click the [Modify] button after any change.

7) Logon Form

‘Logon Form’ is same as ‘Submit a Form’

8) Open a Frame

The property page looks like:

The screenshot shows a dialog box titled "Property". Inside, there is a dropdown menu set to "Open a Frame". Below this, there is a text field labeled "Url:" containing "<FRAME>1", followed by a dropdown menu set to "String". Underneath, there is a text field labeled "Description:" containing "Get a frame named 'http://localhost/ethobank/title.html'". At the bottom right, there are two buttons: "Modify" and "Return".

The Url specifies an HTQL expression to locate the frame. The Url can also be taken from a Variable or an Expression, such as ‘http://mail’ + ‘.yahoo.com’.

Make sure to click the [\[Modify\]](#) button after any change.

9) Sent Email

To be explained.

10) A Schedule

The property page looks like:

Property

A Schedule

[View Schedule](#)

Description: Schedule tasks

Modify

Return

Click on the “[View Schedule](#)” to modify the schedule. The schedule page looks like:

Scheduled Tasks						
	Type	Interval	Base	Action	ActionType	Session
<input type="checkbox"/>	DoEvery	1.5	Minute	check yahoo mails	CallPath	yahoomail1
<input type="checkbox"/>	DoEvery	1.5	Minute	check gmails	CallPath	gmail1
<input type="checkbox"/>	DoEvery	1.5	Minute	check hotmails	CallPath	hotmail1
<div> <div>Modify</div> <div>Return</div> <div>Insert</div> <div>Delete</div> </div>						

To add a new schedule, click the [\[Insert\]](#) button. To modify existing schedules, check the schedules and click the [\[Modify\]](#) button. To delete existing schedules, check the schedules to be deleted and click the [\[Delete\]](#) button. The insert and modify page looks like:

Scheduled Tasks						
	Type	Interval	Base	Action	ActionType	Session
<input type="checkbox"/>	DoEvery	1.5	Minute	check yahoo mails	CallPath	yahoomail1
<input checked="" type="checkbox"/>	Do Every	1.5	Minute	check gmails	CallPath	gmail1
<input checked="" type="checkbox"/>	Do Every	1.5	Minute	check hotmails	CallPath	hotmail1
				Modify	Return	Insert
				Delete		

Customize the Type, Interval, Base, and Action to desirable settings and click [Modify] button to confirm the change.

Leave the Session attribute empty to use a default browser. Give a special Session name if you want the robot to launch a special browser for the scheduled task.

6. Frequently Asked Questions

6.1. Can I create a dummy action as a placeholder?

Yes. You can add a Schedule action, without any associated schedule items, as a dummy action. You can also associate events with this dummy action.

6.2. The robot skips certain actions during navigation!

This most probably is because of the slow Internet connection, where the robot has determined a timeout for the previous action, and continues on the next action. To solve this problem, you can set a slower robot navigation speed from the menu: Settings → Navigation Speed → Very Slow; or you manually introduce a sleep after the previous action (add an event “After each tuple”, and set the value as “sleep(10)”, which sleeps for 10 seconds).

If the previous action is a Click action, it may be because the action was setup not to wait for the target page (this is useful when clicking on some Ajax actions). You can right-click on the action and select the “Property ...” from the popup menu to see if it waits for navigation. Change to “Wait” if it is “No wait”.

6.3. How can I scrape data from pop-up windows?

You can either open pop-up windows in the current browser (using menu: Settings → Popup → Open in Current Browser) or check the “To open in new browser” in the action property.

6.4. Save Variables does not work with MS-Access database, but it works with XML

Your table names or field names in MS-Access database may include special characters, such as spaces or non alpha-numerical characters.

6.5. How to run robots from command line?

You can use a command line like: "irobot.exe robot1.irb /run:taskname /yes /exit" to run the robot and exit after the robot completed.

7. Additional Information

Online Forum (preferred): <http://irobotsoft.org/bb/>

Technical Support: support@irobotsoft.com

FREE donation: <http://irobotsoft.com/buy.htm>

Collaboration: info@irobotsoft.com

Copyright

IrobotSoft.com May. 2005

All rights reserved.