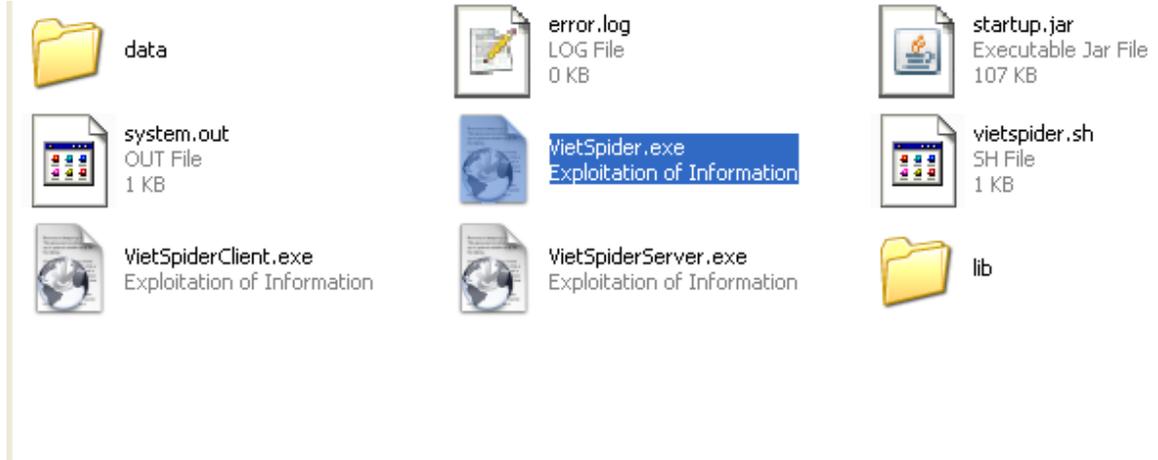# VietSpider Quick Guide – Step by Step – extract data from Android Market
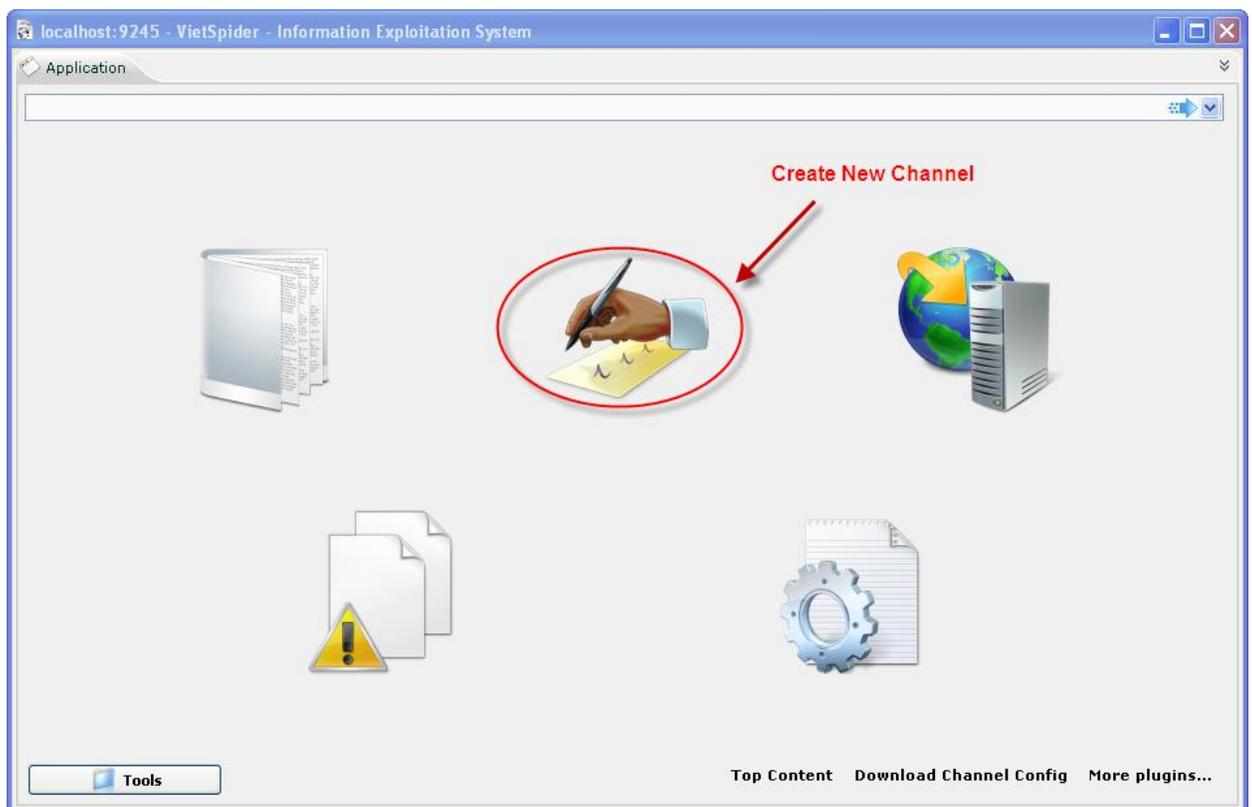
## Part 1: Create new channel

1. **Launch VietSpider**
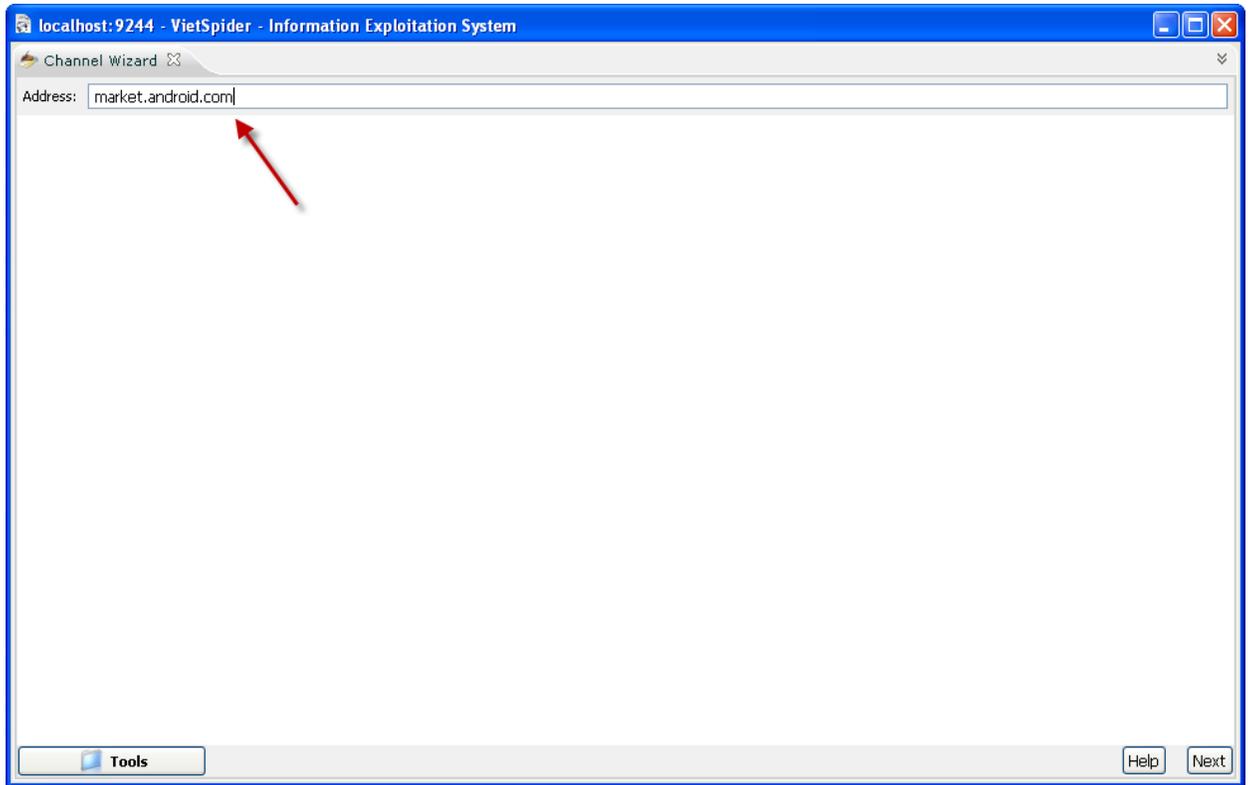


Open VietSpider folder and double click on VietSpider.exe
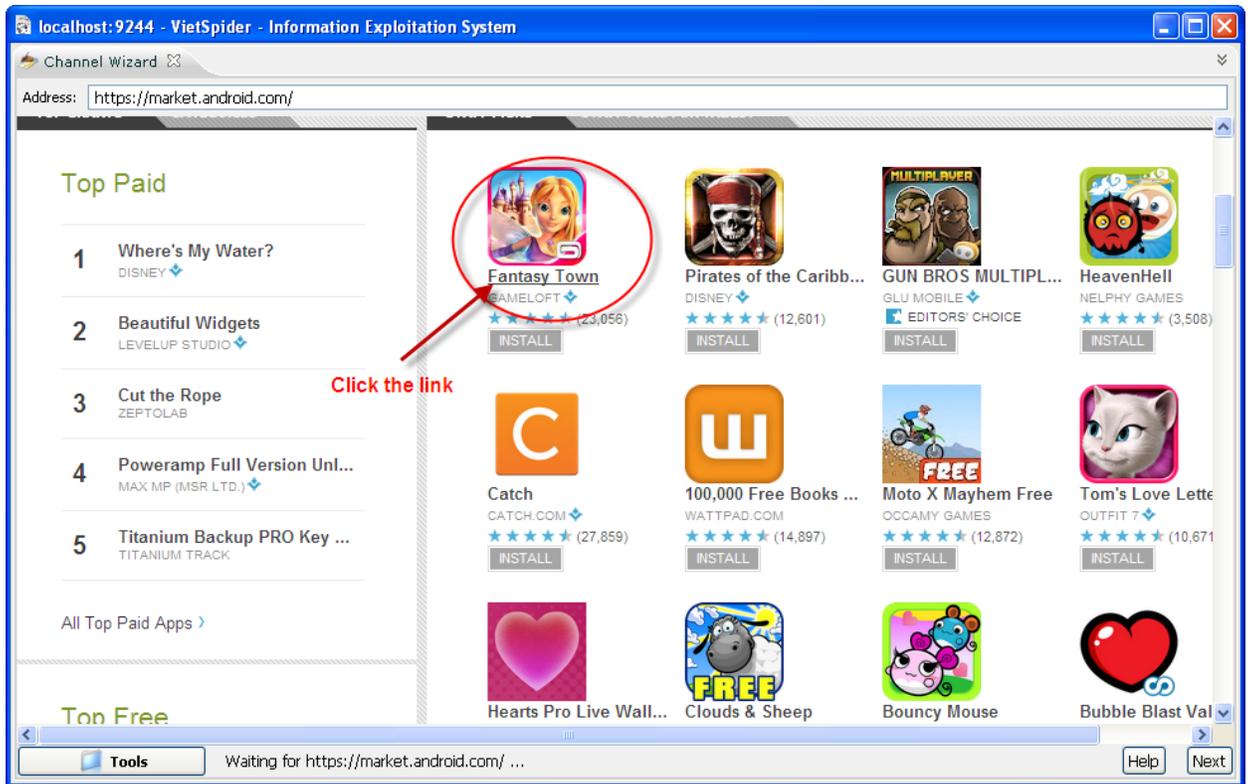
2. **Open Create Channel Wizard**



Click on Create New Channel icon
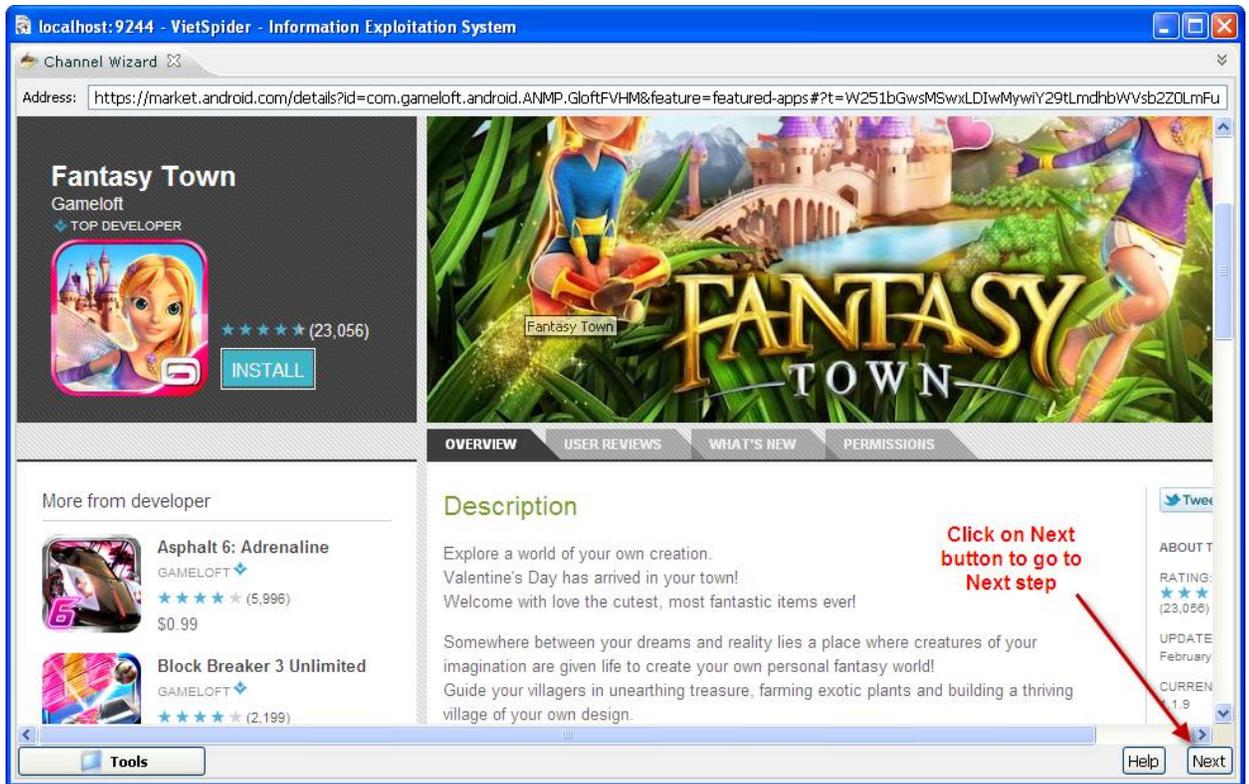
**3.  Browse  website**



Input  your address into the Address text box.

**4.  Select a data page**

Select a sample data page. This page will be used for selecting extract data area.

5. **Go to Next step**

6. **Select Character Encoding for this website**

The character encoding selected here will be used to crawl data from this channel. If the web page displays properly with character encoding, please click on Next button to go to next step.

**7. Define data areas on the data page**



Click on HTML Node to select data area that you want VietSpider extracts information.

**8. Auto select html node**

Please block a passage of text, VietSpider will suggest the node that contains the equivalent data.

**9.    Add extract data node(s)**

Click Add button to add the current path (from selected HTML Node) to the HTML Path list.  The list contains path(s) to  data area(s) that you want to extract. Click Next button to go to next step.

10. **Select category for channel.**

Select a category in the list and click Next button.

11. **Define XML elements of data document.**

You must define XML schame for extracted data document. Please define XML elements by input tag name and add it to the list. Click on Next button when you done.

12. **Define data area for XML elements**

Define data area for XML element, you must select a tag name -> select the node(s) that contains data area then add equivalent path to the list.

There are 4 data type for element:
- Default – VietSpider will auto detect data for equivalent element. That could be CDATA or Text
- Text – the element contains text value.
- CDATA – the element contains HTML tag.
- File  - data type is used for image tag or download file link.

13.  **Define data link template.**

Data Link Template: this is template that will be used for VietSpider filters data link(s) from link collection when the Crawler crawl data.

Summary of regular-expression constructs

Character – Matches
- *       - any character
- @    -  letter character
- $     -  digit character

14. **Add the Start page collection.**

The Start Page(s): the list of links that VietSpider will use to visit and find data link when crawl data. You can drap and drop the link from the browser.

15. **Define channel name**

Input the channel name and click Next button to go to the last step.

16. **Finish and save the channel**

If you want to test data extracting, you can copy another data link from the website and paste it to the Sample Data Page text box and click on Test button. Click Done button to save the channel to the Channel store.

## Part 2: Crawl data by the configured channel.

17. **Open Crawler**

From VietSpider user interface, click Tools button -> click on Crawler icon to open Crawler.

**18. Add channel(s) to crawling list.**

Click on Crawl Channel button and add the channel(s) that you want to download and extract data. Click Start Crawling button for starting crawling data from website(s).

## Part 3: Repair the channel.

19. **Open the Channel Store**



Click Tools button then click on Channel Store icon.

20. **View  website in the browser.**

Select the channel that you want repair it. From Start Page(s) text box, right click then select View In Browser.

21. Copy a sample data link

Copy a sample data link for testing defined extract data area.

**22. Test and repair Data Link Pattern**



Paste the link to Sample Data Page text box, you can see the error message from VietSpider. Please add the below patterns to Data Link Pattern to fix it.

**https://market.android.com/details?id=*&feature=***

**23. Add New Data Link Pattern**

Copy and Paste new data link pattern to the text box and click Add icon at the end text box.

Click Test button

24. **Test extracted data.**

When you click on Test button, you can see extracted data document is empty.

25. **Fix extract data for XML element.**

Please review extract data and extract region, click  icon.

26. **Update incorrect path**

Review the path of XML elements and the path of extract region, update when you see the path is incorrect. Save the channel and re-crawl data when you done.

## Part 4: Browse the crawled Data

Click Tools -> Browse Content

Crawler ⬛ 13.02.2012 ✕

vietspider/DOMAIN/1/13.02.2012

1 | 2  ( 2 )

**Angry Birds**
Use the unique powers of the Angry Birds to destroy the greedy pigs' fortresses!

The survival of the Angry Birds is at stake. Dish out revenge on the greedy pigs who stole their eggs. Use the unique powers of each bird to destroy the pigs' fortresses. Angry Birds features challenging physics-based gameplay and hours of replay value. Each of the 300 levels requires logic, skill, and force to solve.

Terms of Use: http://www.rovio.com/eula
Privacy Policy: http://www.rovio.com/privacy

]]>

Products / App Android / ~market.android.com /22:40:18

**RAIDEN-Sky Force Ace**
FOREVER FREE FULL VERSION DOWNLOAD!! NO IN APP BILLING!!!NO ADS!!!!

DEATH TO ALIEN ENEMIES!!
AN AIR BATTLE IN THE GALAXY IS INEVITABLE!!
TRY TO JOIN THE WAR AND DEFEAT THE ALIEN ENEMIES AS THE HERO OF THE GALAXY ALLIANCE!!

INTRODUCTION
The galaxy is on fire!
The evil force is conquering the kind
HOT ANDROID AIR BATTLE GAME
It will let you recall the old school games and bring you to the good old days!

**Tools**    Done      **Export Data to CSV**   **Export to Excel**   **More plugins...**

**Tools** window:
- Browse Content
- Create New Channel
- Channel Store
- Crawler

---

Crawler ⬛ http://thuannd:9245/vietspider/DETAIL/20... ✕

http://thuannd:9245/vietspider/DETAIL/201202132240180010

13/02/2012
Previous  1/1  Next
**XML.Products**

```xml
<?xml version="1.0" encoding="UTF-8"?>
<document>
    <name>Angry Birds</name>
    <desc>
        <![CDATA[ <div id="doc-description-container"class="doc-description-collapsed"> <div class="doc-description toggle-overflow-contents"> <div id="doc-original-text"> Use the unique powers of the Angry Birds to destroy the greedy pigs' fortresses! <p> The survival of the Angry Birds is at stake. Dish out revenge on the greedy pigs who stole their eggs. Use the unique powers of each bird to destroy the pigs' fortresses. Angry Birds features challenging physics-based gameplay and hours of replay value. Each of the 300 levels requires logic, skill, and force to solve. </P> <p> Terms of Use: <A href="http://www.google.com/url?q=http%3A//www.rovio.com/eula&#38;usg=AFQjCNFE_-7fBP3ZKkk3RXAV1sVMGcHqEQ" target="_blank"> http://www.rovio.com/eula </A> <br> Privacy Policy: <A href="http://www.google.com/url?q=http%3A//www.rovio.com/privacy&#38;usg=AFQjCNGyZ7cMPyY_EYqqduZOUHdQP3PV3Q" target="_blank"> http://www.rovio.com/privacy </A> </P> </DIV> </DIV> <div class="doc-description-overflow"> </DIV> </DIV>]]>
    </desc>
    <size>15M</size>
    <rating>Low Maturity</rating>
    <src>
        <![CDATA[https://market.android.com/details?id=com.rovio.angrybirds&feature=top-free]]>
    </src>
</document>
```

**Tools**