

Introduction

This is an informal manual on the gpu search engine 'gpuse'. There are some other documents available, this one tries to be a practical how-to-use manual.

Table of Contents

Introduction.....	1
GPU Search Engine (GPUSE).....	2
Contact.....	2
So, you wanted to crawl the web?.....	3
Consider a few things:.....	3
Minimum System Specifications.....	3
Is GPUSE safe to run?.....	3
Is it safe for others (third party). Will it not 'DoS' sites?.....	3
Screw your isp's dns server.....	4
Can i view my crawled results immediately?.....	4
Can i search in real-time over the p2p network?.....	4
Can i search without crawling?.....	4
Sure you can. Just launch the frontend.....	4
So you think you are ready? Ok. Let's go.....	5
Some screen shots that may help guide you trough the process of enabling the crawler.....	5
Enabling the GPU Search engine plugin.....	5
The Search engine front end.....	6
Check what your crawler is doing.....	7
Enabling the webserver.....	8
Security considerations.....	8
What is more.....	9
And, how to see if my gpu returns results?.....	9
With gpu task list you can see the search is executed:.....	10
After waiting a few seconds, you can see the results:.....	11
You can sort the results by clicking the column name:.....	12
Of course, you can also use the web interface on localhost:.....	13
The reverse index.....	14
Investigating the gpu network with netmapper.....	15
A numeric representation:.....	15
A Graphical representation:.....	15
The crawler's progress.....	16
Concluding.....	17
Many thanks to all that keep the GPU network alive!.....	17

GPU Search Engine (GPUSE)

GPUSE comes with the GPU “Distributed Computing over a peer-to-peer network” package. GPU is the glue that holds several plugins, distributes requests, shows statistics and hosts the chat.

Contact

The chat is our main system for users and developers to exchange information. There is also a mailing list hosted at sourceforge “gpu-world” but you may notice direct contact may be very productive and interesting. The GPU chat is very interactive, the whiteboard can be used to exchange screenshots etc.

So, you wanted to crawl the web?

Consider a few things:

- Good hardware. Especially if your local database grows, it will take some system resources. Especially harddisk access, and system cache, are a factor. Another one is a good internet connection. A stable DNS server isn't a luxury either.
- Permanent connection. Preferably you have a permanent internet connection, so you can leave the crawler running overnight. Unlimited usage or a Fair Use Policy will help. Some ISP's will have another interpretation of fair use than you have, and with continuous crawling while time is ticking the amount of data easily cumulates. Be sure you know what you are doing.
- Patience. Continues crawling with 1 or 2 crawlers is much more efficient than short bursts with many crawlers.
- Controlled startup/shutdown of the gpu main application. With the search engine enabled, it may take (much) longer for gpu main application to completely shutdown. In case of a shutdown, first, all crawler threads need to be stopped. Then all recent collected data have to be flushed to disk. Normally this is done in 10-30 seconds or so, but if you have a busy crawler it may take longer.

Minimum System Specifications

The minimum recommended specs are:

- broadband. 50Kb/s or more.
- Reasonable cpu (133MHz or more).
- Modern harddrive (6Gb is really absolute minimum, 17 Gib a practical minimum, 120Gb or more will perform a lot better and is recommended.).
- Sufficient memory; 512Mb is no luxury. With 256Mb things will work, but you will notice the system is less usable as workstation in that case. 512 Mb with 2-3 crawlers running is reasonable stable on long term.
- Operating systems: XP, Windows 2000, 98 (not really recommended but should work), Linux: 2.4 or 2.6 kernel using wine. Most common releases should be working (Debian. Redhat) & wine (most recent releases) working. Wine may have issues.

Is GPUSE safe to run?

No and yes. With recent changes, we think yes, it is. There have been done a few modifications that significantly reduce disk access. Your harddisk should stay almost completely silent (no extraordinary head movements).

Also, there are some protection mechanisms to avoid network overload. In case a thread detects network 'malfunction' or slowness, it will sleep for half a minute before continuing with the next job, hereby avoiding abuse, overloading your DNS server, remote server or just your local system.

No because this is explicitly a beta test. Loads of situations might occur that somehow or another leads to undesired behavior. Although tested, we cannot guarantee an undesired behavior not to occur.

Is it safe for others (third party). Will it not 'DoS' sites?

We believe the crawler behaves very correctly. First of all, it examines the robots file that may be available, and strictly obeys instructions. Robots.txt may hold specific instructions for agent

'gpuse'. Then, a single crawler will never in short time access more than 5 documents of a single domain. Normally this will be less, but since there is some random factor involved, we limited number of pages from one site to 5 each time the crawler fetches a new set of urls. This number of 5 is a weighted efficiency factor that is introduced to reduce both database access and DNS Server requests. A local crawler gets urls on certain domains in batches of maximum 5 and cumulates until it reaches 100 or more urls or times out.

On startup with an empty database, the crawler has a few urls's hardcoded built in. Those urls typically point to pages that link to loads of other pages. Once the database grew to reasonable size, randomness is the factor that distributes the load.

Screw your isp's dns server.

Well probably you will not screw it but may get slow. "Crashing" a DNS server is definitively not impossible. Both scenario's may affect you and possibly other innocent users. Take care what you do. A single crawler is like 5 cyberpunks on speed surfing random sites non-stop 24h a day.

Can i view my crawled results immediately?

The first result set will be finished after about 1-2 hours of crawling with 1 or 2 crawlers on a typical broadband connection. You can view results of other crawlers on the network though.

Can i search in real-time over the p2p network?

Yes, you can, using the frontend. The search frontend also has a simple web server built in. For a collective reverse index, we are experimenting on a system based on mysql at <http://search.dubaron.com>

Can i search without crawling?

Sure you can. Just launch the frontend.

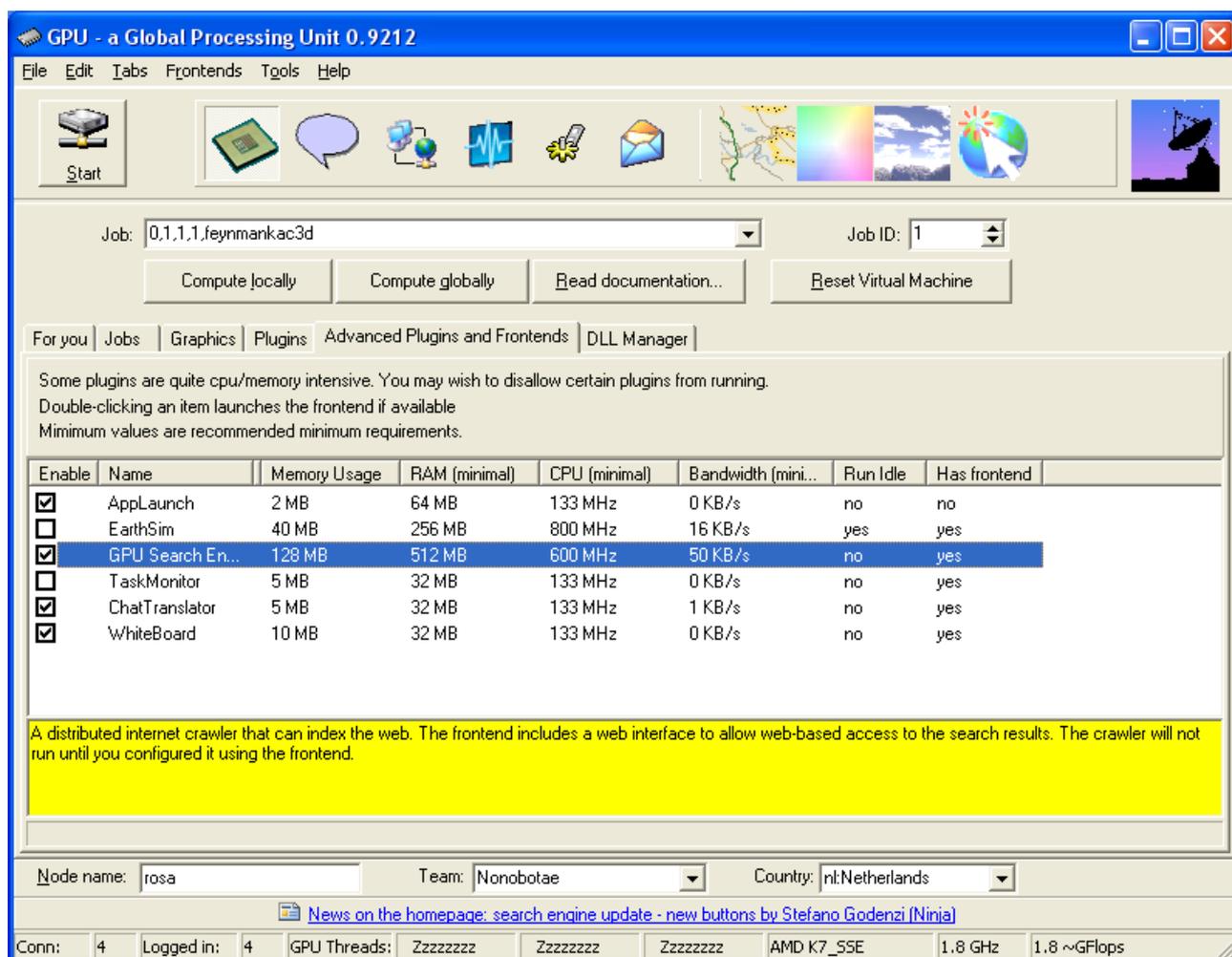
So you think you are ready? Ok. Let's go...

Some screen shots that may help guide you through the process of enabling the crawler

Enabling the GPU Search engine plugin

GPUSE needs to be enabled on two places. First, within the GPU main application, you have to enable the search engine plugin. Then **restart GPU** to activate it.

In the screenshot below, you can see the activated tabsheets. The line with the search engine is selected. Make sure the check box is checked.



The screenshot shows the GPU application window titled "GPU - a Global Processing Unit 0.9212". The interface includes a menu bar (File, Edit, Tabs, Frontends, Tools, Help), a toolbar with various icons, and a main workspace. The "Plugins" tab is active, displaying a table of installed plugins. The "GPU Search Engine" plugin is highlighted in blue, and its "Enable" checkbox is checked. Below the table, a yellow box contains a description of the crawler. At the bottom, there are fields for "Node name", "Team", and "Country", along with a status bar showing system information.

Enable	Name	Memory Usage	RAM (minimal)	CPU (minimal)	Bandwidth (mini...	Run Idle	Has frontend
<input checked="" type="checkbox"/>	AppLaunch	2 MB	64 MB	133 MHz	0 KB/s	no	no
<input type="checkbox"/>	EarthSim	40 MB	256 MB	800 MHz	16 KB/s	yes	yes
<input checked="" type="checkbox"/>	GPU Search En...	128 MB	512 MB	600 MHz	50 KB/s	no	yes
<input type="checkbox"/>	TaskMonitor	5 MB	32 MB	133 MHz	0 KB/s	no	yes
<input checked="" type="checkbox"/>	ChatTranslator	5 MB	32 MB	133 MHz	1 KB/s	no	yes
<input checked="" type="checkbox"/>	WhiteBoard	10 MB	32 MB	133 MHz	0 KB/s	no	yes

A distributed internet crawler that can index the web. The frontend includes a web interface to allow web-based access to the search results. The crawler will not run until you configured it using the frontend.

Node name: rosa Team: Nonobotae Country: nl:Netherlands

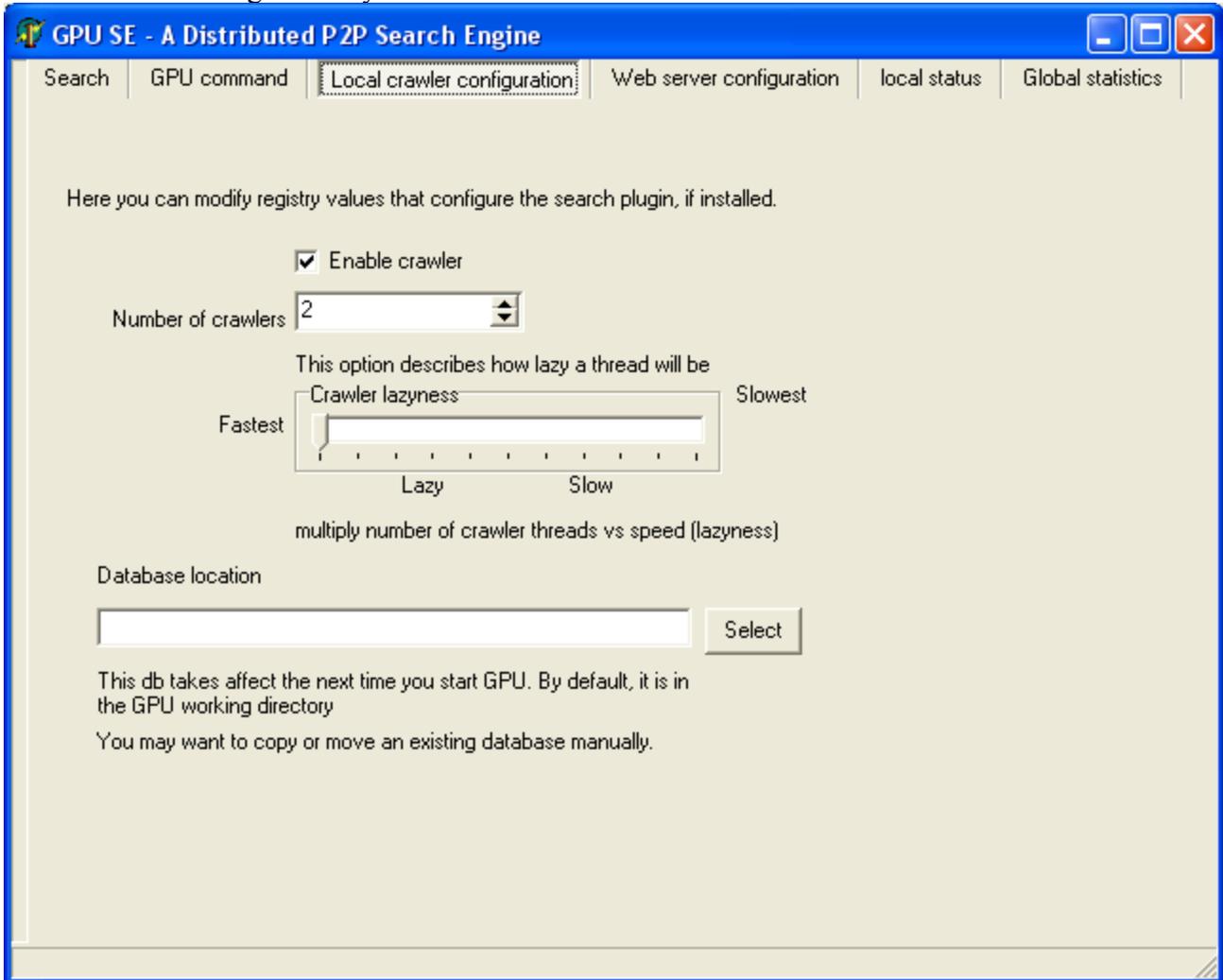
News on the homepage: [search engine update - new buttons by Stefano Godenzi \(Ninja\)](#)

Conn: 4 Logged in: 4 GPU Threads: Zzzzzzzz Zzzzzzzz Zzzzzzzz AMD K7_SSE 1.8 GHz 1.8 ~GFlops

The Search engine front end

After restarting the gpu application, enable the search engine frontend. Here you can control the crawler in real-time. There may be some latency before the crawler responds, varying from seconds to minutes, depending on the setting changed.

With tab sheet 'local status' you see logbooks of what the crawler is doing. With tab sheet 'Local crawler configuration' you can enable it:



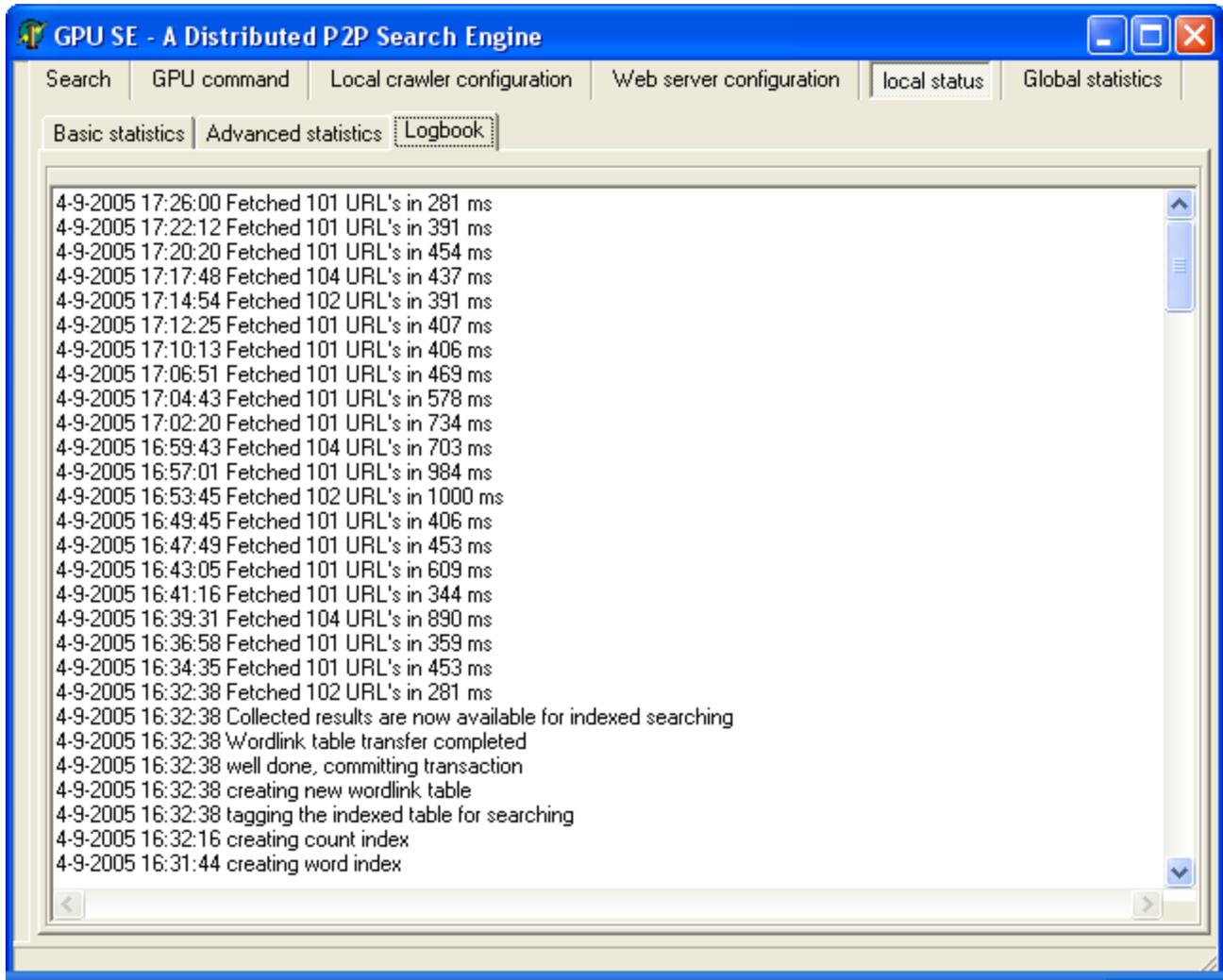
The default number of crawlers is 2. Use of 1 crawler is also very suitable. Remember that continues crawling with 1 or 2 crawlers is much more efficient than short bursts with many crawlers.

We set no maximum to the number of crawlers, but a real-life maximum is about 6. In general, it seems having 3 or 4 crawlers running will already take most of your available bandwidth.

Any number above 12 would be absurd, even on high end hardware.

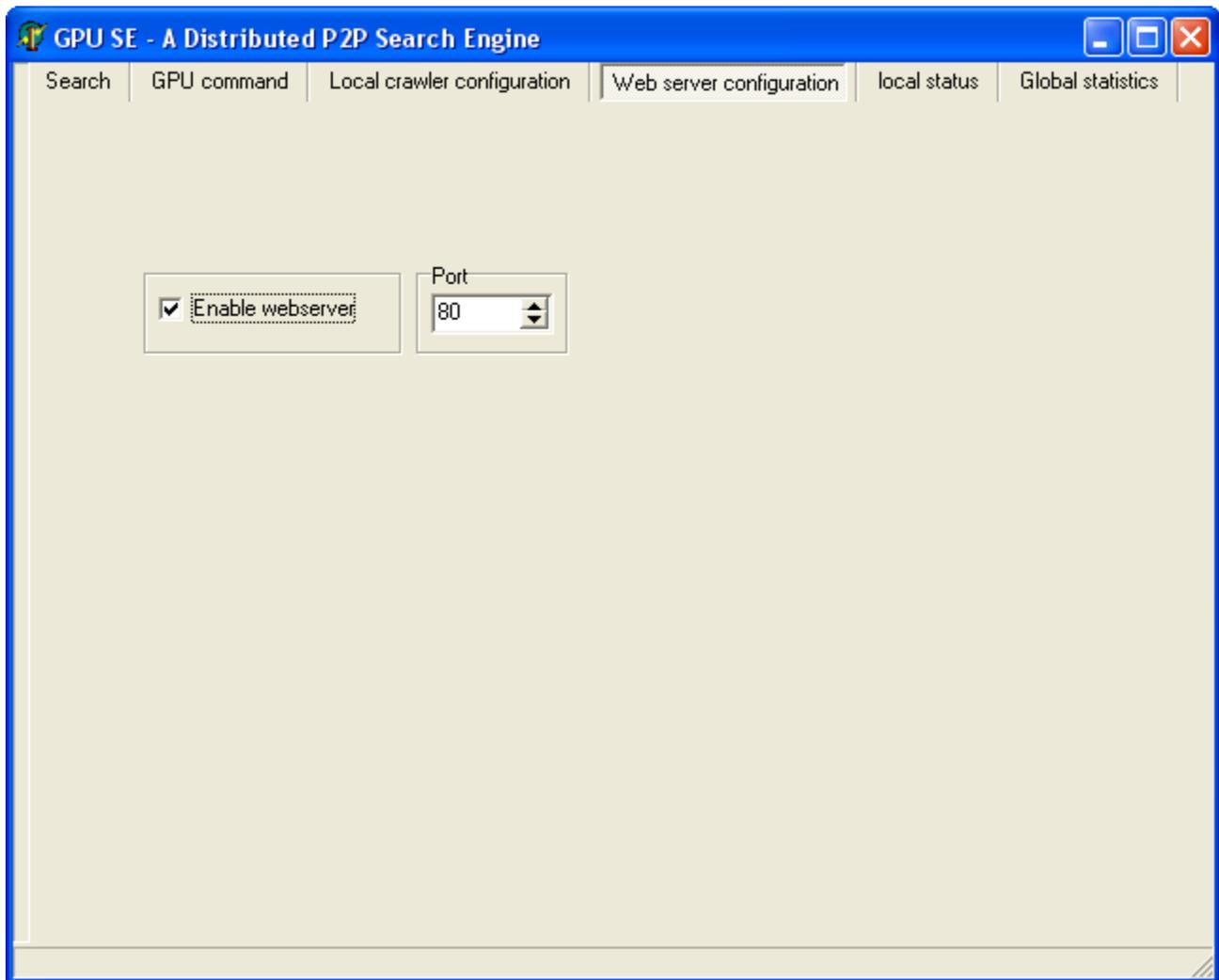
Check what your crawler is doing.

Check the tabsheet logbook:



If it says “Fetched 5 url's” or so, this is normal during initial startup phase. However, quite soon you should see numbers above 100. If it stays below, you have too few robots threads crawling. Reduce the number of crawlers in that case, and give the robots threads a chance to build a nice database of crawlable url's.

Enabling the webserver



Make sure the checkbox is checked

Use your browser to navigate to <http://localhost/>

If you use another port number than 80, for example 81, tell your browser:

<http://localhost:81/>

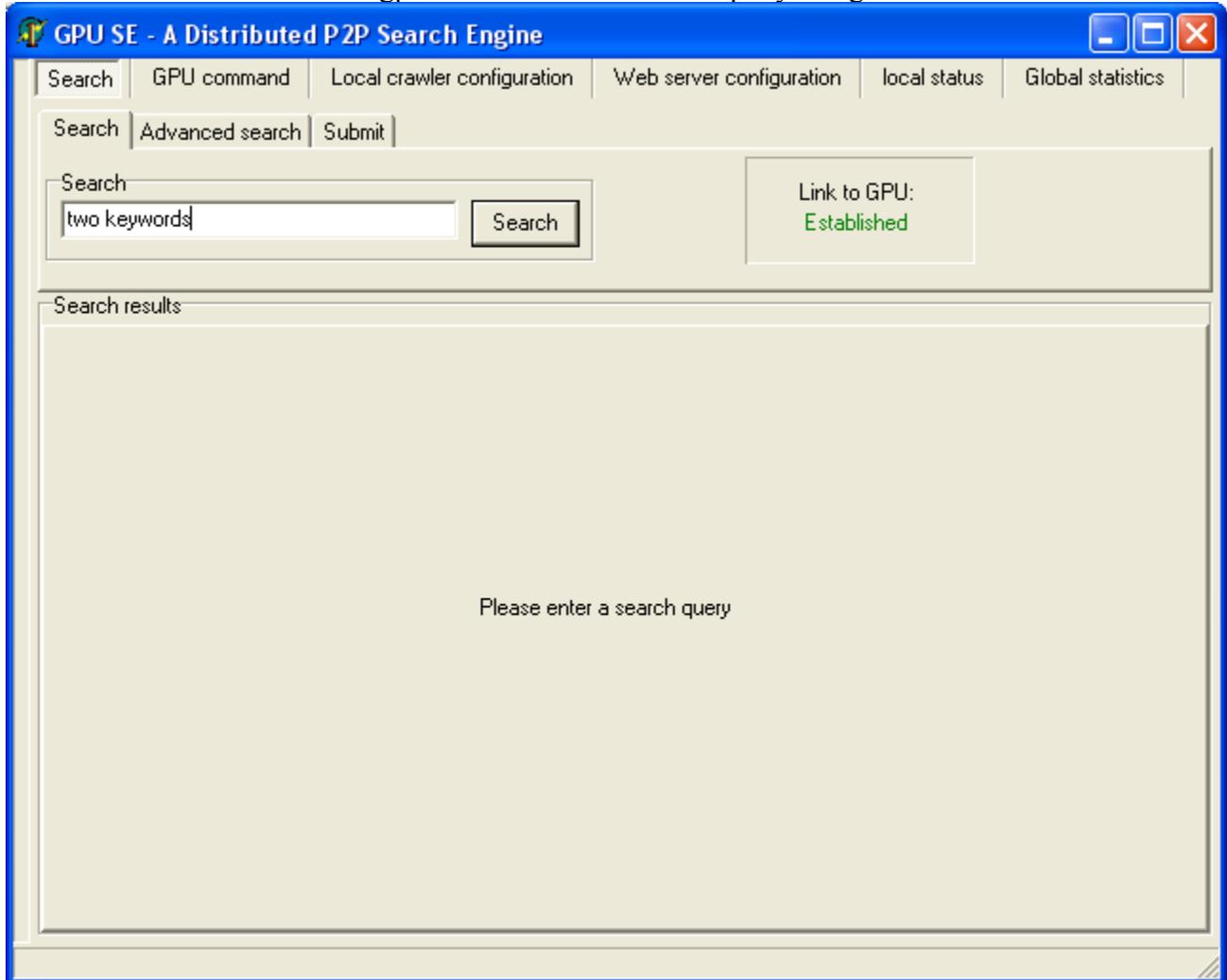
Security considerations

Please keep in mind that if you enable the web interface others can use it as well. Although it is multi-threaded, there is a chance of overloading the gpu network, so take care.

What is more

And, how to see if my gpu returns results?

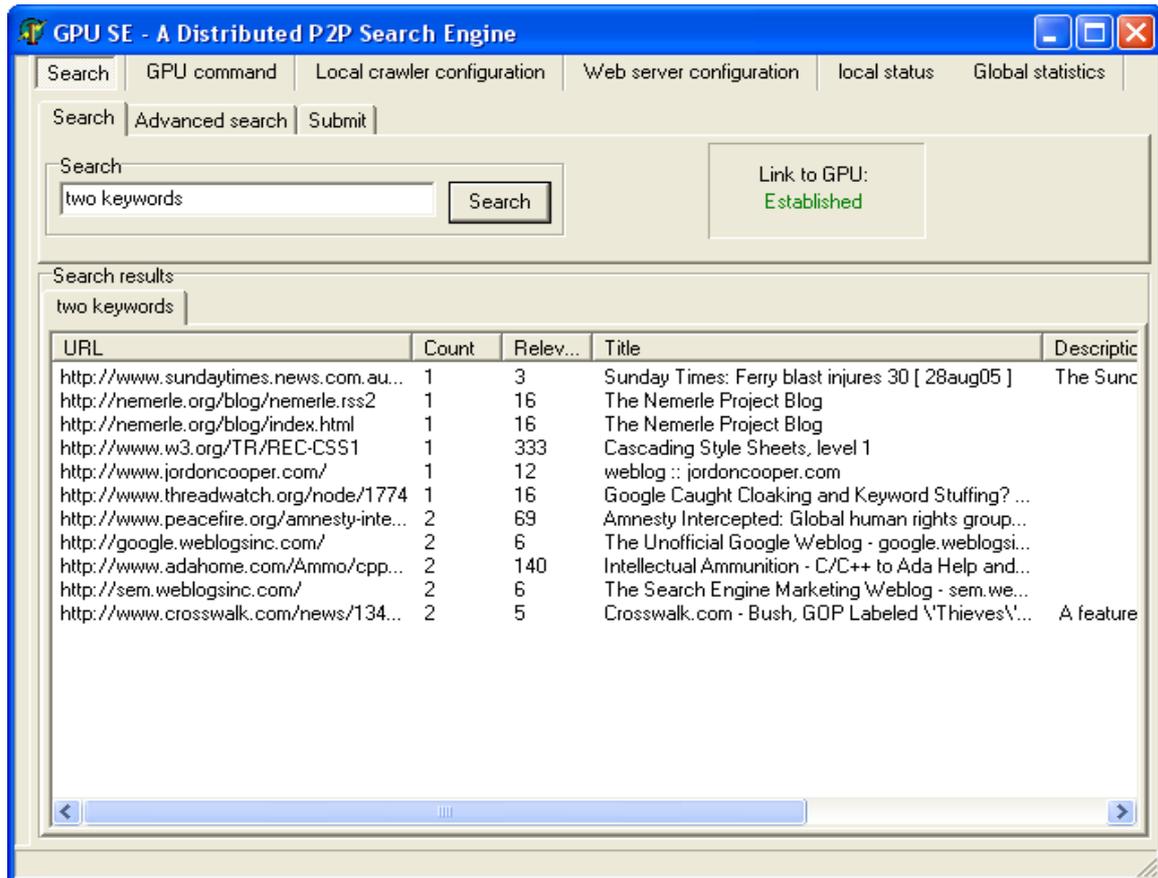
You can see that at the gpu task list. Enter a search query using the search frontend:



With gpu task list you can see the search is executed:

You see it twice due to an (unfiltered) roundtrip. You got the same request back from another client. This is a gnutella issue we are working on.

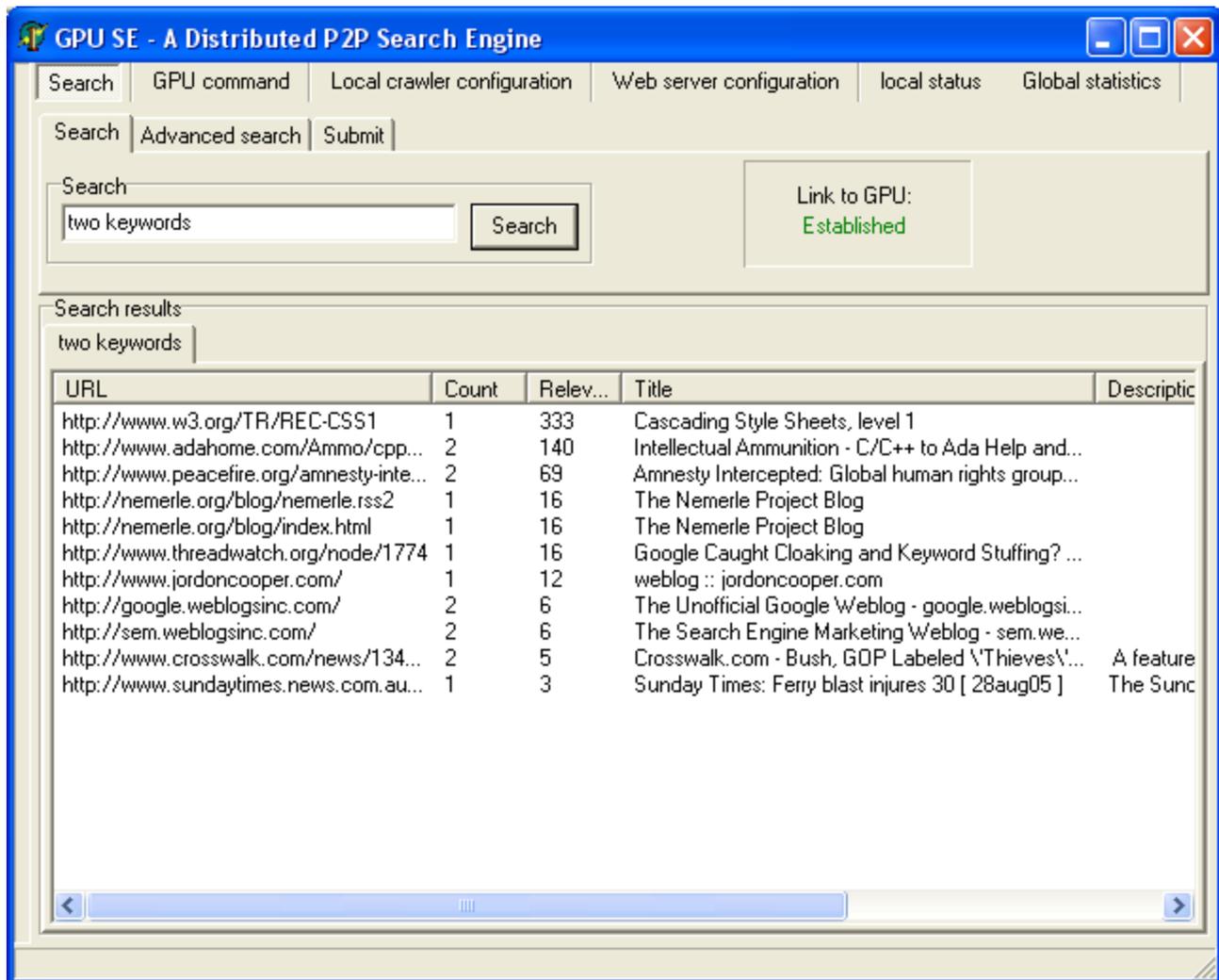
After waiting a few seconds, you can see the results:



The screenshot shows the GPU SE - A Distributed P2P Search Engine interface. The window title is "GPU SE - A Distributed P2P Search Engine". The interface includes a search bar with the text "two keywords" and a "Search" button. To the right of the search bar, there is a "Link to GPU: Established" button. Below the search bar, there is a "Search results" section with a sub-section for "two keywords". The search results are displayed in a table with the following columns: URL, Count, Relev..., Title, and Descriptic.

URL	Count	Relev...	Title	Descriptic
http://www.sundaytimes.news.com.au...	1	3	Sunday Times: Ferry blast injures 30 [28aug05]	The Sunc
http://nemerle.org/blog/nemerle.rss2	1	16	The Nemerle Project Blog	
http://nemerle.org/blog/index.html	1	16	The Nemerle Project Blog	
http://www.w3.org/TR/REC-CSS1	1	333	Cascading Style Sheets, level 1	
http://www.jordoncooper.com/	1	12	weblog :: jordoncooper.com	
http://www.threadwatch.org/node/1774	1	16	Google Caught Cloaking and Keyword Stuffing? ...	
http://www.peacefire.org/amnesty-inte...	2	69	Amnesty Intercepted: Global human rights group...	
http://google.weblogsinc.com/	2	6	The Unofficial Google Weblog - google.weblogsi...	
http://www.adahome.com/Ammo/cpp...	2	140	Intellectual Ammunition - C/C++ to Ada Help and...	
http://sem.weblogsinc.com/	2	6	The Search Engine Marketing Weblog - sem.we...	
http://www.crosswalk.com/news/134...	2	5	Crosswalk.com - Bush, GOP Labeled \Thieves\...	A feature

You can sort the results by clicking the column name:



The screenshot shows the GPU SE - A Distributed P2P Search Engine interface. The window title is "GPU SE - A Distributed P2P Search Engine". The interface includes a search bar with the text "two keywords" and a "Search" button. To the right, there is a "Link to GPU: Established" button. Below the search bar, the "Search results" section displays a table with the following data:

URL	Count	Relev...	Title	Descriptic
http://www.w3.org/TR/REC-CSS1	1	333	Cascading Style Sheets, level 1	
http://www.adahome.com/Ammo/cpp...	2	140	Intellectual Ammunition - C/C++ to Ada Help and...	
http://www.peacefire.org/amnesty-inte...	2	69	Amnesty Intercepted: Global human rights group...	
http://nemerle.org/blog/nemerle.rss2	1	16	The Nemerle Project Blog	
http://nemerle.org/blog/index.html	1	16	The Nemerle Project Blog	
http://www.threadwatch.org/node/1774	1	16	Google Caught Cloaking and Keyword Stuffing? ...	
http://www.jordoncooper.com/	1	12	weblog :: jordoncooper.com	
http://google.weblogsinc.com/	2	6	The Unofficial Google Weblog - google.weblogsi...	
http://sem.weblogsinc.com/	2	6	The Search Engine Marketing Weblog - sem.we...	
http://www.crosswalk.com/news/134...	2	5	Crosswalk.com - Bush, GOP Labeled \Thieves\...	A feature
http://www.sundaytimes.news.com.au...	1	3	Sunday Times: Ferry blast injures 30 [28aug05]	The Sunc

Of course, you can also use the web interface on localhost:



The reverse index

Or, even more sophisticated, using one of the web interfaces that maintain a permanent reverse index (discussion of how to set up this mysql webinterface in another document):



Web Search using GPU - Mozilla Firefox

Bestand Bewerken Beeld Ga Bladwijzers Extra Help

http://search.dubaron.com/index.php?offset=1&numresults=50&query=twc Ga

Beginnen Laatste nieuws

GPU

make LOVE -not WAR

Search the web with P2P search technology

Reverse index

Statistics:
Indexed urls: 15042
Indexed queries: 462
Number links: 29431

Enter your search request:

For best results, use a maximum of two keywords

Questions? Read our [FAQ](#)

Results as fetched from the P2P network in real time combined with cached results:

Real-time index

Statistics:
Known urls: 1019281
Known domains: 320665
Crawled pages: 200281

Top 20 searches:
[gpu](#)
[test](#)
[torrent](#)
[nova-ts](#)
[paid grid computing](#)
[open source](#)
[mp3](#)
[kek](#)
[wiki](#)

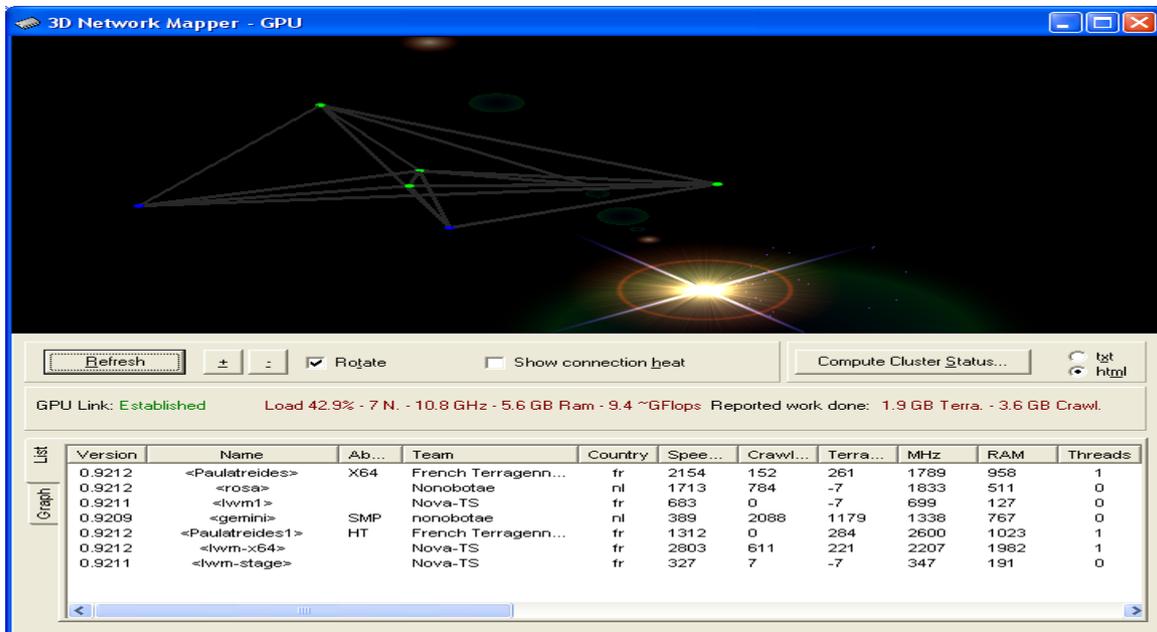
- ◆ **Intellectual Ammunition - C/C++ to Ada Help and Discussion** (126502 bytes, score: 140)
<http://www.adahome.com/Ammo/cpp2ada.html>
- ◆ **SEO - Search Engine Optimization and Position Ranking** (17382 bytes, score: 410) [cached](#)
Search Engine Optimization and Position Ranking
<http://www.addme.com/optimization.htm>
- ◆ **Amnesty Intercepted: Global human rights groups blocked by Web censoring software** (29659 bytes, score: 69)
<http://www.peacefire.org/amnesty-intercepted/>
- ◆ **SEARCH** (1877 bytes, score: 410) [cached](#)
SEARCH
<http://www.wholewheatradio.org/jbb/search/index.php>
- ◆ **The Nemerle Project Blog** (111623 bytes, score: 16)
<http://nemerle.org/blog/nemerle.rss2>

Klaar

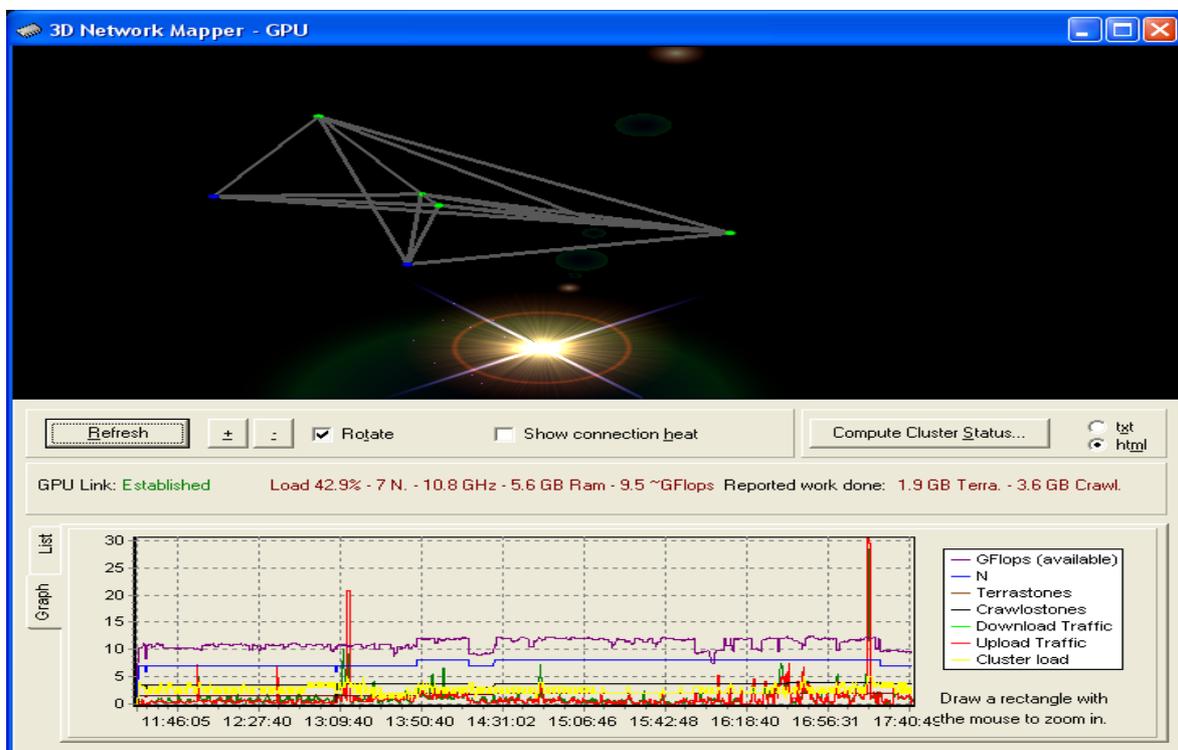
Investigating the gpu network with netmapper

With netmapper, we can see what the nodes are crawling. The newest software has just been installed and is being tested:

A numeric representation:

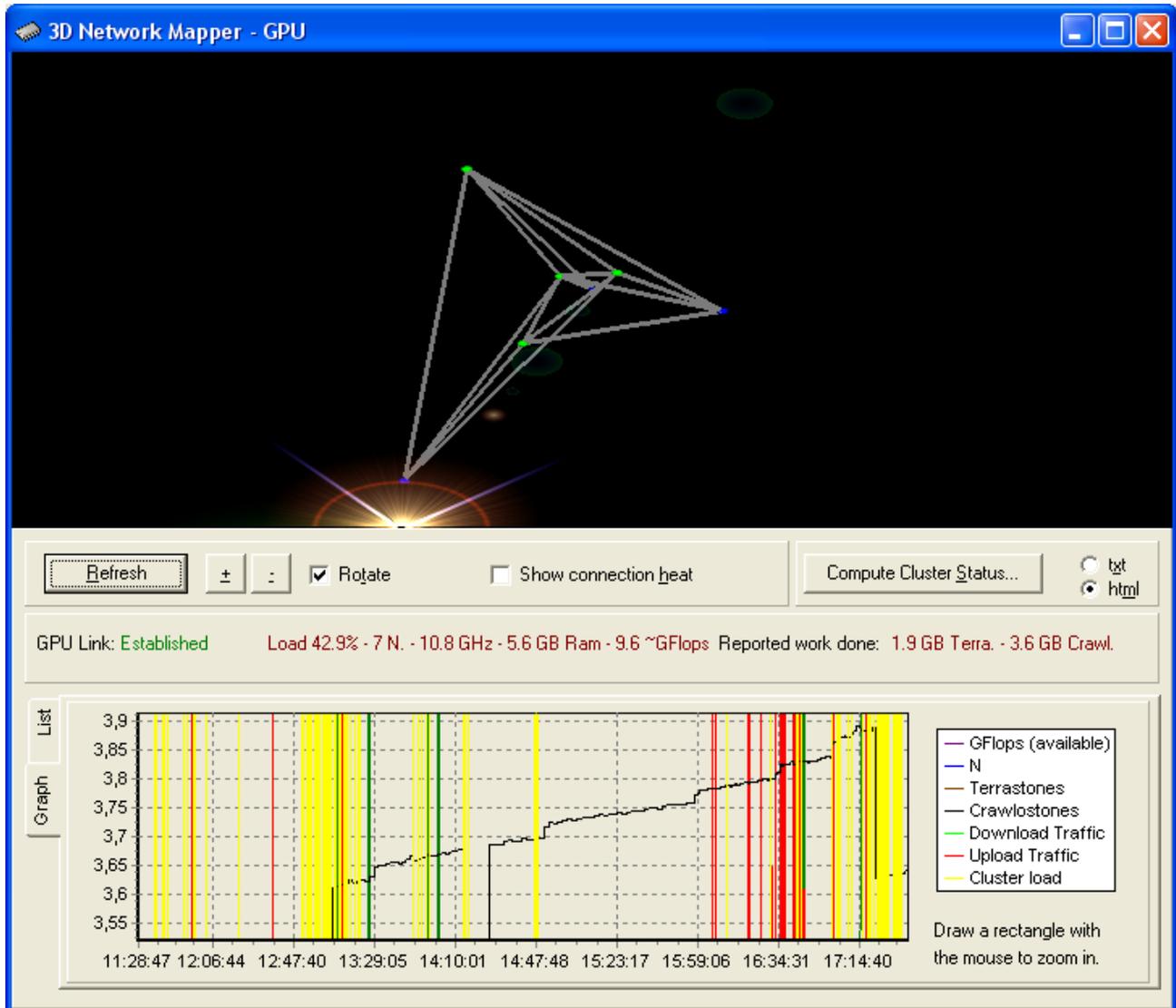


A Graphical representation:



The crawler's progress

If we zoom in on the crawler line (the black one) we can see the database slowly growing:



Concluding

Although not a scientific project, we believe the GPU Search Engine is a working proof of concept. If this network would scale to tens or even hundreds of crawlers, we can expect excellent search results and a large portion of the web crawled. Improvements may be made to sorting algorithms, but this KISS concept also has advantages. For many queries, page relevance sorting shows much similarities to other search engines, like google.

Many thanks to all that keep the GPU network alive!