# Crawlzilla - A Toolkit for Deploying Cluster Search Engine Quickly and Easily

**Shun-Fa Yang、Wei-Yu Chen、Wen-Chieh Kuo**

**Free Software Lab.@ NCHC**
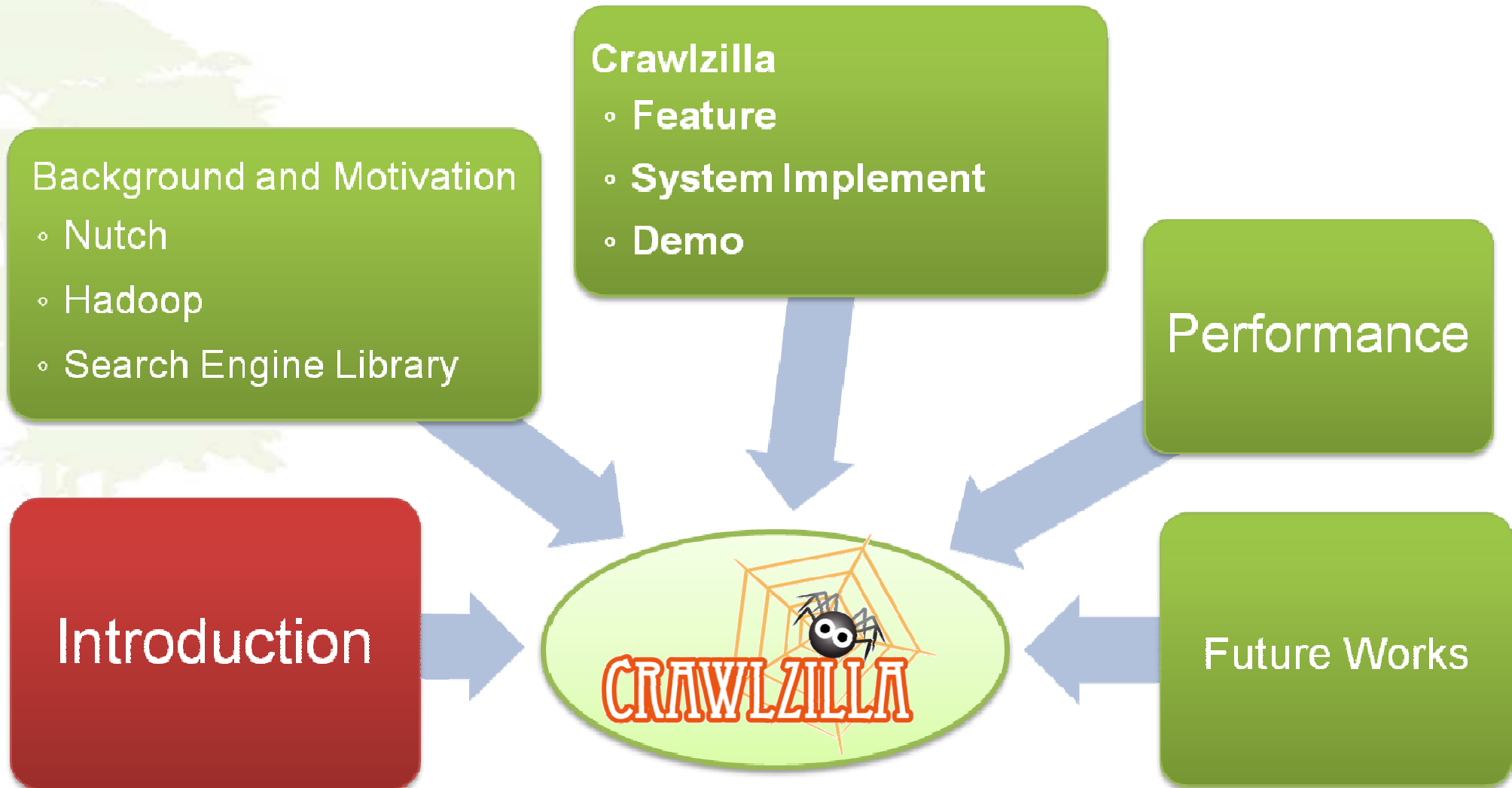
INVENSIVE 2011 May 23, 2011

TAIWAN

www.nchc.org.tw

National Applied Research Laboratories

NCHC

# Outline

Background and Motivation
- Nutch
- Hadoop
- Search Engine Library

Crawlzilla
- Feature
- System Implement
- Demo

Performance
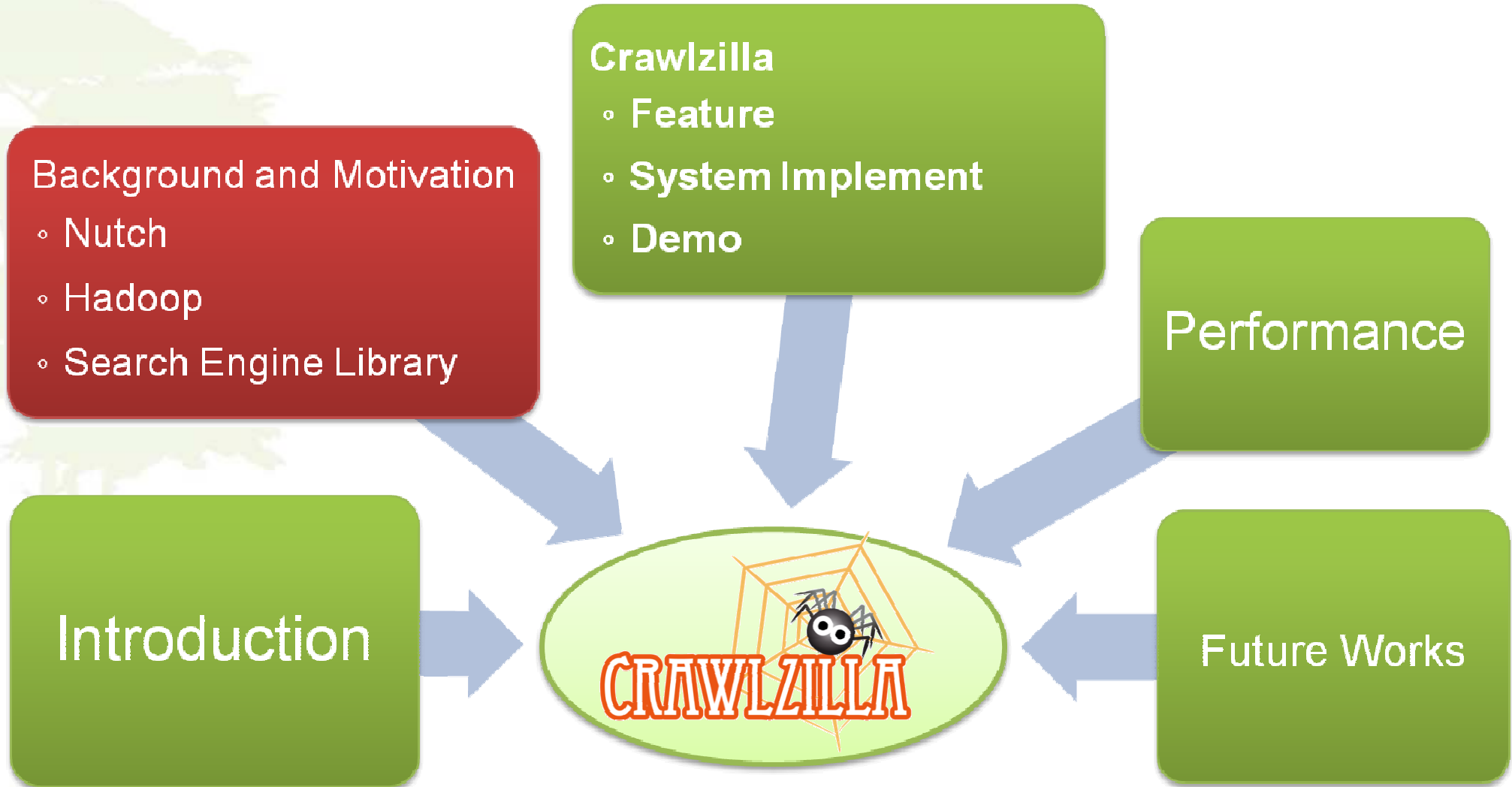
Introduction

Future Works

# Introduction

- **The Information Explosion**

- **Increase Filter Efficiency by Search Engines**

- **Intranet also need Search Engines**

- **Build Search Engines isn't very Easy**

- **Crawlzilla can help You!**

# Outline

**Background and Motivation**
- Nutch
- Hadoop
- Search Engine Library

**Crawlzilla**
- Feature
- System Implement
- Demo

**Performance**

**Introduction**

**Future Works**

# Background and Motivation

**Search Engine workflow**

**Related Open Source Projects**

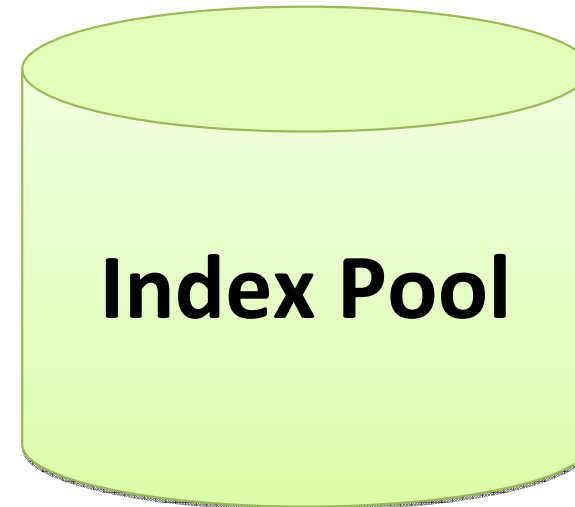**Compare with Other Projects**
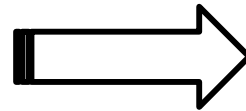
# Search Engine workflow – Phase 1

- **Crawling the Web**

List of Links → [spider/crawler image] → Page Contents

**Crawler visits the web pages of the links**

# Search Engine workflow – Phase 2

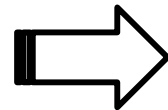- ## Building the Index Pool

Page Contents → Index Pool

Parse Contents

- **Serving Queries**



**User Sent a Query**    **Search from Index Pool**

# Background and Motivation

- **Related Open Source Projects**
  - Search Engine - Ntuch
  - Distributed Computing Platform – Hadoop
  - Search Engine Library – Lucene

# Background and Motivation

- **If Build Search by Yourself …**
  - Setup Hadoop
  - Deploy  System Configure Files
  - Debug Errors…
  - …
  - …
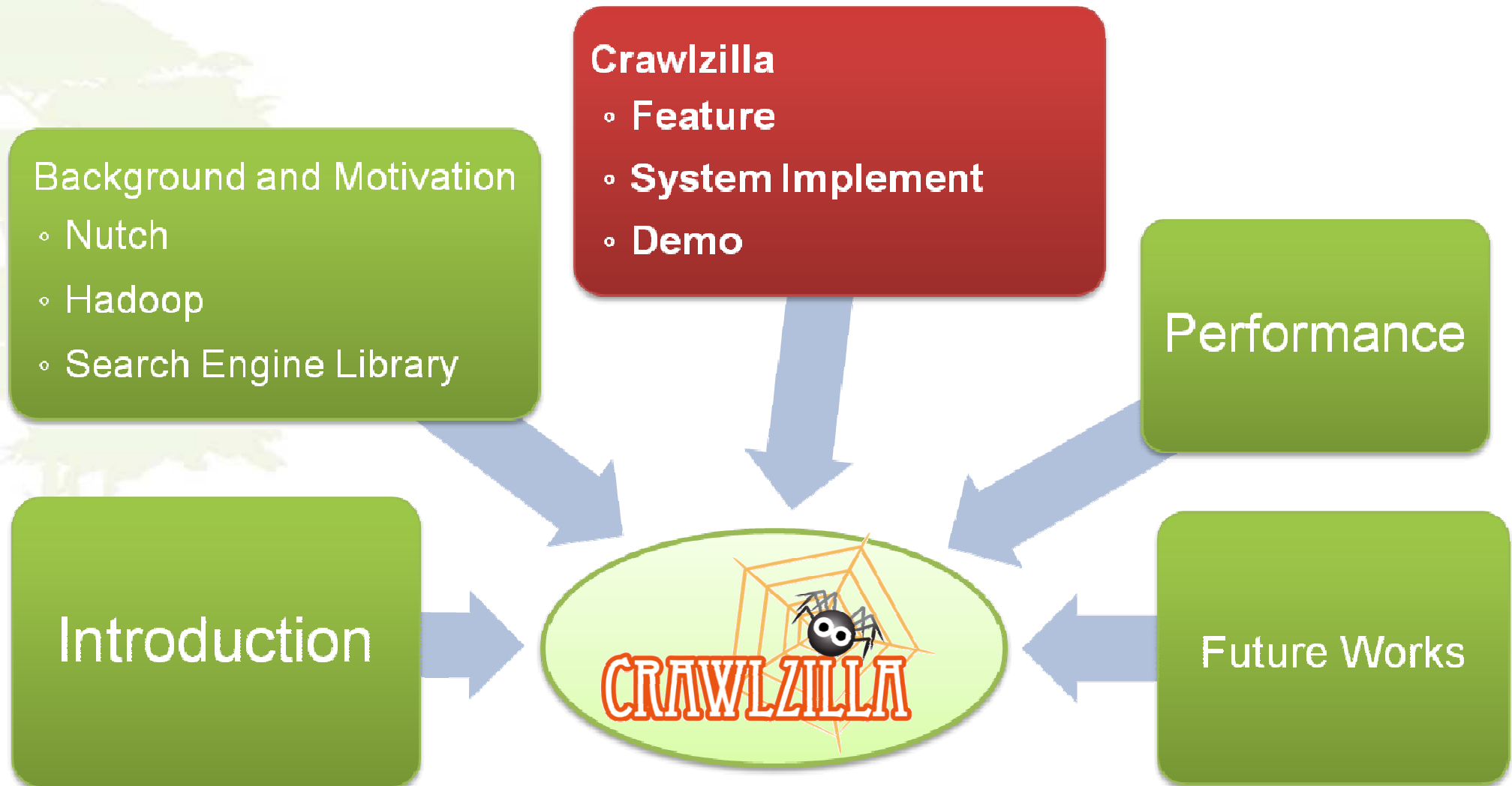  - …

www.nchc.org.tw

National Applied
Research Laboratories

# Compare with Other Projects

| | Spidr | Larbin | Jcrawl | Nutch | Crawlzilla |
|---|---|---|---|---|---|
| Install | Rube Package Install | Gmake Compiler and Install | Java Compiler and Install | Deploy Configure Files | **Provide Auto Installation** |
| Crawl website pages | O | O | O | O | **O** |
| Parser Content | X | X | X | O | **O** |
| Cluster Computing | X | X | X | O | **O** |
| Interface | Command | Command | Command | Command | **Web-UI** |
| Support Chinese Segmentation | X | X | X | X | **O** |

# Goal

- **To Help Users to Build Search Engines Easily!**

- **To Help Users to Operate System Easily!**

- **Crawlzilla doesn't improve the algorithm of Nutch and Hadoop!**

- **Crawlzilla Provides Friendly Operating Interface and an Easy Way to Deploy Cluster Computing Environment!**
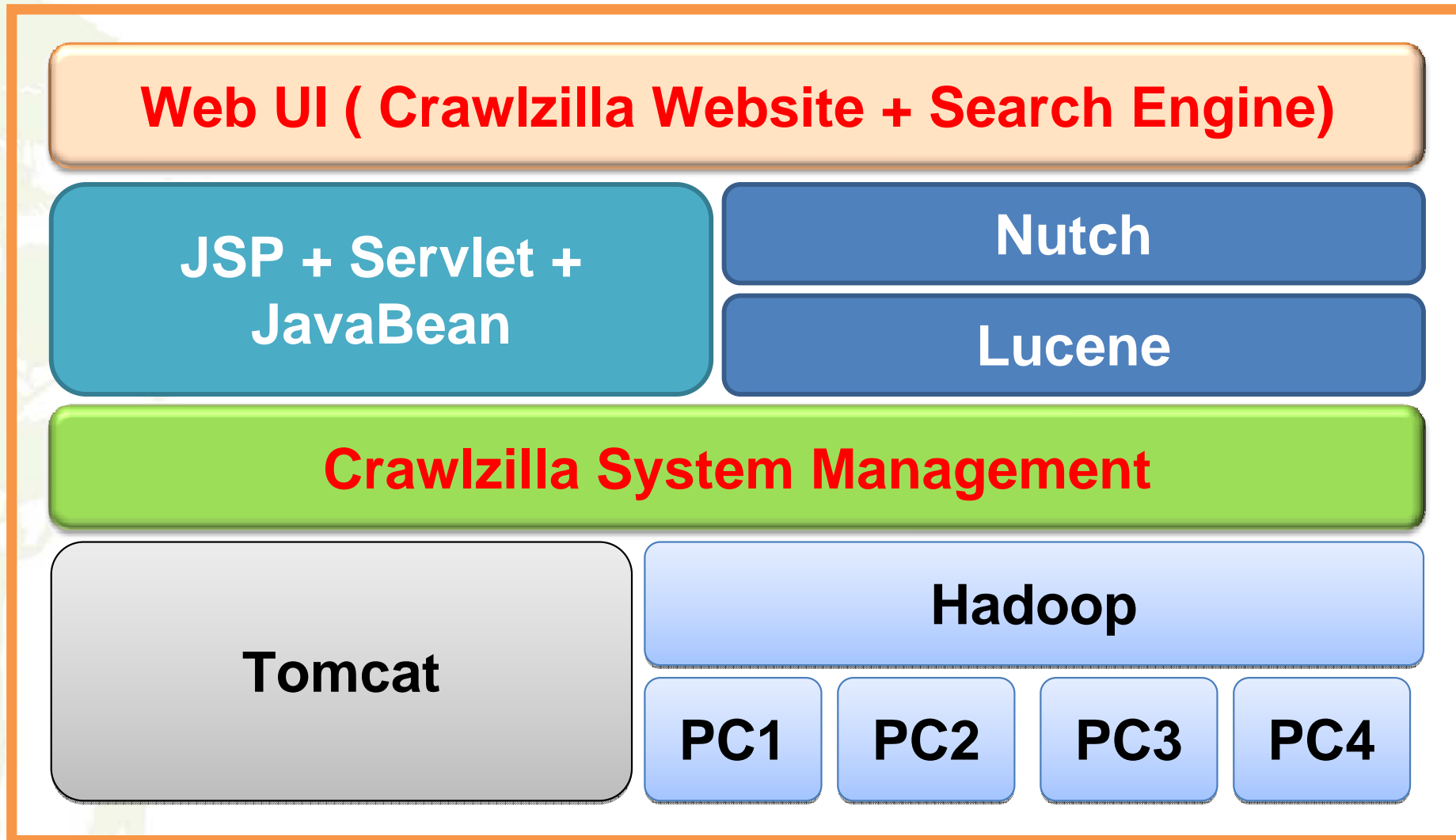
# Outline

**Background and Motivation**
- Nutch
- Hadoop
- Search Engine Library

**Crawlzilla**
- **Feature**
- **System Implement**
- **Demo**

**Performance**

**Introduction**

**Future Works**

# Crawlzilla Feature

- **Simply Install and Easy to Operate**

  – Customize user interface

- **More Powerful**

  – Support multiple search engines

- **More Search Engine Info.**

- **Developers to focus more**

  – Data mining tools

www.nchc.org.tw
National Applied
Research Laboratories

# Crawlzilla Architecture

**Web UI ( Crawlzilla Website + Search Engine)**

**JSP + Servlet + JavaBean**

**Nutch**

**Lucene**

**Crawlzilla System Management**

**Tomcat**

**Hadoop**

**PC1** **PC2** **PC3** **PC4**

# System Implement

| | JSP | Shell Script |
|---|---|---|
| Function | User UI | Admin and MIS UI |
| Security | Website Session | Crawler password with RSA Keys |
| Environment | Browser | Terminal with SSH –Client |
| Architecture | MVC | Module |
| Multi Language | i18n | Language parameters |
| | Default language is depend on O.S. Env. | |

# Web Management

(Model 2)

**Setup PW**

**Capture Lucene index pool**

**Setup and Drive Crawl Procedure**

**Servlet**

**Control**

**Model**

**JavaBean**

**View**

**JSP**

**System Status**

**Session Certification**

**i18N language setup**

# System Implement – Web Parser

crawlzilla

**Crawl Page**

http://localhost:8080/crawl.html

URL

Depth 1
2
3
...

submit

**Servlet**

**Nutch**

**Java Beans**

**JSP**

**Lucene Index DB**

**Hadoop**

**JobTracker**

**NameNode**

T  D  T  D  T  D  T  D  T  D  T  D  T  D

**T** TaskTracker: Job Executor

**D** DataNode: Data Storage Node

索引庫管理

| 索引庫名稱 | 建立時間 | 刪除索引庫 | 預覽統計資料 | 嵌入搜尋引擎到網頁的語法 |
|---|---|---|---|---|
| nchc-en_3 | 2010-08-24_16:16:14 | Delete | Preview | embed code |
| nchc-tw_3 | 2010-08-24_15:22:48 | Delete | Preview | embed code |

資料總覽

| 起始URL | http://www.nchc.org.tw/tw/ |
|---|---|
| 本機索引路徑 | /home/crawler/crawlzilla/archieve/nchc-tw_3/index |
| 總共文字數 | 37095 | 文件檔數量 | 1036 |
| 索引庫更新日期 | Tue Aug 24 15:22:46 CST 2010 | 使用者名稱 | crawler |

被搜尋分析到的網址:

| 排序 | 內容 | 引用次數 | 排序 | 內容 | 引用次數 |
|---|---|---|---|---|---|
| 0 | site:www.nchc.org.tw | 336 | 1 | site:pccluster.nchc.org.tw | 87 |
| 2 | site:bioinfo.nchc.org.tw | 66 | 3 | site:www.narl.org.tw | 57 |
| 4 | site:edu.nchc.org.tw | 53 | 5 | site:service.nchc.org.tw | 35 |
| 6 | site:accta.nchc.org.tw | 28 | 7 | site:colife.nchc.org.tw | 14 |
| 8 | site:wlanrc.nchc.org.tw | 13 | 9 | site:elib.nchc.org.tw | 13 |
| 10 | site:www.medicalgrid.org | 13 | 11 | site:volunteer.nchc.org.tw | 9 |
| 12 | site:www.stpi.org.tw | 7 | 13 | site:noc.twaren.net | 7 |
| 14 | site:ecogrid.nchc.org.tw | 6 | 15 | site:www.sipa.gov.tw | 3 |
| 16 | site:asp.104ehr.com.tw | 3 | 17 | site:viml.nchc.org.tw | 3 |
| 18 | site:www.ym.edu.tw | 2 | 19 | site:www.tnu.edu.tw | 2 |
| 20 | site:www.usc.edu.tw | 2 | 21 | site:www.ssvs.tp.edu.tw | 2 |
| 22 | site:www.smelearning.org.tw | 2 | 23 | site:ecocam.nchc.org.tw | 2 |

搜尋引擎快速連結

CrawlZilla 搜尋引擎範例

nchc-en_3

nchc-tw_3

系統功能

修改管理者密碼

相關資源

CrawlZilla@GoogleCode

# Friendly Interface!



**Admin**

**MIS**

**USER**

# Live Demo I

## Crawlzilla Install

(1) Master Install
(2) Cluster Slave Install

www.nchc.org.tw

# Live Video Demo

- **Master Install ([Demo Video also @ YouTube](#))**

www.nchc.org.tw

National Applied
Research Laboratories

# Live Video Demo

- **Slave Install (Demo Video also @ YouTube)**

www.nchc.org.tw
National Applied
Research Laboratories

# Live Demo II

## Dialog Management

# Live Demo III
## Web Management

(1) Crawl Setup

(2)Search Engine Index Pool

(3)Search it!

# Outline

**Crawlzilla**
- Feature
- System Implement
- Demo

Background and Motivation
- Nutch
- Hadoop
- Search Engine Library

Performance

Introduction

CRAWLZILLA

Future Works

www.nchc.org.tw

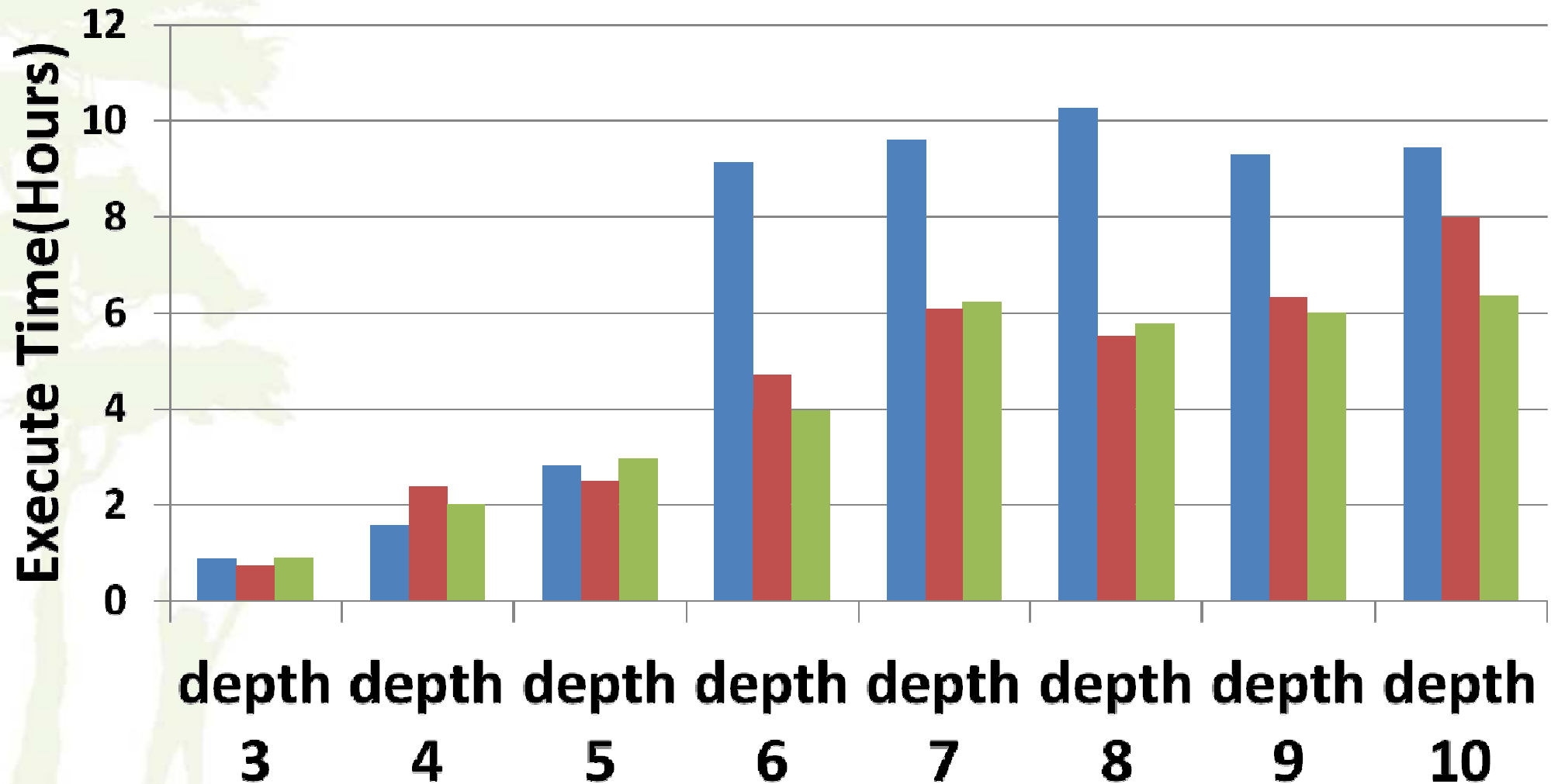National Applied
Research Laboratories

NCHC

# Performance

## Experiment Environment

- CPU
  - Intel(R) Core(TM)2 Quad CPU Q9550 2.83GHz
- Memroy
  - 8 GigaBytes
- Operation System
  - Ubuntu 10.04 Lucid(x86)
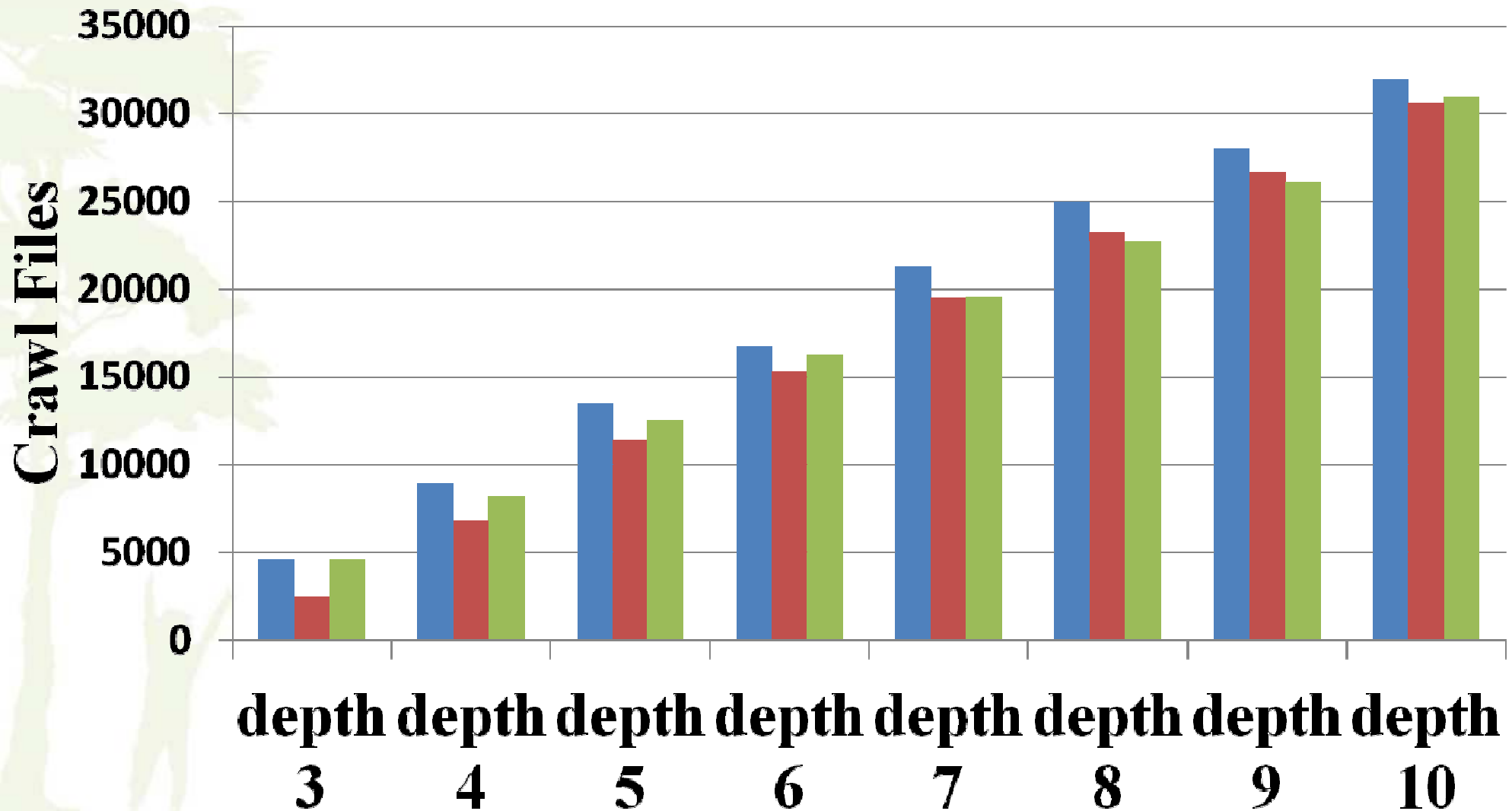- Crawlzilla Version
  - 0.3.0-101116

# Execute Time

# Crawl Words

# Outline

**Crawlzilla**
- Feature
- System Implement
- Demo

Background and Motivation
- Nutch
- Hadoop
- Search Engine Library

Performance

Introduction

Future Works

# Future Works

- **New Version**
  - Support Multi User
  - Support Schedule
  - Update the Kernel
  - More Easily to deploy Slave Computing Nodes
  - Now is testing!
  - Release Day See http://crawlzilla.info

# Reference

- J. Dean and S. Ghemawat, MapReduce: Simplified Data Processing on Large Clusters, In Proceedings of the 6th Conference on Symposium on Opearting Systems Design & Implementation - Volume 6, San Francisco, CA, December 06 - 08, 2004.

- S. Ghemawat, H. Gobioff and S. T. Leung, The Google File System, 19th ACM Symposium on Operating Systems Principles, Lake George, NY, October, 2003.

- The Apache Software Foundation, Nutch, available at: http://nutch.apache.org/ , accessed 5 June 2010.

- The Apache Software Foundation, Hadoop, available at: http://hadoop.apache.org/ , accessed 5 June 2010.

- The Apache Software Foundation, Lucene, available at: http://lucene.apache.org/ , accessed 5 June 2010.

- Crawlzilla @ Google Code Project Hosting, available at: http://code.google.com/p/crawlzilla/, accessed 15 Sep 2010.

# Enjoy your search engines!!! Start from Here!

- **Crawlzilla @ Google Code Project Hosting (Tutorials in Chinese)**

  - **http://code.google.com/p/crawlzilla/**

- **Crawlzilla @ Source Forge (Tutorials in English)**

  - **http://sourceforge.net/p/crawlzilla/home/**

- **Crawlzilla User Group @ Google**

  - **http://groups.google.com/group/crawlzilla-user**

- **NCHC Cloud Computing Research Group**

  - **http://trac.nchc.org.tw/cloud**

# Thank You!
# Q & A

www.nchc.org.tw