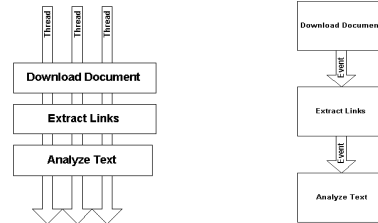# CrawlBuddy

The web's best friend

---

# Design Decisions

- Two paradigms of design to choose from: multi-threaded or event-driven
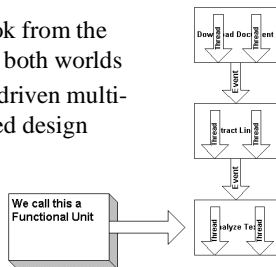


---

# Multi-Threaded Programming

- **Advantages**
- Programming is easier because threads are linear and we (usually) think linearly
- Threads can take advantage of multiprocessors easily
- Threads are synchronous i.e. it is okay for a thread to block because there are many of them running at once
- Debugging a threaded program is considerably easier than an event based program
- **Disadvantages**
- Threads are limited by the underlying operating system (operating systems can only efficiently handle so many threads)

---

# Event-Driven Programming

- **Advantages**
- Handles well under heavy load, the queues act as a buffer to soften the load
- Simple to add new functionality and process in parallel
- Easy to split up and run on multiple machines
- Modular
- **Disadvantages**
- Not as intuitive as Thread programming
- Harder to debug system level errors (but easier to debug individual pieces)

---

# What CrawlBuddy Does

- We took from the best of both worlds
- Event-driven multi-threaded design
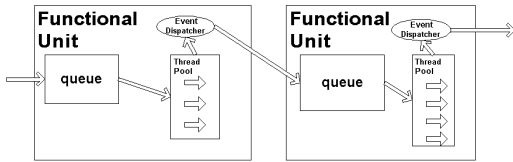


We call this a Functional Unit

---

# Functional Units Are Our Friends

- Each Functional Unit has a …
- Queue – holds events to be processed
- Thread Pool – takes events off the queue and processes them
- Event Dispatcher – sends events to other Functional Units
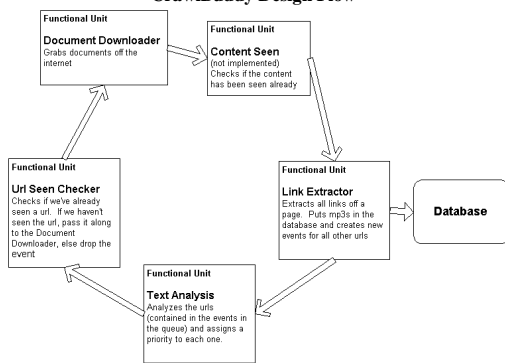
## Design of a Functional Unit

- Arrows represent flow of a task



## CrawlBuddy Design

- Basically, events are passed between Functional Units
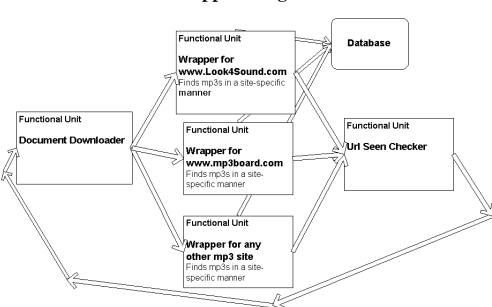- The arrows (on the next slide) represent event flow

## CrawlBuddy Design Flow



**Document Downloader** — Grabs documents off the internet

**Content Seen** (not implemented) — Checks if the content has been seen already

**Url Seen Checker** — Checks if we've already seen a url. If we haven't seen the url, pass it along to the Document Downloader, else drop the event

**Link Extractor** — Extracts all links off a page. Puts mp3s in the database and creates new events for all other urls

**Text Analysis** — Analyzes the urls (contained in the events in the queue) and assigns a priority to each one.

**Database**

## Wrapper Design

- Our wrapper crawler targets specific sites and uses site-specific format to find mp3s and record information about them (song name, artist name, etc)
- The wrapper Functional Units can be run in parallel and the each use the same database
- The Document Downloader passes each event to each of the wrappers. If the event does not apply to the wrapper (i.e. the document comes from a different site), the wrapper will simply drop the event

## Wrapper Design Flow



**Wrapper for www.Look4Sound.com** — Finds mp3 in a site-specific manner

**Wrapper for www.mp3board.com** — Finds mp3s in a site-specific manner

**Wrapper for any other mp3 site** — Finds mp3s in a site-specific manner

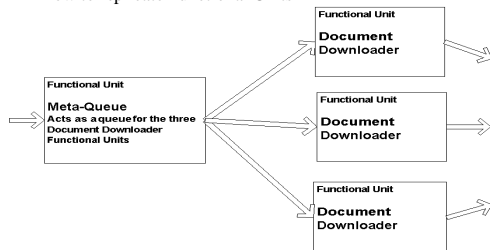**Document Downloader**

**Url Seen Checker**

**Database**

## Design Advantages

- Code re-use (Functional Units shared across CrawlBuddy and the wrapper)
- Expandable
- Checkpointing is simple (save the queues)
- Easy to run on multiple machines
- Queues buffer the load on threads
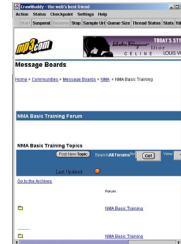- Functional Units Replicable (see next slide)

## Meta Queue
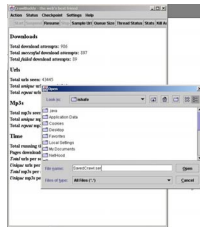
- How to replicate Functional Units

**Functional Unit**
Meta-Queue
Acts as a queue for the three Document Downloader Functional Units

**Functional Unit**
Document Downloader

**Functional Unit**
Document Downloader

**Functional Unit**
Document Downloader

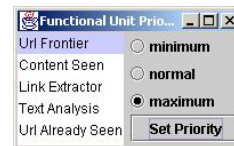

---

## CrawlBuddy Features

- GUI



---

## CrawlBuddy Features (cont)

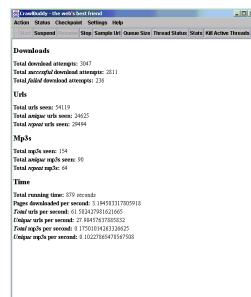- Checkpointing



---

## CrawlBuddy Features (cont)

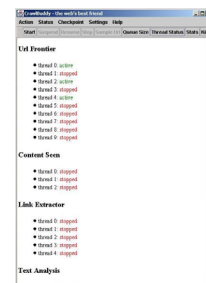- Dynamic control of Functional Unit priority



---

## CrawlBuddy Features (cont)

- Real-time stats
- Total downloads
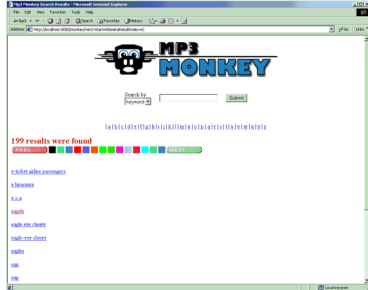- Total mp3
- Downloads / sec
- Etc.



---

## CrawlBuddy Features (cont)

- Thread status monitor

## Mp3Monkey

- Search for all 'e' artists



## Mp3Monkey Features

- Self-maintaining database – if a user attempts to download non-existent mp3, that url is marked for deletion
- Statistics are kept of how many searches and what has been downloaded