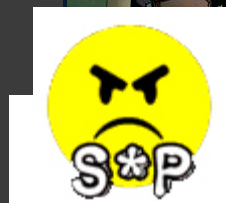
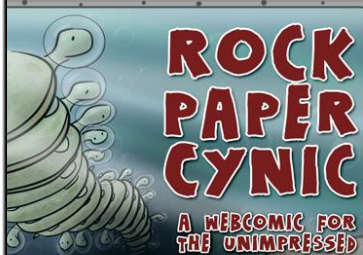
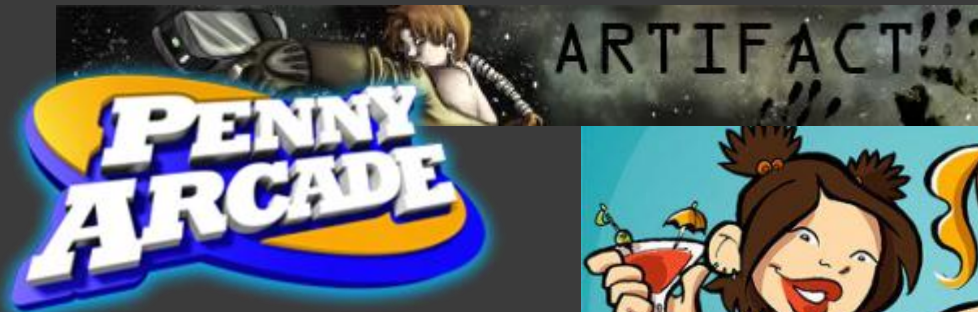


Suzanna Kangas  
Daniel Leblanc  
Lucas Berge



Fredrik Fostvedt  
Joseph Schutz  
Eric Smith

**COMIC ROCKET™**





# Comic Rocket

- ◎ Over 12,000 titles indexed
  - 4.1% more titles per month
- ◎ More than 1.5 million pages
  - 5.3% more pages per month
- ◎ Unique visitors climbing by 30% per month





# Comic Rocket

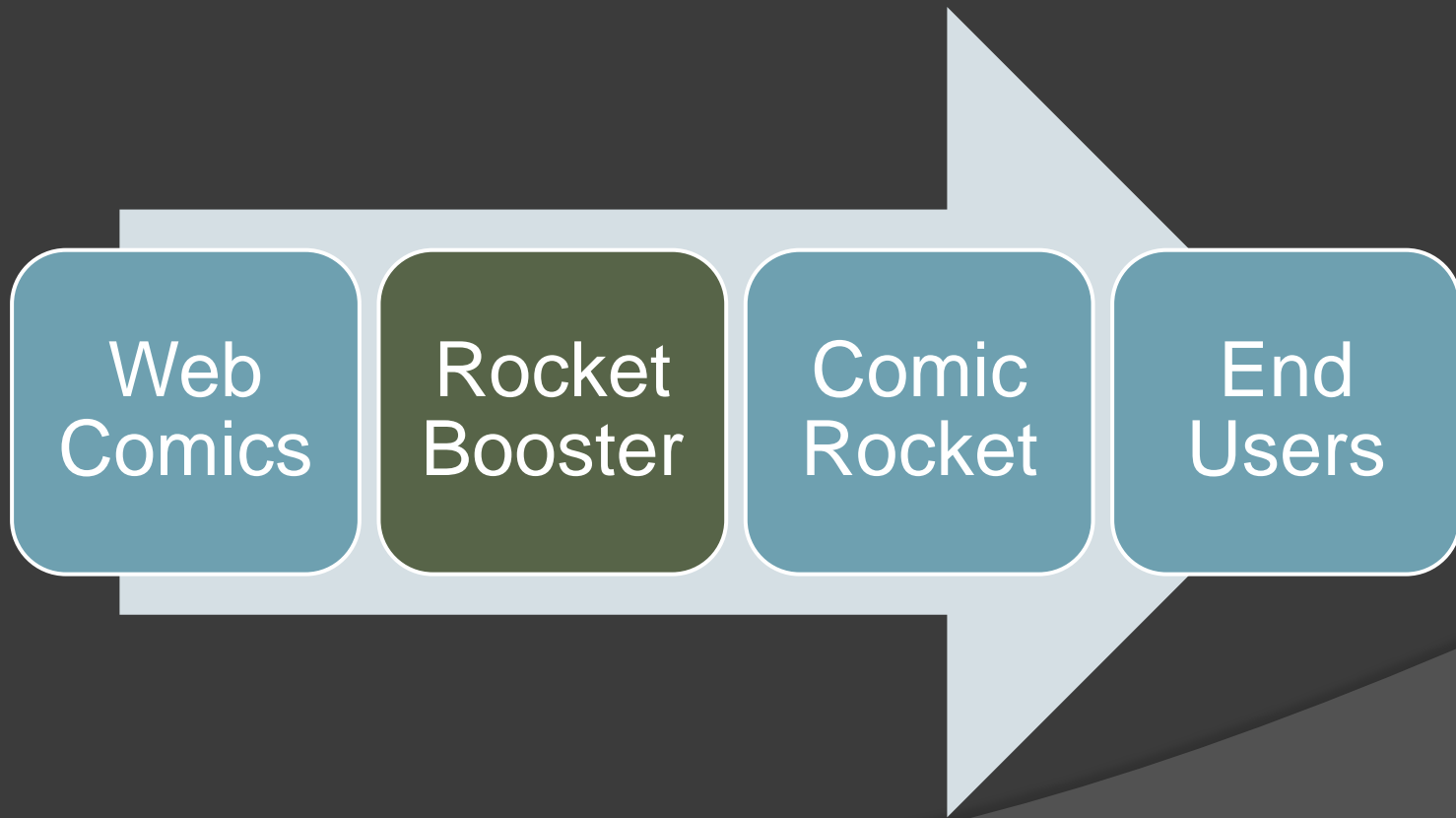
Web  
Comics

Comic  
Rocket

End  
Users



# Rocket Booster





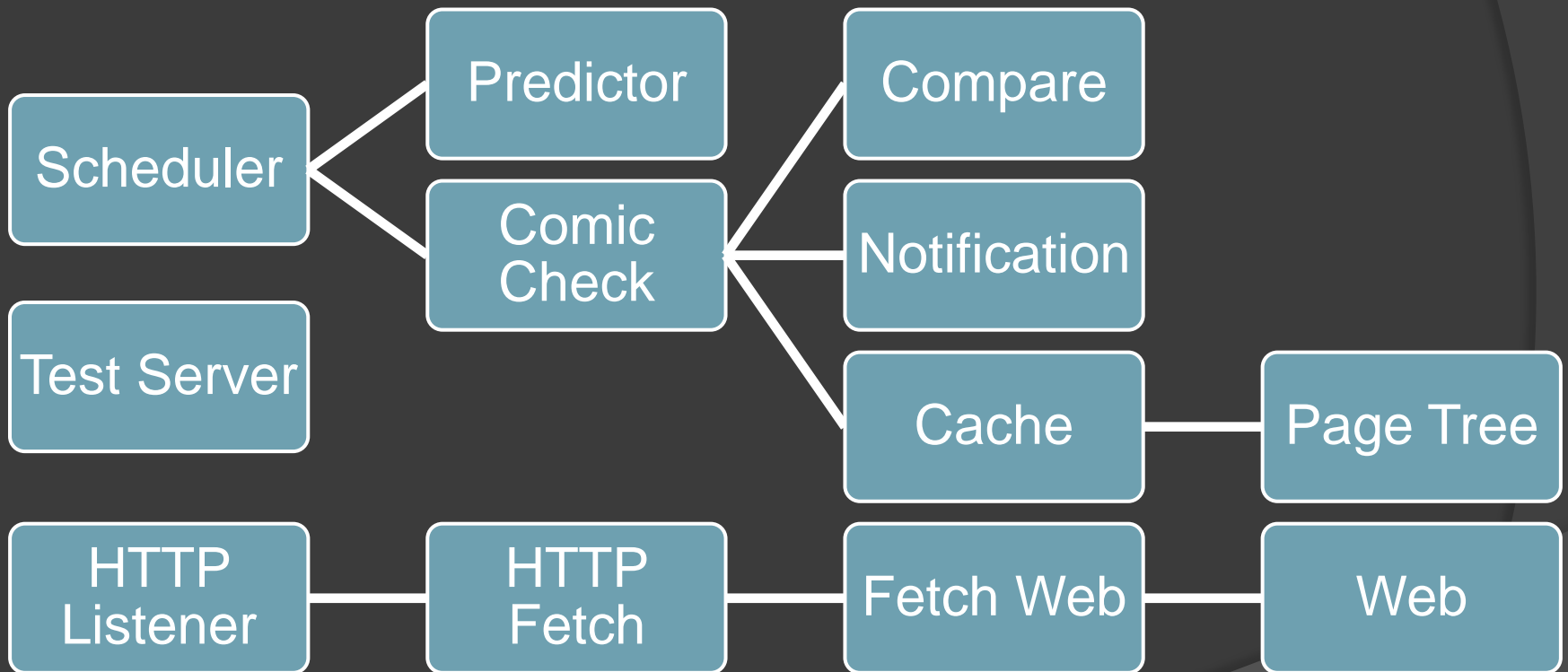
# Rocket Booster

- ◉ Manage Comic Rocket's interaction with the web
- ◉ Cache all pages
- ◉ Check for new comics
- ◉ Check for updates to old comics
- ◉ Limit number of checks



COMIC ROCKET™

# Rocket Booster Architecture





COMIC ROCKET™

# Planned Timeline

## Planning

Oct 15, 2012 –  
Nov 19, 2012

Initial design  
Project analysis  
Task breakdown  
Risk  
assessment

## Initial setup

Nov 19, 2012 –  
Nov 26, 2012

Git repository  
Coding  
standards  
Function names

## Coding and Unit Testing

Nov 26, 2012 –  
Jan 28, 2013

Individual  
coding  
Handling issues  
Testing code  
Scripted Tests

## Project Assembly

Jan 28, 2013 –  
Feb 11, 2013

Combining units  
into a whole  
Dealing with  
issues that arise

## Project Testing

Feb 11, 2013 –  
Mar 11, 2013

Test Server  
Cache  
compression  
Additional  
functionality



# Actual Timeline



COMIC ROCKET™

## Planning

Oct 15, 2012 –  
Nov 19, 2012

Initial design  
Project analysis  
Task breakdown  
Risk  
assessment

## Initial setup

Nov 19, 2012 –  
Nov 26, 2012

Git repository  
Coding  
standards  
Function names

## Coding and Unit Testing

Nov 26, 2012 –  
Feb 11, 2013

Individual  
coding  
Handling issues  
Testing code  
Scripted Tests

## Project Assembly

Feb 11, 2013 –  
Feb 25, 2013

Combining units  
into a whole  
Dealing with  
issues that arise

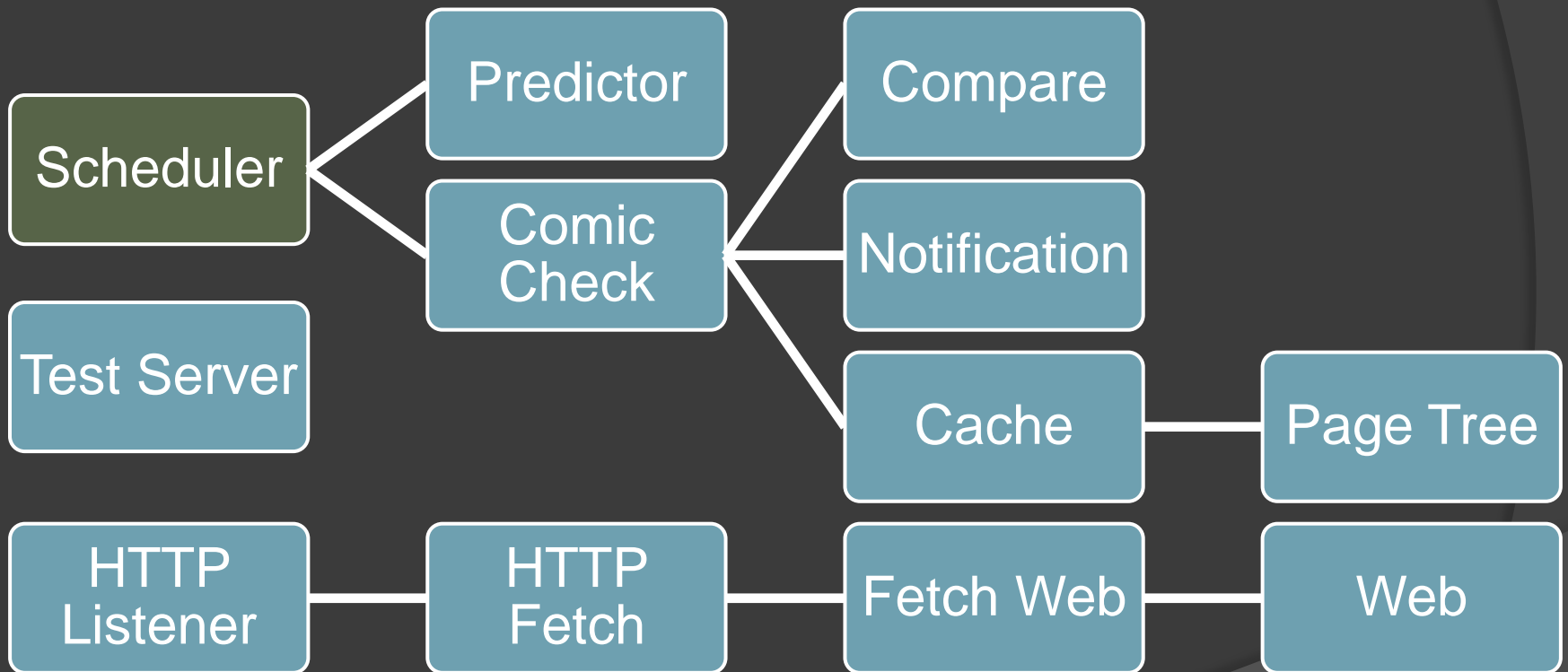
## Project Testing

Feb 25, 2013 –  
Mar 13, 2013

Time  
compression



# Daniel Leblanc





# Scheduler

- ⦿ Keeps track of history checking
  - Continuous loop of all comics
  - Checks at most 20 per hour per domain
- ⦿ Checks the list of comics provided by the Predictor every hour
- ⦿ Getting it to terminate was a challenge



# HTTP Requester

- ⦿ Pretends to be Jamey
- ⦿ Receives notifications and responds with HTTP requests
- ⦿ Keeps track of bad notifications and missed updates



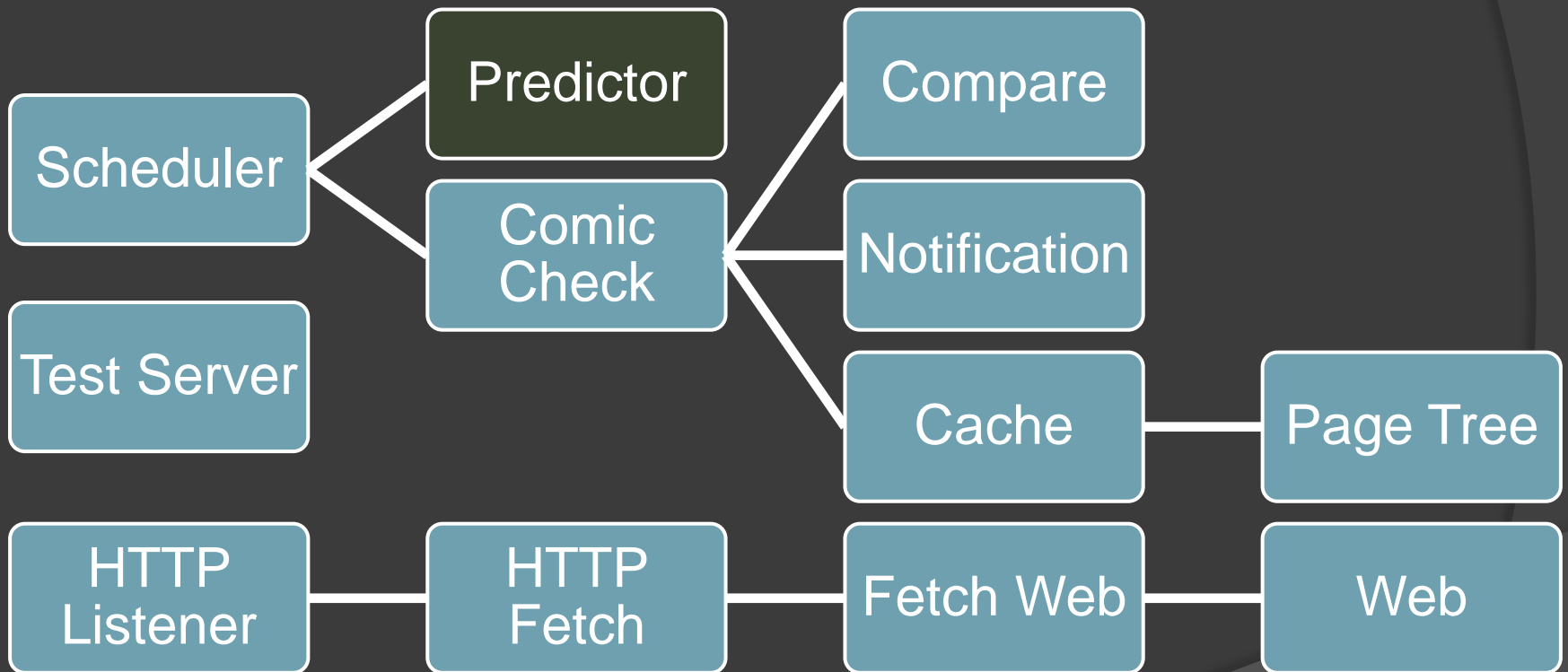
# What I Learned

- ⦿ Mitigating risks
- ⦿ Explaining a technical idea can be challenging





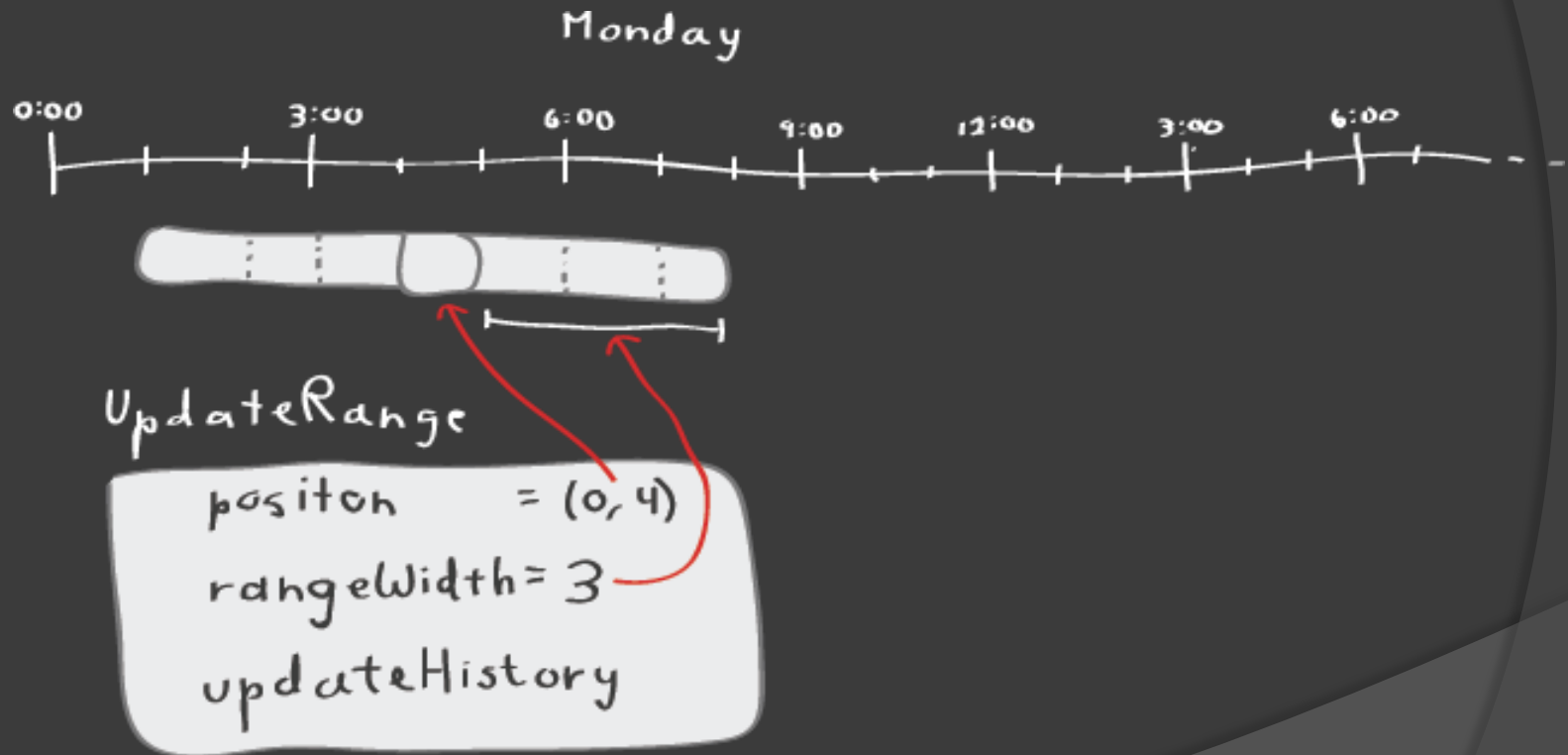
# Fredrik Fostvedt





COMIC ROCKET™

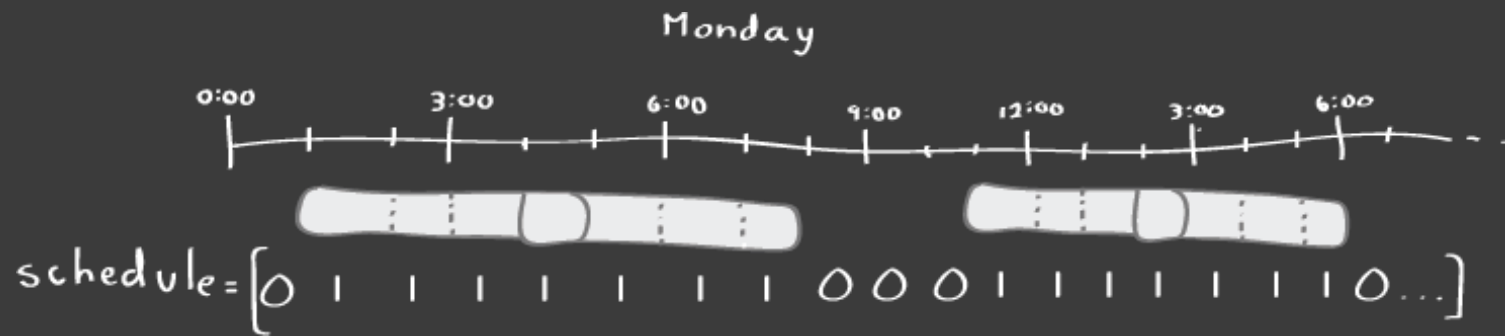
# Predictor Structure





COMIC ROCKET™

# Predictor Structure





# Predictor Building

Mon Tu We Th Fr Sa Su



schedule = [ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | ]

change

found

uRange





COMIC ROCKET™

# Predictor Learning

UpdateRange

updateHistory = [ (day, hour), ... ]

position = average(updateHistory)





COMIC ROCKET™

# Predictor Learning

Monday

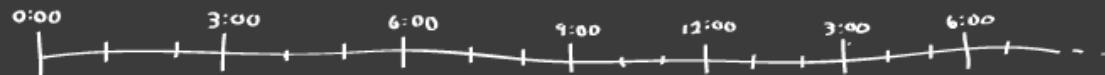


schedule = [0 1 1 1 1 1 1 0 0 0 0 0 ...]

updateHistory



Monday



schedule = [0 0 1 1 1 1 1 1 0 0 0 0 ...]

updateHistory





# Predictor Scheduling

Monday



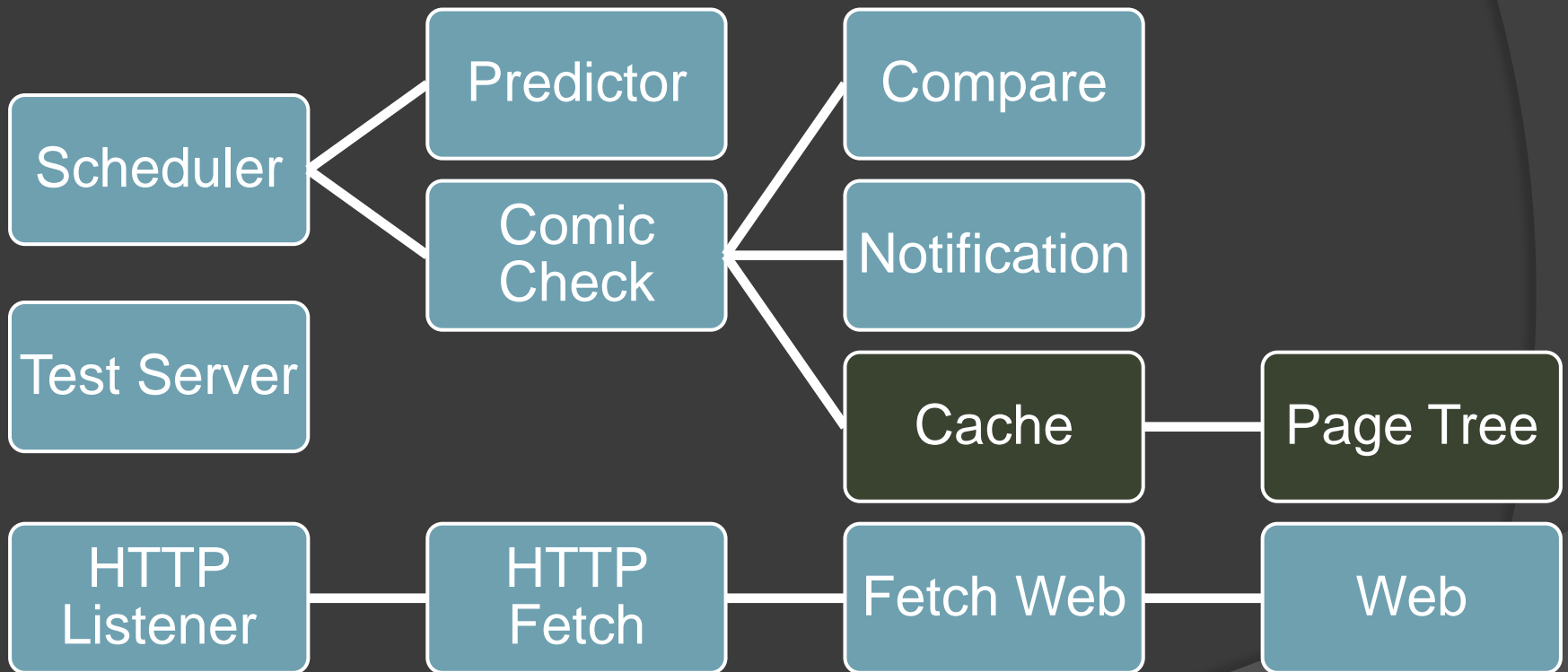
comic 1  
schedule  $[0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \dots]$

comic 2  
schedule  $[0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 0 \ 0 \dots]$

comic list  $(\square) (\square) (\square) [1,2] [1,2] [1,2] (\square) (\square) (\square) (\square) \dots$



# Suzanna Kangas





# Cache

- ◉ Stores page versions for comparison
- ◉ Retrieve any past version
- ◉ Revision pushing
- ◉ File structure of storage
- ◉ No compression currently



# Page Tree

- ⦿ Internal representation of a web page
- ⦿ Readily extensible



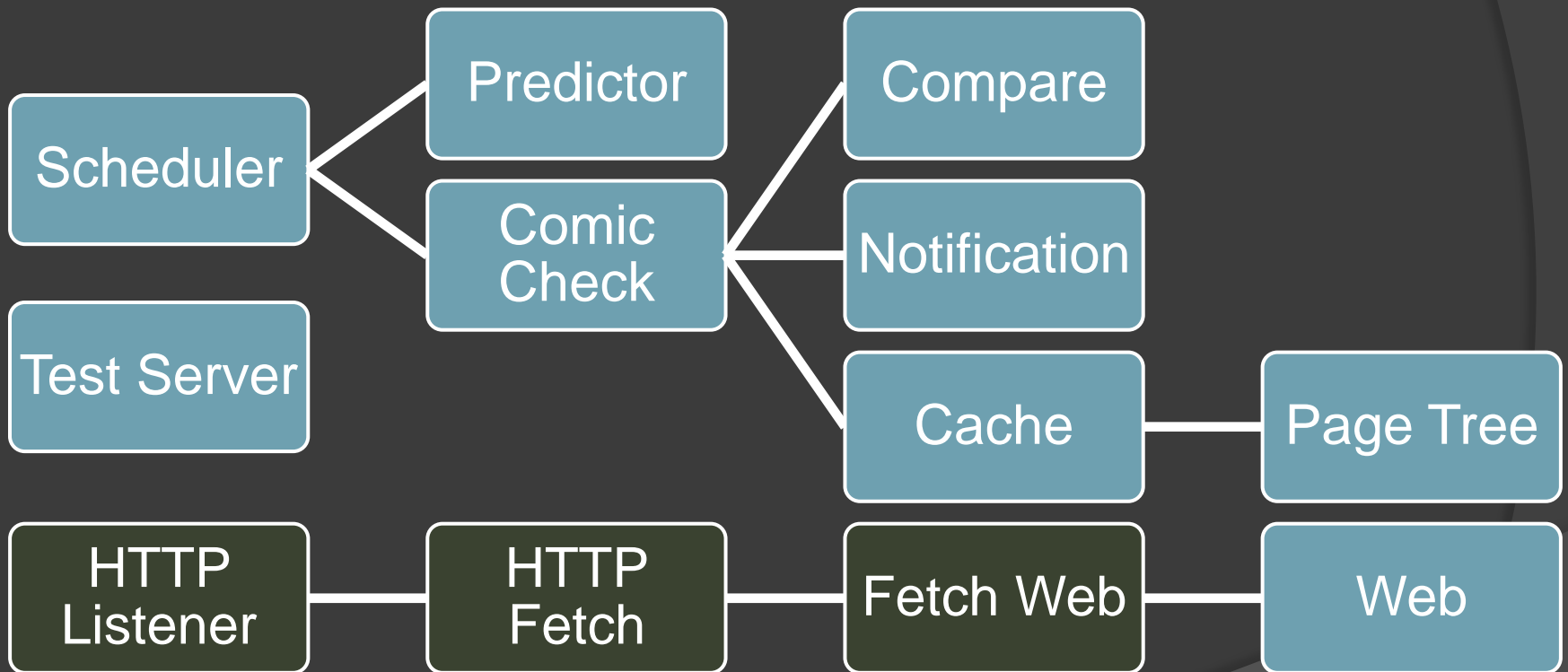


# What I Learned

- ◎ Python
- ◎ Working with a team
- ◎ Got a glimpse of how important testing can be



# Lucas Berge





# HTTP Listener

## ⦿ What it is

- Modified SimpleHTTPRequestHandler

## ⦿ What it does

- Runs on a separate thread waiting for requests
- Searches Cache for comicID, gets from Web if not found
- Allows a user to send commands to predictor
- Internal notifications



# Fetch HTTP

## ⦿ What it is

- Checks cache for a page, requests from Web if not found

## ⦿ What it does

- Asks for a pointer to the object, if none call FetchWeb



# Fetch Web

## ◉ What it is

- HTTP Client that constructs pageTree objects and stores them in the cache

## ◉ What it does

- Requests objects with urllib2
- Parses out all links with BS4
- Constructs a pageTree object and caches it with ComicID key
- Internal notifications





# What I Learned

- ◉ GitHub basics
  - Not the most user friendly
- ◉ Python
  - Great language
- ◉ Libraries
  - More useful when you know how they work



# GitHub

- ⦿ Commit Often!
  - Even if you don't like what you have
  - Conflicts
- ⦿ Stash and pull
- ⦿ Use the shell
- ⦿ Don't be afraid to break everything



# Python

- ◉ No prior experience
  - Easy to understand syntax
  - Lots of support and documentation
  - Human readable libraries!
- ◉ 2.7 vs 3.0
  - Urllibs
  - Whitespace
- ◉ My new language of choice

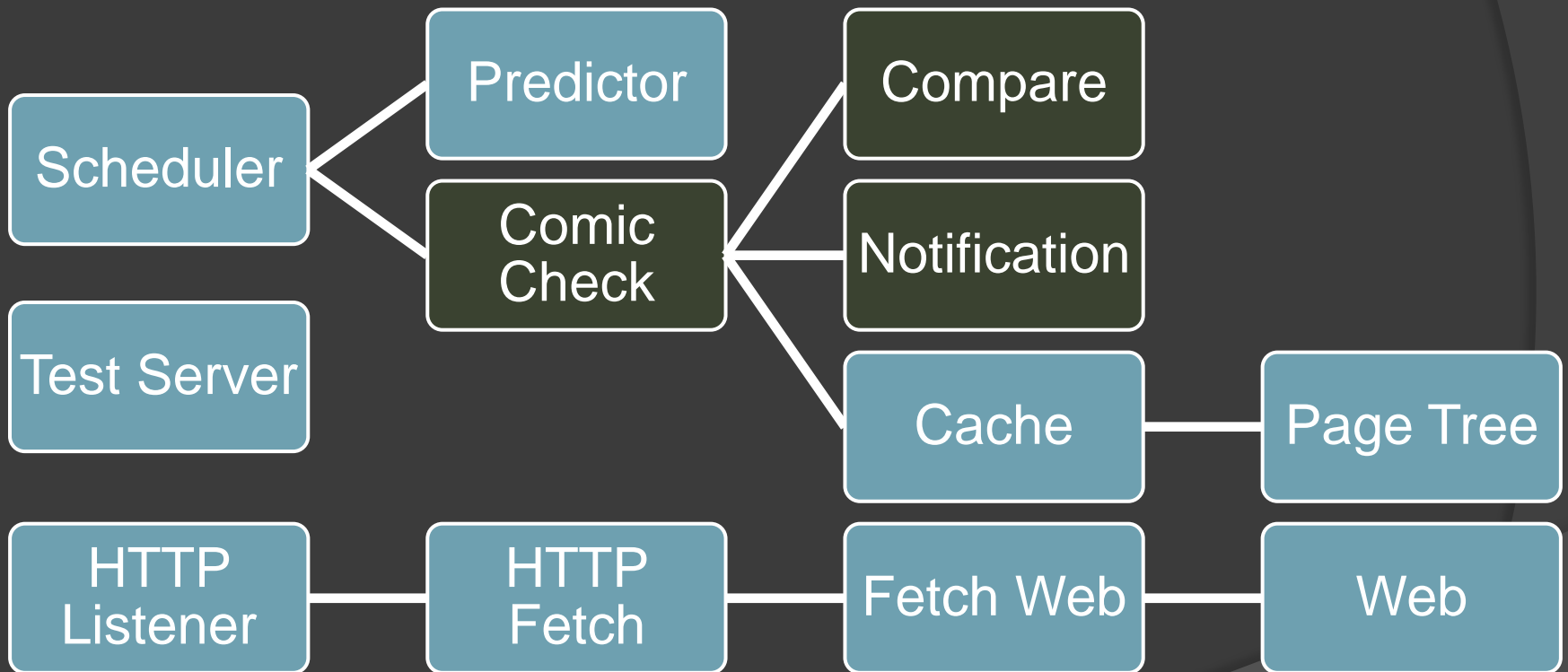


# Libraries

- ◉ Urllib2
  - ◉ BaseHTTPServer
  - ◉ BeautifulSoup
  - ◉ Threading
- 
- ◉ Read Manuals, Look at examples, debug objects, delve into library code



# Eric Smith





# Compare

- ⦿ Compares two pages and reports if they are different
- ⦿ Three iterations:
  - Page Data
  - Hash value
  - Links



# Notification

- ◉ Sends a message to Jamey via Rabbitmq
- ◉ Two iterations:
  - New comic notifications only, list of notifications sent
  - General notifications (new comic, error, etc...), no list



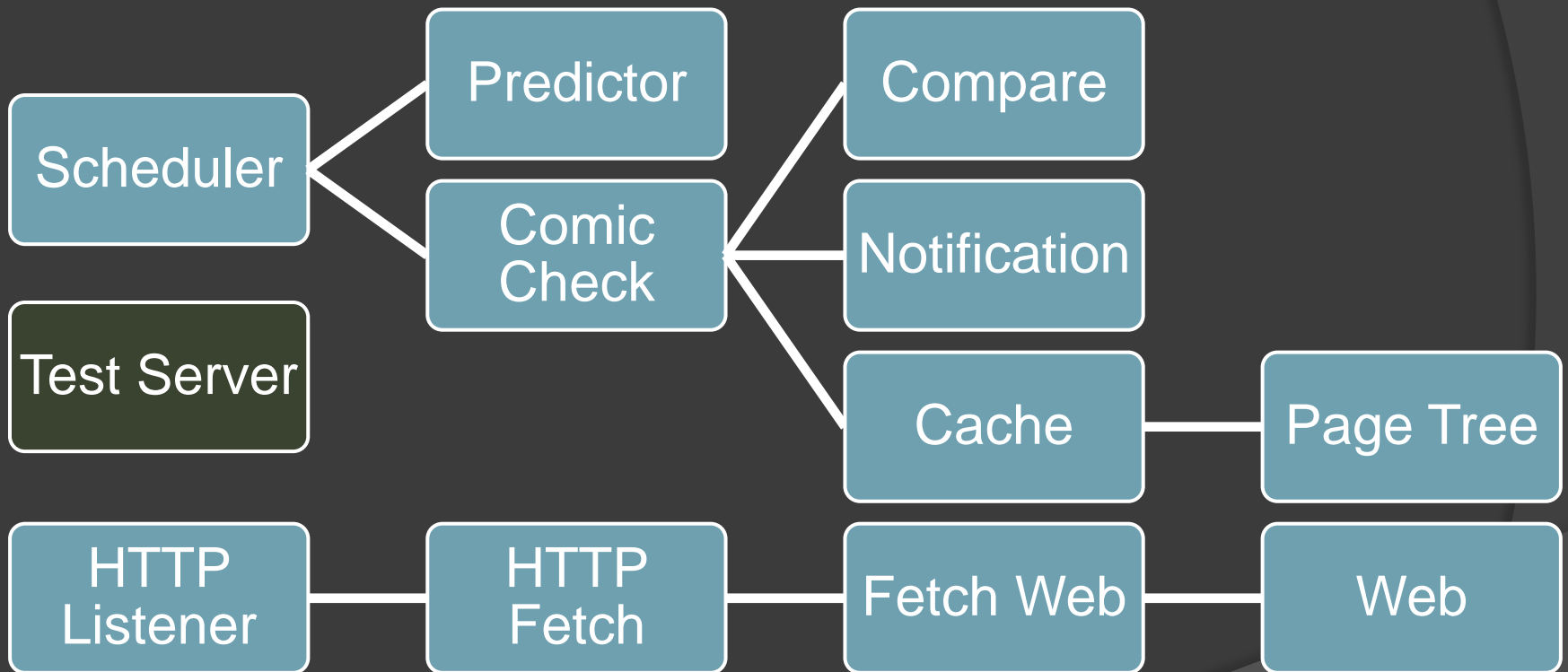
# Comic Check

- ◉ newComic & histComic
- ◉ Wrapper functions that call the other functions
- ◉ Three iterations:
  - Latest three
  - One by one
  - List based





# Joseph Schutz





# Challenges

---

Web Comics updated slowly.

---

Web Comics update in uncontrollable ways

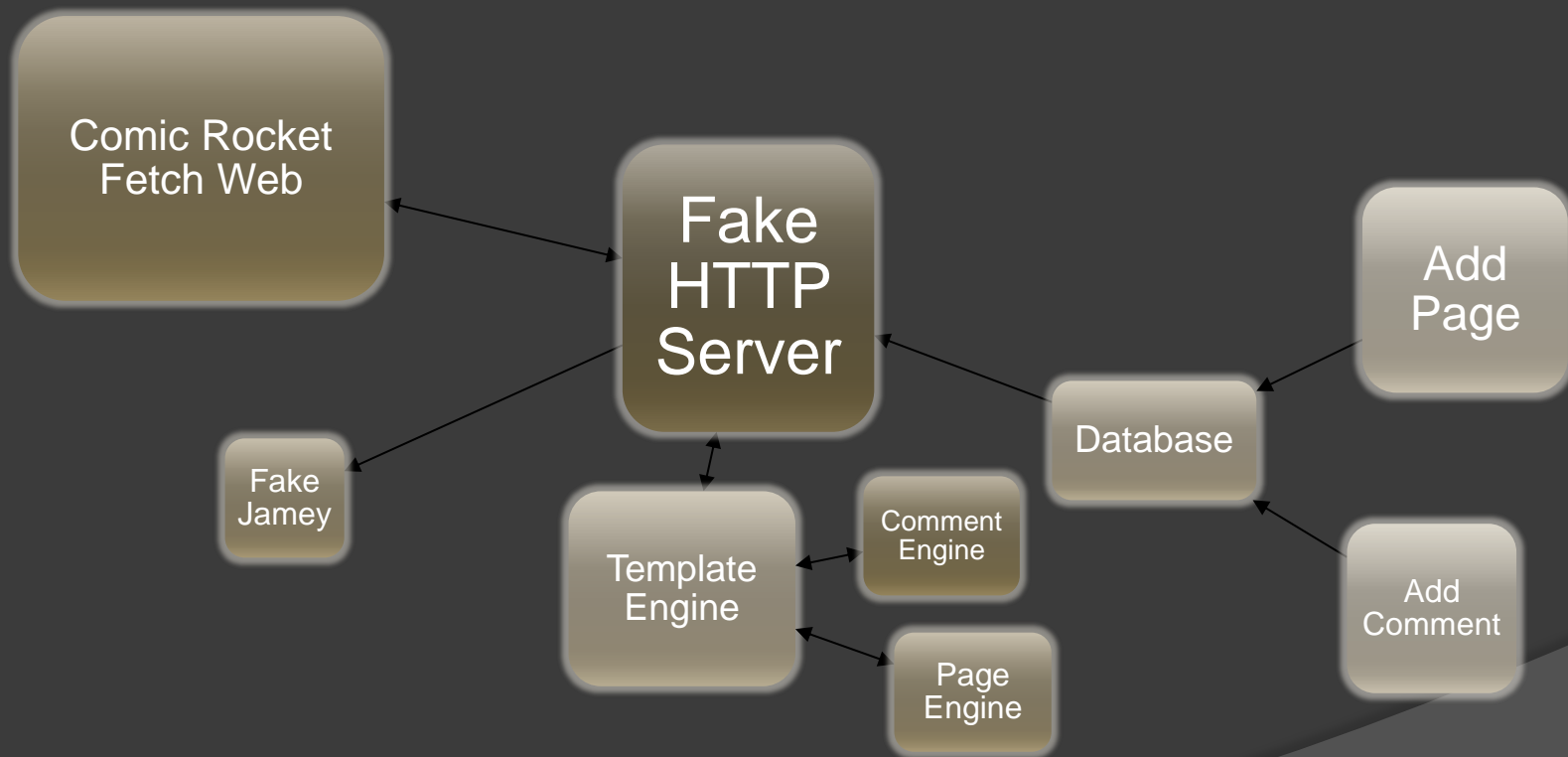
---

Most Web Comics don't appreciate extra bandwidth from "testing"

---



# Solution





COMIC ROCKET™

# Solution

## Http Server

- Parses specific domains, and calls the template engine to build the page

## Template Engine

- Builds a page, with {key = value} pairs
- Builds comments in a similar manner

## Crontab

- Scheduler to run scripts at specific times, ie updates for comics or comment insertion

## Hosts File

- Redirects specific domains to 127.0.0.1 (or localhost)



COMIC ROCKET™

# Testing the Rocket Booster

Time Scale

First Scale: 1 hour = 75 sec

Second Scale: 1 hour = 5 sec



Update Schedules

Uniform Schedule

Variable Schedule



Comments

HTML Diff

Link Diff



COMIC ROCKET™

# Future Enhancements

- ◉ Should be easily extensible to store more page information
- ◉ Compression of the cache
- ◉ Filter out unimportant changes