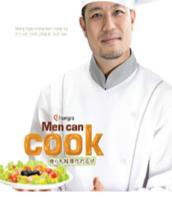# Framgia Vietnam Co.,ltd

# Professional Company

# Friendly working environment

# Recruitment

- Ruby Developer
- Senior PHP Developer
- Senior Java Developer
- Senior .NET Developer
- BrSE
- New Developer

● hr_team@framgia.com

# Extract Content from Articles / News Pages

*Giang Nguyen @nguyenducgiang*

# background

- API for a news reader app

  - web pages content extraction is vital

# unify content structure

- title

- description

- readability content

- images

# every site is different

- frontend styles

- backend configurations

- not every site is of news / article type

# approach

- a general parser to extract desire information

    - any kind of pages, focus on articles / news

- custom parsers for specific cases

- post-proceed extracted content for extra goodness

# notice when fetching urls

- mime types

- redirections

- user agents

**NET::HTTP is good enough!**

- http / https

- ~~authentication~~

- ~~proxies~~

# tools & techniques

- Nokogiri FTW

- ruby-readability

- HTML metas

  - w3c specified / Open Graph / Twitter Card

# title – what could be wrong?

- length (truncated title)

| | |
|---|---|
| og_title | Rails Conf 2014 Concerns, Decorators, Presenters, Service-objects, He… |
| twitter_title | |
| meta_title | Rails Conf 2014 Concerns, Decorators, Presenters, Service-objects, He… |
| match_title | Rails Conf 2014 Concerns, Decorators, Presenters, Service-objects, Helpers, Help Me Decide-april-22-2014 |

- site name in title

| | |
|---|---|
| meta_title | Northern Ireland GPs could have saved £73m says PAC - BBC News |

# title - approach

- pattern to recognize visible title

- longest common substring

- regex to remove site name from title

# description

- og:description || meta description

- yes, it just this simple!

# content - headaches

- encoding

  - work for a JP company? I know that feel, mate.

- ReactJS & heavy AJAX

- complex / alien structured articles

  - multiple pages content, slides, etc.

# content - approach

- ruby-readability (monkey-patched) for general use

- custom methods for specific sites

- reject overcomplicated sites

# images - issues

- no images!!!

- relative path

- size

- fixed og:image

# select feature image

- og:image || first image in content || snapshot of page

- some sites need specific tune up

# a possible (better) approach

- machine learning?

- feed sample correct content (per site)

- let the machine

  - guess the element it should get content from the base

    samples

  - improve itself from guess results

# what next!

- further manipulation of extracted content

  - category & type recognition

  - language recognition

- involve machine learning

- increase users' experience

  - related & recommend news / articles

# conclusion

- there is no absolute / perfect solutions

    - requirements / priorities / desire

- understand the general issues make it easier to work with specific one

- there are gems for almost everything, but writing you own lib helps to understand the issues in a higher level