

Whirlpool

Data Acquisition using N-node Distributed Web Crawler

Rihan Pereira, MSCS

Advisor: Dr. Michael Soltys
Department of Computer Science
MSCS Graduate 2018-2019

November 24, 2019



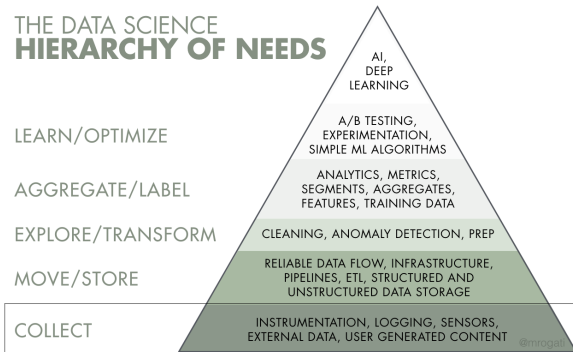
Channel Islands

CALIFORNIA STATE UNIVERSITY

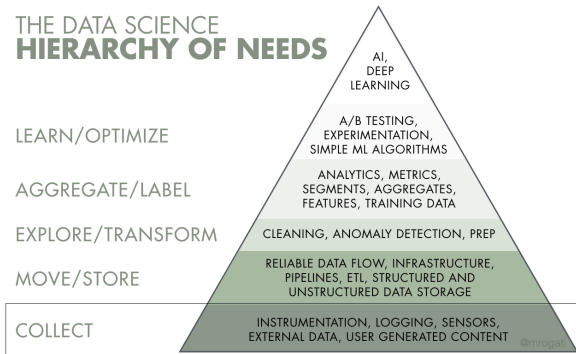
- 1 Motivation & Contribution
- 2 Crawler characteristics & history
- 3 Mercator 1999(Heydon & Najork)
- 4 Software Design Principles
- 5 Whirlpool: Event-driven architecture
- 6 Whirlpool: Parser
- 7 Whirlpool: Near-Deduplication
- 8 Whirlpool: Distributed Crawling
- 9 Whirlpool: Operations
- 10 Future work

Motivation & Contribution

Motivation



Motivation



Self-actualization (AI) is great, but you first need food, water, and shelter (data literacy, collection, and infrastructure).”

Contributions

to be completed

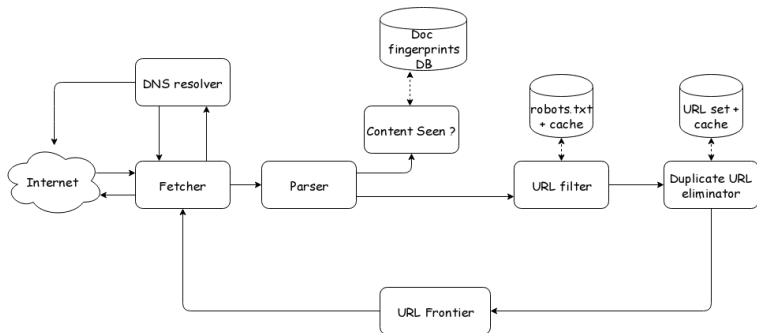
Crawler characteristics & history

Crawl char & hist.

to add something

Mercator 1999(Heydon & Najork)

Mercator background



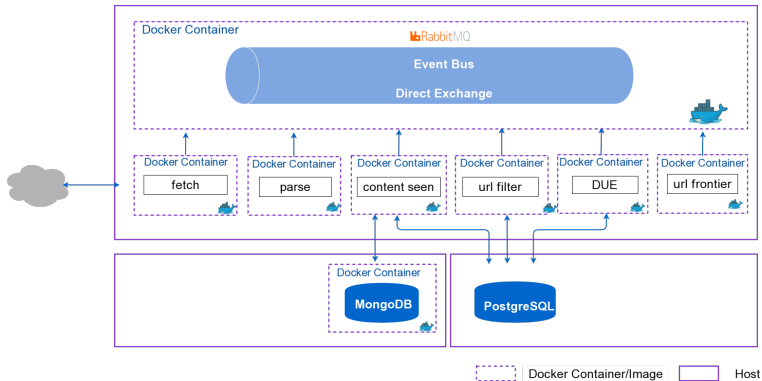
Software Design Principles

Soft. Design

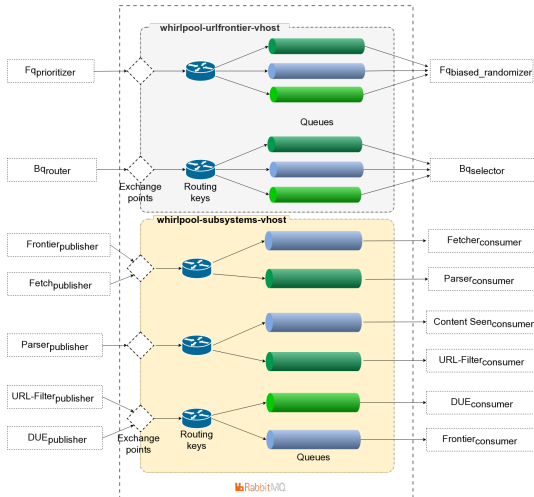
to add something

Whirlpool: Event-driven architecture

RabbitMQ: Message bus



Direct Worker Queue Data Flow



Whirlpool: Parser

Parser

to add something

Whirlpool: Near-Deduplication

Dedupe

to add something

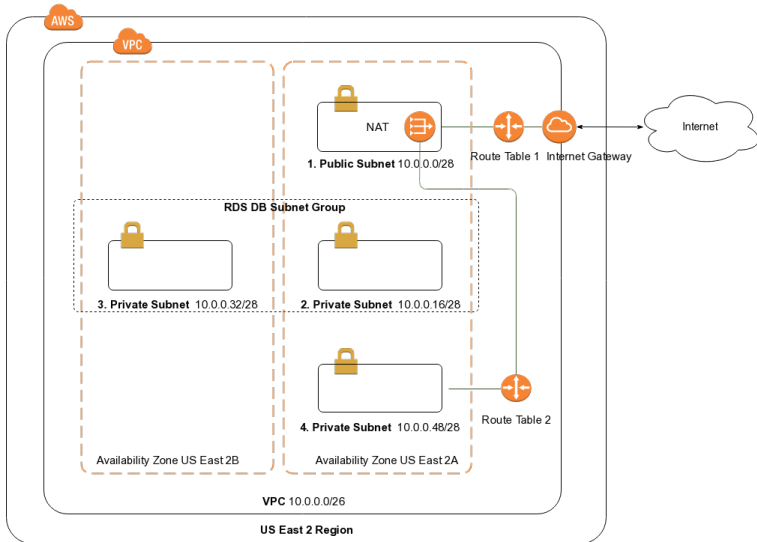
Whirlpool: Distributed Crawling

Dist. crawl

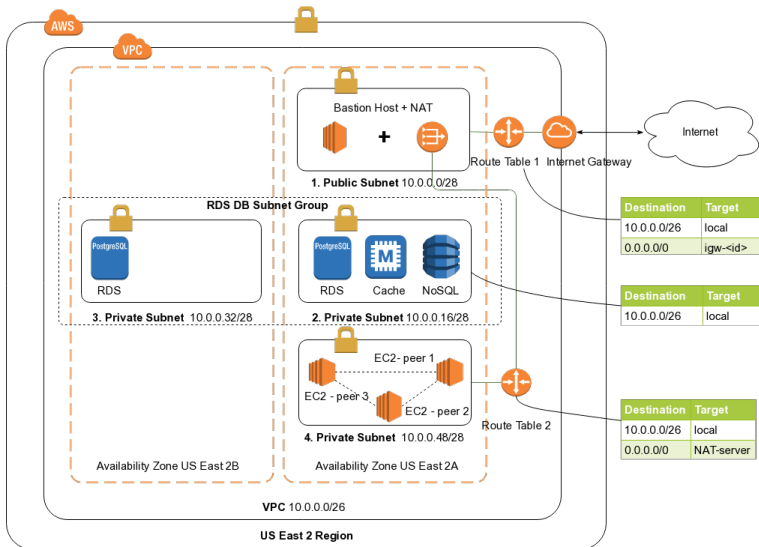
to add something

Whirlpool: Operations

From 10,000 ft.



From 5,000 ft.



Future work

future to do

to add something

Thank you! Questions ?