# 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 27, 2019

- Motivation & Contributions

- Motivation & Contributions
- Crawler characteristics & history

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)          } background
- Software Design Principles

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)   } background
- Software Design Principles
- Whirlpool: Event driven architecture

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork) ⎫
- Software Design Principles        ⎬ background
                                    ⎭
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)      } background
- Software Design Principles
- Whirlpool: Event driven architecture
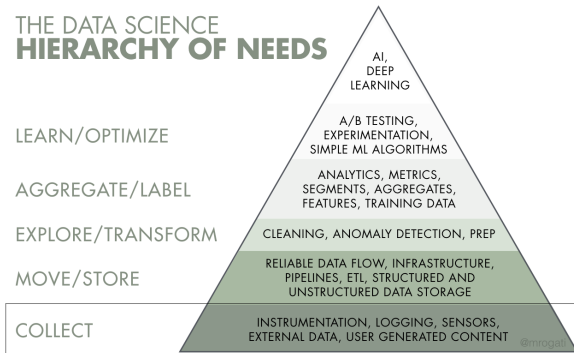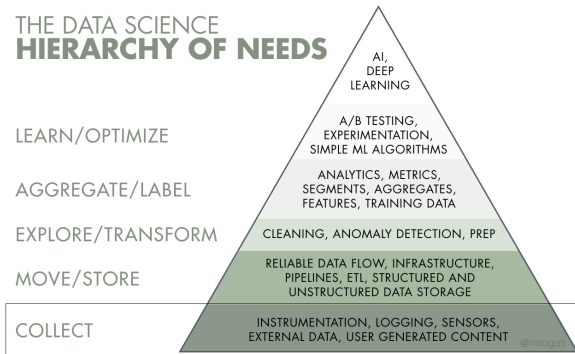- Whirlpool: Fetcher
- Whirlpool: Parser

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication

} background

- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)       } background
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication
- Whirlpool: Distributed Crawling

- Motivation & Contributions
- Crawler characteristics & history ⎫
- Mercator 1999 (Heydon & Najork) ⎬ background
- Software Design Principles ⎭
- Whirlpool: Event driven architecture ⎫
- Whirlpool: Fetcher ⎪
- Whirlpool: Parser ⎪
- Whirlpool: Deduplication ⎬ implementation & results
- Whirlpool: Distributed Crawling ⎪
- Whirlpool: Operations ⎭

- Motivation & Contributions
- Crawler characteristics & history ⎫
- Mercator 1999 (Heydon & Najork) ⎬ background
- Software Design Principles ⎭
- Whirlpool: Event driven architecture ⎫
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication ⎬ implementation & results
- Whirlpool: Distributed Crawling
- Whirlpool: Operations ⎭
- Future work

Motivation & Contribution

## Motivation



THE DATA SCIENCE
**HIERARCHY OF NEEDS**

AI, DEEP LEARNING

LEARN/OPTIMIZE — A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

AGGREGATE/LABEL — ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

EXPLORE/TRANSFORM — CLEANING, ANOMALY DETECTION, PREP

MOVE/STORE — RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

COLLECT — INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

## Motivation



THE DATA SCIENCE
**HIERARCHY OF NEEDS**

AI, DEEP LEARNING

LEARN/OPTIMIZE — A/B TESTING, EXPERIMENTATION, SIMPLE ML ALGORITHMS

AGGREGATE/LABEL — ANALYTICS, METRICS, SEGMENTS, AGGREGATES, FEATURES, TRAINING DATA

EXPLORE/TRANSFORM — CLEANING, ANOMALY DETECTION, PREP

MOVE/STORE — RELIABLE DATA FLOW, INFRASTRUCTURE, PIPELINES, ETL, STRUCTURED AND UNSTRUCTURED DATA STORAGE

COLLECT — INSTRUMENTATION, LOGGING, SENSORS, EXTERNAL DATA, USER GENERATED CONTENT

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

## Coverage & Freshness

# Web crawlers (1990 - 2019)

| 1993 | 1994 | 1996 | 1999 | 2000 | 2001 | 2004 | 2006 | 2009 | 2016 |
|------|------|------|------|------|------|------|------|------|------|
| Wanderer | MOMSpider | Internet Archive | Mercator | Polybot | IBM Webfountain | Ubicrawler | Multicrawler | IRLbot | Software Heritage |

Mercator 1999 (Heydon & Najork)

## basic crawling algorithm

1: Let I $\leftarrow$ {1,2,3,4,5} such that seed set S = {$U_i$| i $\in$ I}
2: $U_f \leftarrow$ S; where $U_f$ is a Frontier queue
3: **procedure** Spider($U_f$)
4:     **while** $U_f \neq \emptyset$ **do**
5:         $u \leftarrow$ Pop($U_f$)
6:         $p \leftarrow$ Fetch($u$)
7:         $T \leftarrow \exists p$ [{Extract($p, t$) | $t$ is a text }]
8:         $L \leftarrow \exists p$ [{Extract($p, l$) | $l$ is a link }]
9:         $U_f \leftarrow U_f \cup L$
10:        $\exists u$ [{Delete($U_f, u$) | $u$ is a already fetched URL}]
11: **end**

# Mercator background



Figure: Mercator building blocks (Heydon & Najork)

URL Frontier Scheme

Front queue (Frontier Queue)

# Back queue (Frontier Queue)

Software Design Principles

# Designing scalable systems

# Designing scalable systems

- Adding identical copies of components

Designing scalable systems

- Adding identical copies of components
- Functional partitioning

## Designing scalable systems

- Adding identical copies of components
- Functional partitioning



- Data partitioning

## State Management



node 1 halts and comes back online, making node 2 hash url1 again and register causing multiples copies of cached object

Figure: identical copies of same cached object

## State Management



Figure: Using local locks to access shared resources

## State Management



Figure: using shared locks to access shared resources
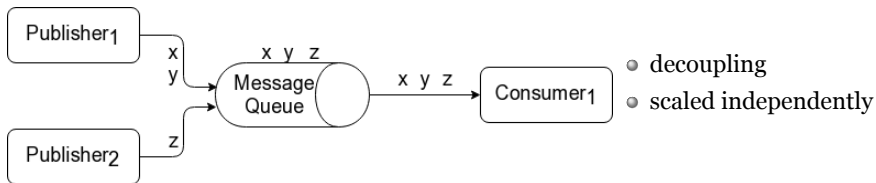
Whirlpool: Event-driven architecture

## Message Queue(MQ)



Publisher₁ — x, y →
Publisher₂ — z →

Message Queue: x y z

→ x y z → Consumer₁

## Message Queue(MQ)



- decoupling

# Message Queue(MQ)



- decoupling
- scaled independently

# Message Queue(MQ)



- decoupling
- scaled independently
- balancing traffic

# Message Queue(MQ)



- decoupling
- scaled independently
- balancing traffic
- fault-tolerance

## MQ: Routing mechanisms

## MQ: Routing mechanisms

- Direct Worker Queue Data Flow

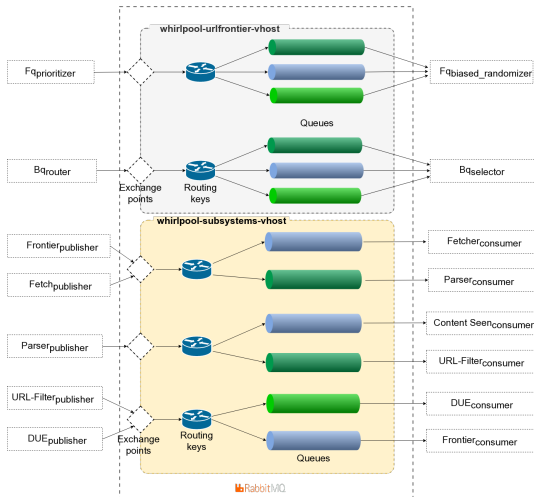## MQ: Routing mechanisms

- Direct Worker Queue Data Flow
- Fanout

## MQ: Routing mechanisms

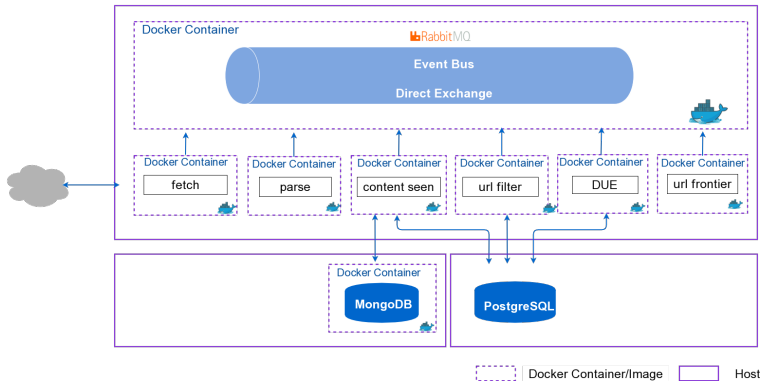- Direct Worker Queue Data Flow
- Fanout
- Topic

## MQ: Routing mechanisms

- Direct Worker Queue Data Flow
- Fanout
- Topic
- Header

# Direct Worker Queue Data Flow

# RabbitMQ: Message bus

development vs. production docker containers

things to add

Whirlpool: Fetcher

Whirlpool: Parser

## Parser

to add something

Whirlpool: Near-Deduplication
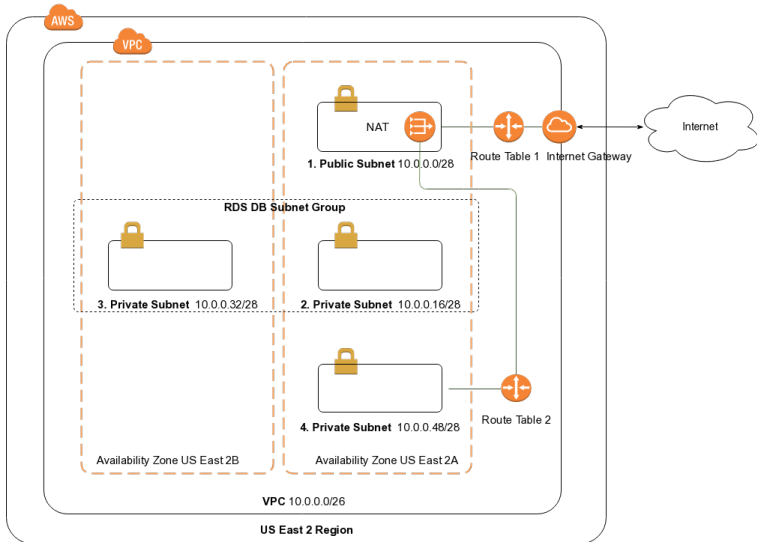
## Dedupe

to add something

Whirlpool: Distributed Crawling
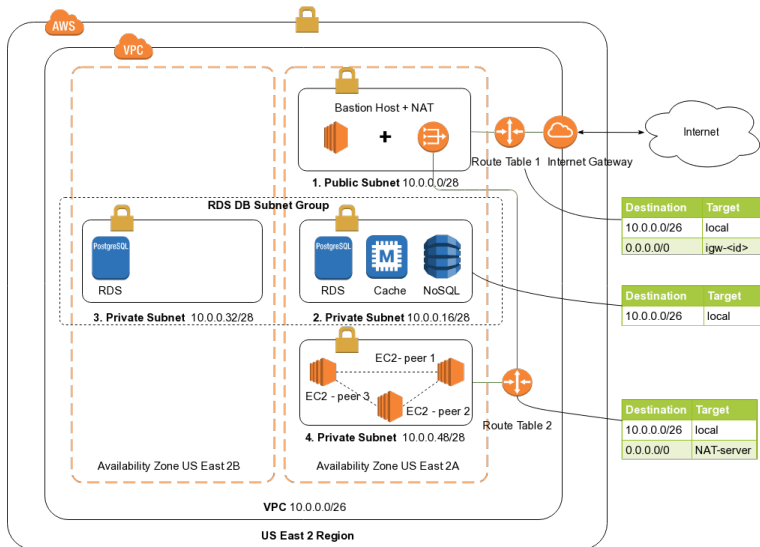
# Dist. crawl

to add something

Whirlpool: Operations

# From 10,000 ft.

# From 5,000 ft.

Future work

# future to do

to add something

Thank you! Questions ?