# Whirlpool

## Data Acquisition using N-node Distributed Web Crawler

Rihan Pereira, MSCS

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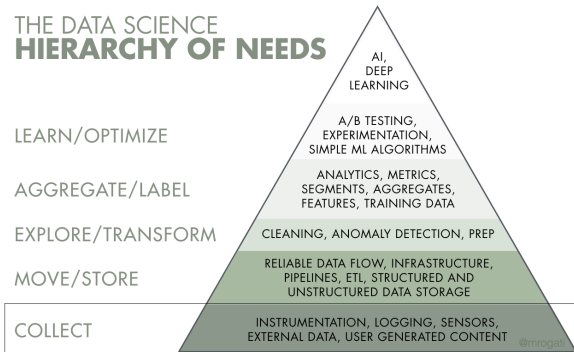MSCS Graduate 2018-2019

November 27, 2019

Motivation & 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

@mrogati

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



Comprehensive crawl

Scoped crawl  Topical crawl

Coverage

Freshness

# Web crawlers (1990 - 2019)

to add something

Mercator 1999 (Heydon & Najork)

# basic crawling algorithm
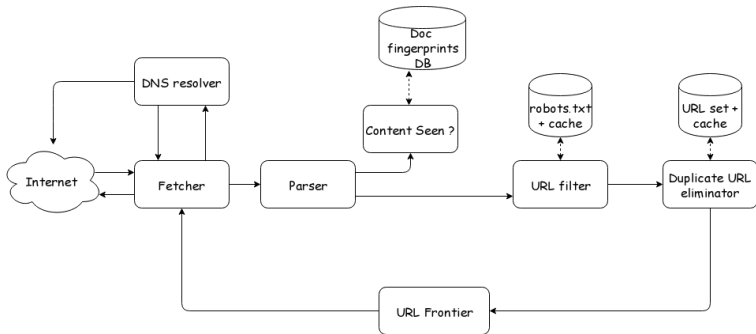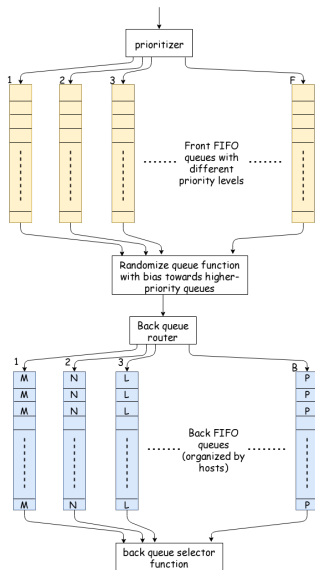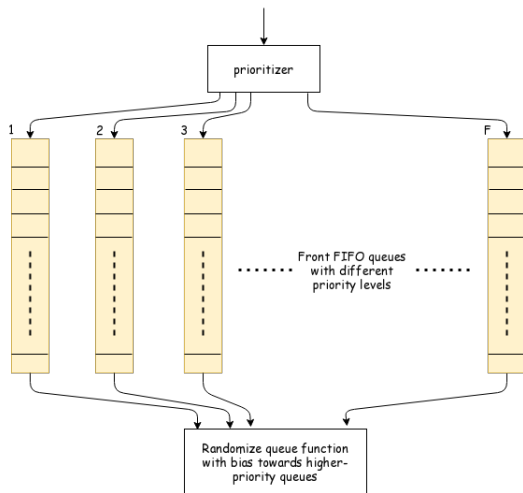
to add content

## Mercator background



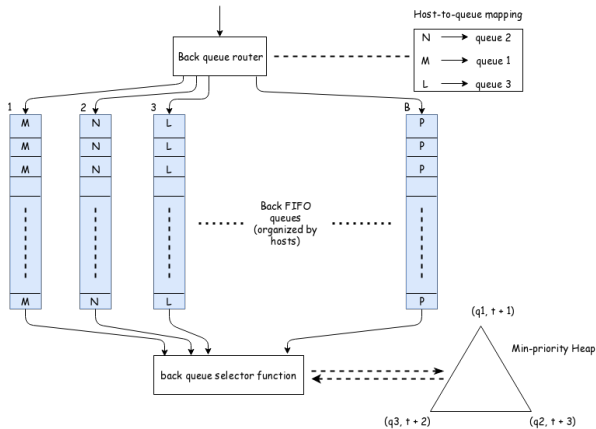Figure: Mercator building blocks (Heydon & Najork)

## URL Frontier Scheme

# Front queue (Frontier Queue)



Front FIFO queues with different priority levels

# Back queue (Frontier Queue)

Software Design Principles

Designing scalable systems

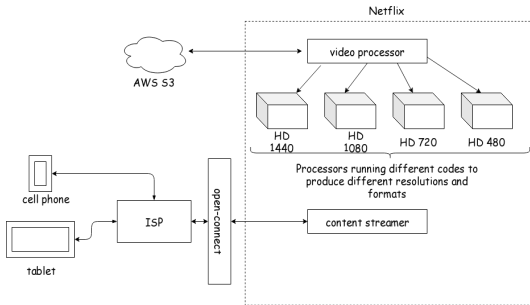Designing scalable systems

- Adding identical copies of components

Designing scalable systems

- Adding identical copies of components
- Functional partitioning

Designing scalable systems

- Adding identical copies of components
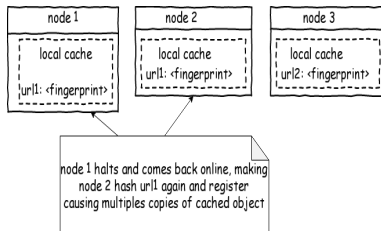- Functional partitioning



- Data partitioning

## State Management



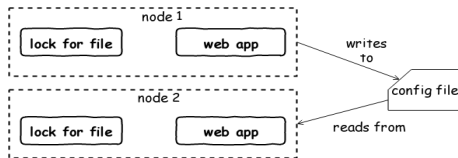Figure: identical copies of same cached object

## State Management



Figure: Using local locks to access shared resources

## State Management
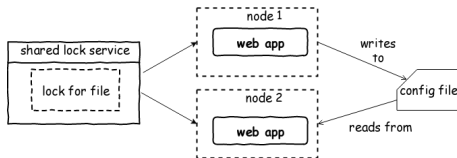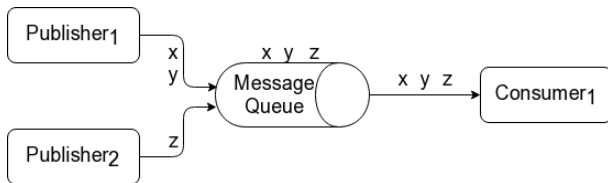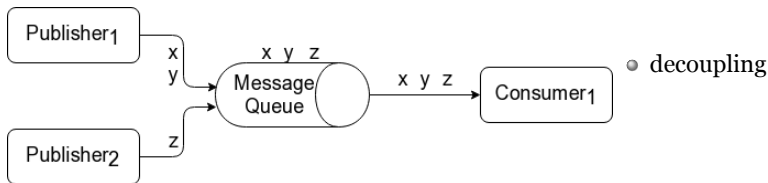


Figure: using shared locks to access shared resources
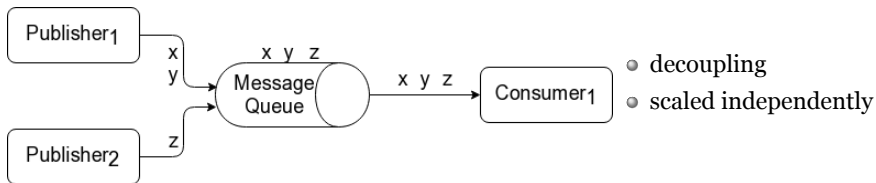
Whirlpool: Event-driven architecture
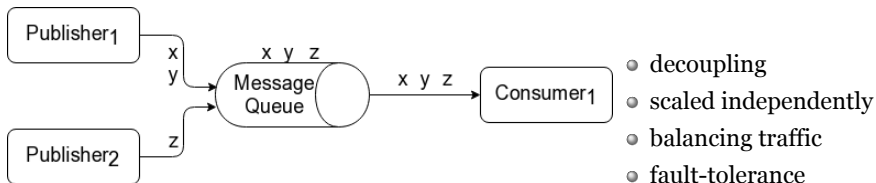
# Message Queue(MQ)

## 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

Motiv. & Contrib \> Crawler history \> Mercator \> Soft. design \> Event-driven \> Parser \> Deduplication \> Dist. Crawling \> Opworks \> Future

MQ: Routing mechanisms

- Direct Worker Queue Data Flow

## MQ: Routing mechanisms

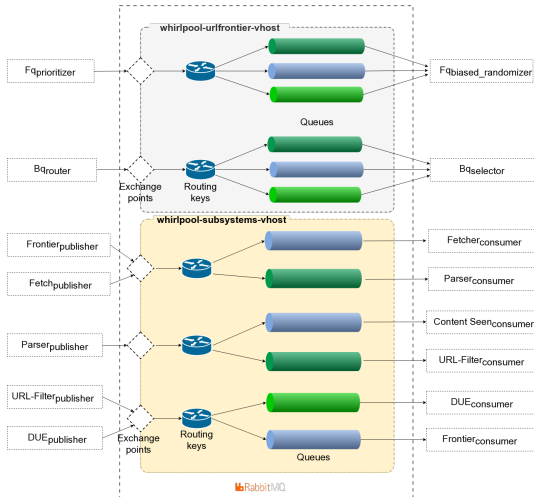- Direct Worker Queue Data Flow
- Fanout

## MQ: Routing mechanisms

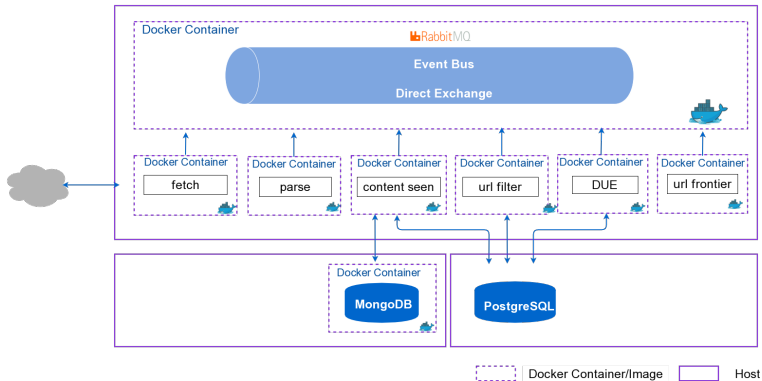- Direct Worker Queue Data Flow
- Fanout
- Topic

## 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: Parser

# Parser

to add something

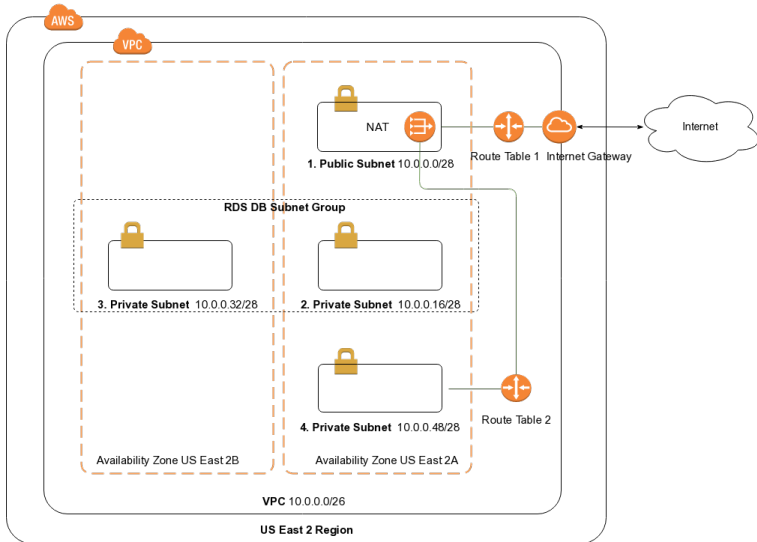Whirlpool: Near-Deduplication

# Dedupe

to add something

Whirlpool: Distributed Crawling
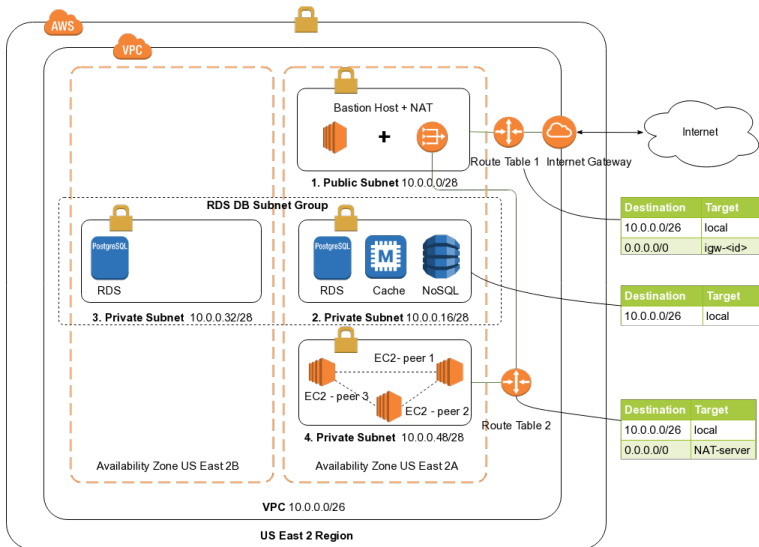
## Dist. crawl

to add something

Whirlpool: Operations

# From 10,000 ft.

# From 5,000 ft.

Future work

future to do

to add something

Thank you! Questions ?