

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 30, 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)
 - Software Design Principles
- }
- background

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

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

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

}

background

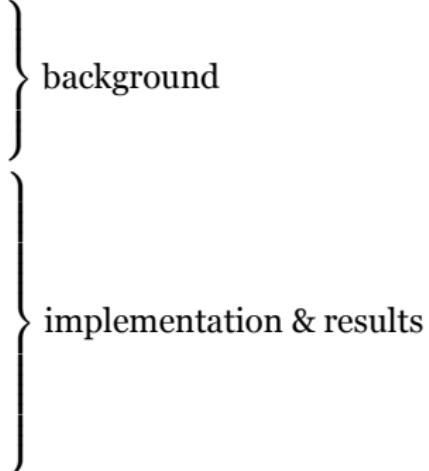
- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication

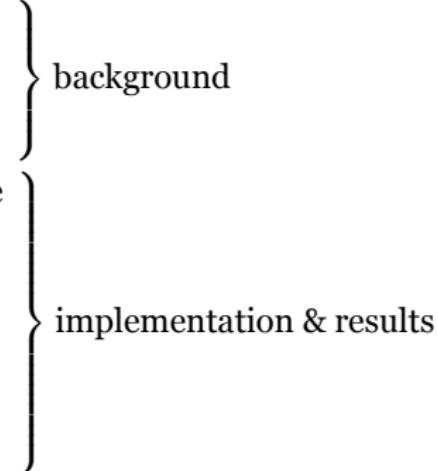


- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication
- Whirlpool: Distributed Crawling

}

background

- Motivation & Contributions
 - Crawler characteristics & history
 - Mercator 1999 (Heydon & Najork)
 - Software Design Principles
 - Whirlpool: Event driven architecture
 - Whirlpool: Fetcher
 - Whirlpool: Parser
 - Whirlpool: Deduplication
 - Whirlpool: Distributed Crawling
 - Whirlpool: Operations
- 
- The list of items is grouped into two main categories: 'background' and 'implementation & results'. The 'background' category includes the first four items: Motivation & Contributions, Crawler characteristics & history, Mercator 1999 (Heydon & Najork), and Software Design Principles. The 'implementation & results' category includes the remaining six items: Whirlpool: Event driven architecture, Whirlpool: Fetcher, Whirlpool: Parser, Whirlpool: Deduplication, Whirlpool: Distributed Crawling, and Whirlpool: Operations.

- Motivation & Contributions
 - Crawler characteristics & history
 - Mercator 1999 (Heydon & Najork)
 - Software Design Principles
 - Whirlpool: Event driven architecture
 - Whirlpool: Fetcher
 - Whirlpool: Parser
 - Whirlpool: Deduplication
 - Whirlpool: Distributed Crawling
 - Whirlpool: Operations
 - Future work
- 
- The list of items is grouped into two main categories: 'background' and 'implementation & results'. The 'background' category includes the first four items: Motivation & Contributions, Crawler characteristics & history, Mercator 1999 (Heydon & Najork), and Software Design Principles. The 'implementation & results' category includes the remaining seven items: Whirlpool: Event driven architecture, Whirlpool: Fetcher, Whirlpool: Parser, Whirlpool: Deduplication, Whirlpool: Distributed Crawling, Whirlpool: Operations, and Future work.

Motivation & Contribution

Motivation

THE DATA SCIENCE HIERARCHY OF NEEDS

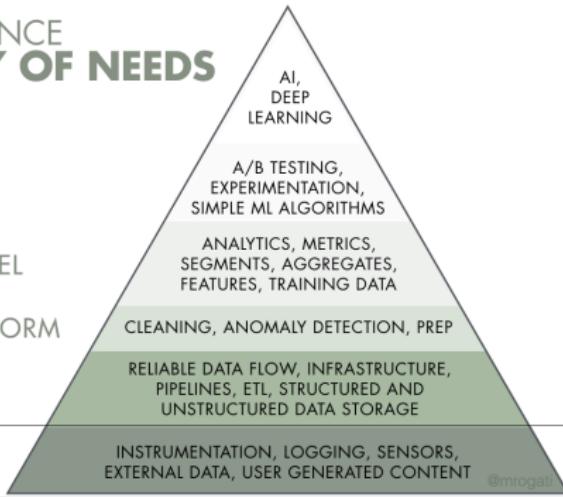
LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



Motivation

THE DATA SCIENCE HIERARCHY OF NEEDS

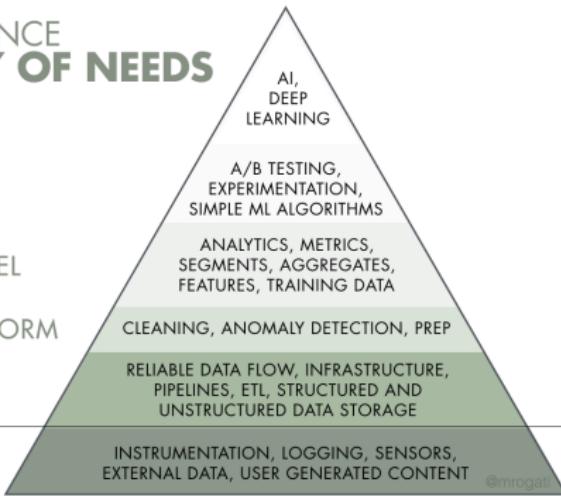
LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



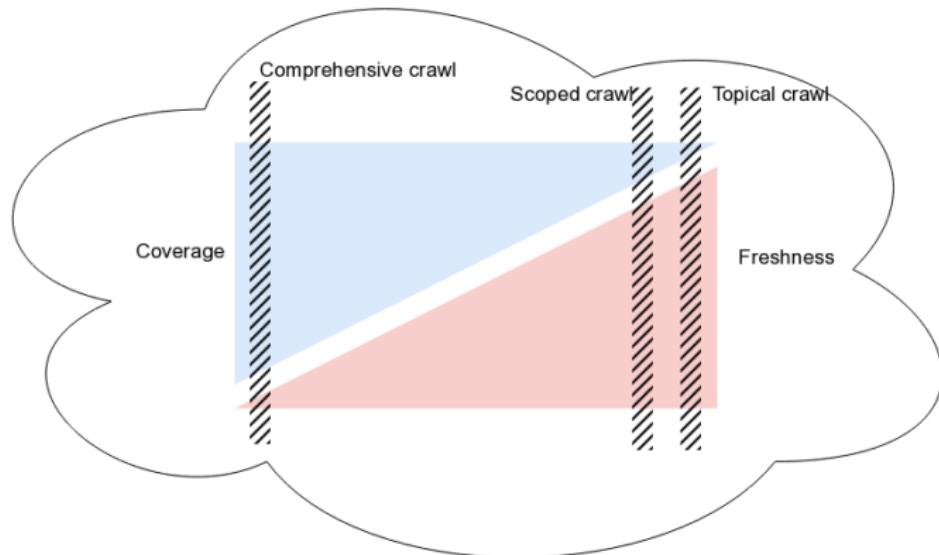
Self-actualization (AI) is great, but you first need food, water, and shelter (data literacy, collection, and infrastructure)."

Contributions

to be completed

Crawler characteristics & history

Coverage & Freshness



Seedlist



Web crawlers (1990 - 2019)



Mercator 1999 (Heydon & Najork)

basic crawling algorithm

```
1: Let I  $\leftarrow \{1,2,3,4,5\}$  such that seed set S =  $\{U_i \mid i \in I\}$ 
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:      $u \leftarrow \text{Pop}(U_f)$ 
6:      $p \leftarrow \text{Fetch}(u)$ 
7:      $T \leftarrow \exists p \ [\{\text{Extract}(p, t) \mid t \text{ is a text}\}]$ 
8:      $L \leftarrow \exists p \ [\{\text{Extract}(p, l) \mid l \text{ is a link}\}]$ 
9:      $U_f \leftarrow U_f \cup L$ 
10:     $\exists u \ [\{\text{Delete}(U_f, u) \mid u \text{ is a already fetched URL}\}]$ 
11:  end
```

Mercator background

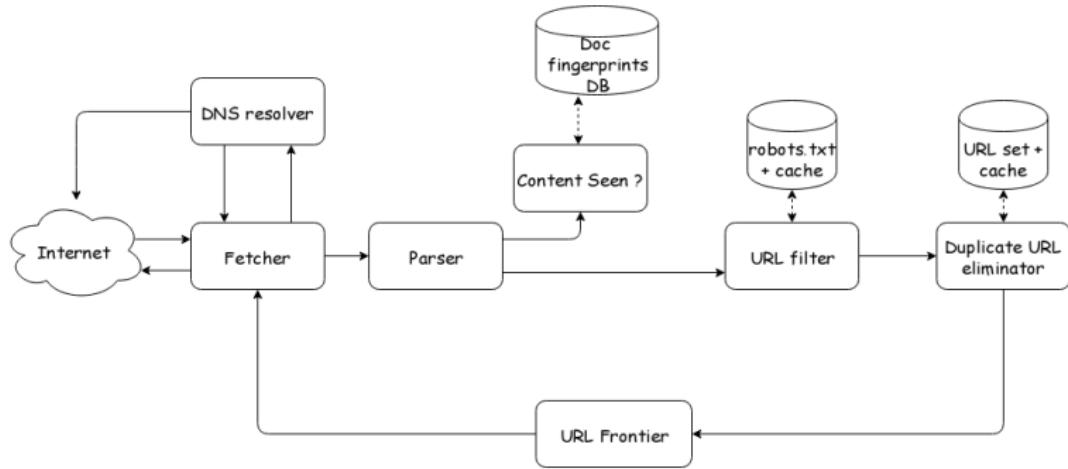
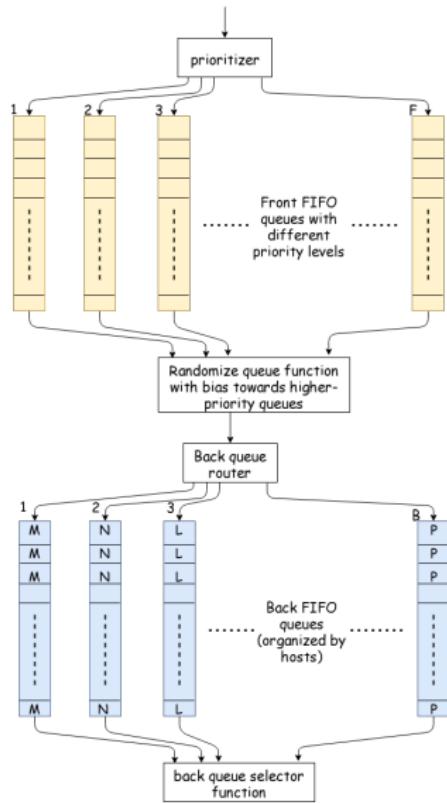
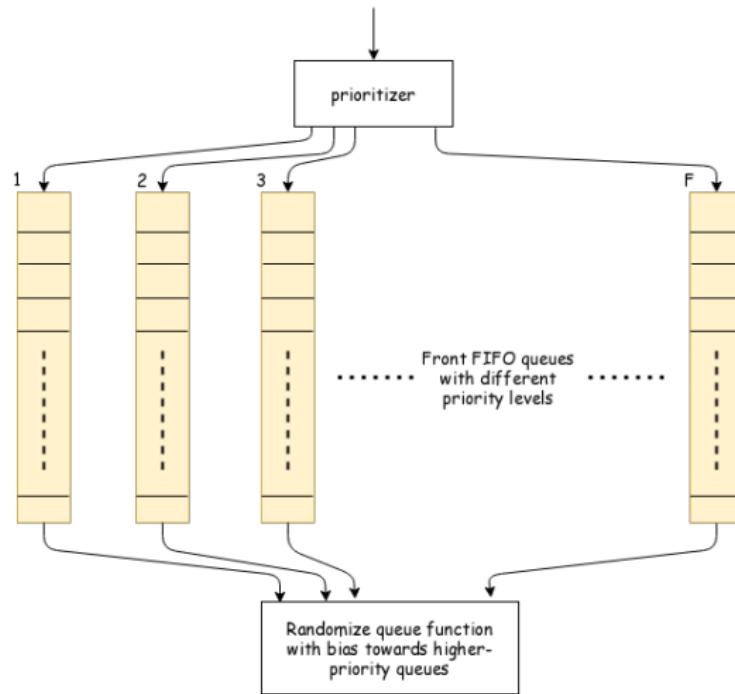


Figure: Mercator building blocks (Heydon & Najork)

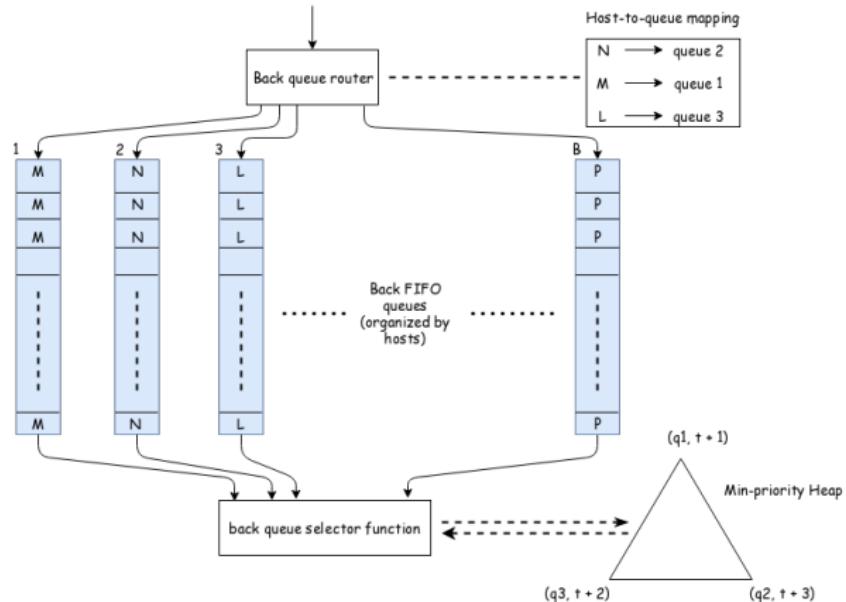
URL Frontier Scheme



Front queue (Frontier Queue)



Back queue (Frontier Queue)



Software Design Principles

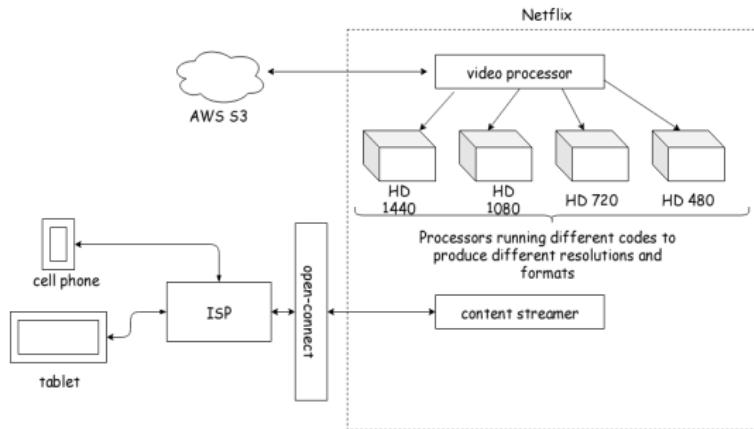
Designing scalable systems

Designing scalable systems

- Adding identical copies of components

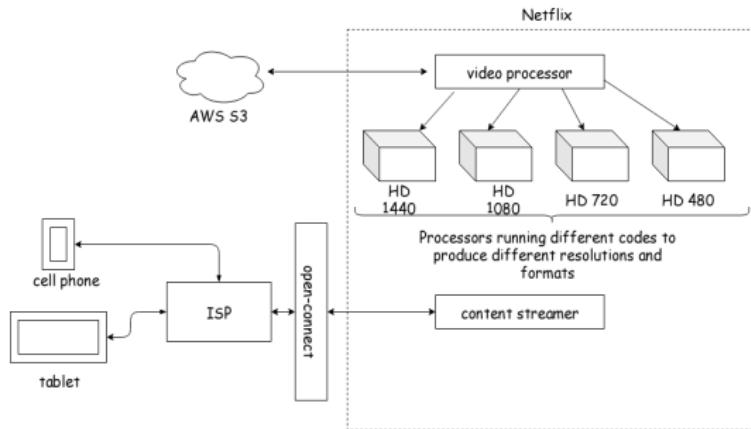
Designing scalable systems

- Adding identical copies of components
- Functional partitioning



Designing scalable systems

- Adding identical copies of components
- Functional partitioning



- Data partitioning

State Management

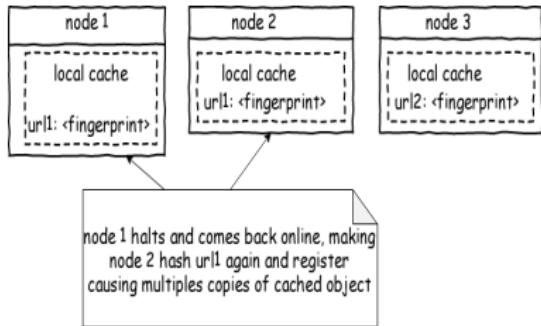


Figure: identical copies of same cached object

State Management

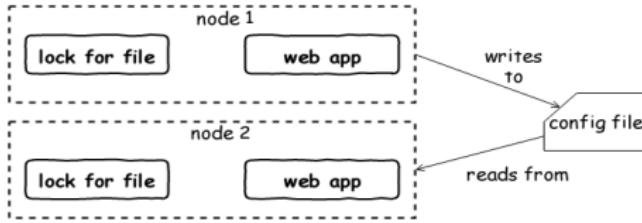


Figure: Using local locks to access shared resources

State Management

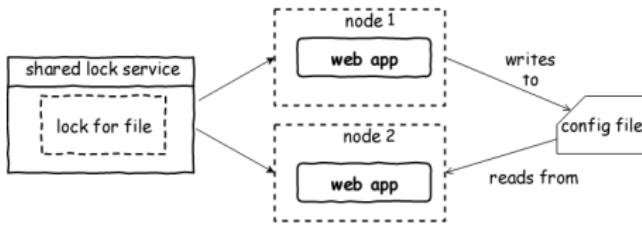


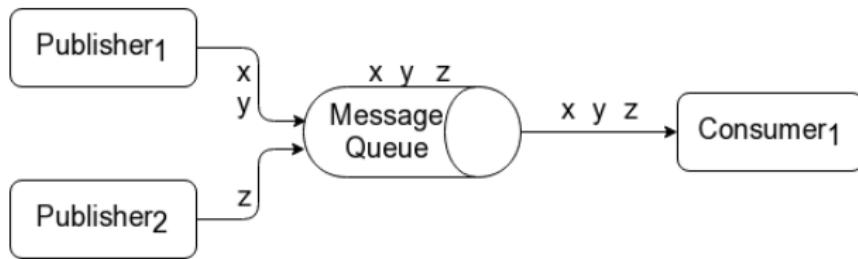
Figure: using shared locks to access shared resources

Whirlpool: Event-driven architecture

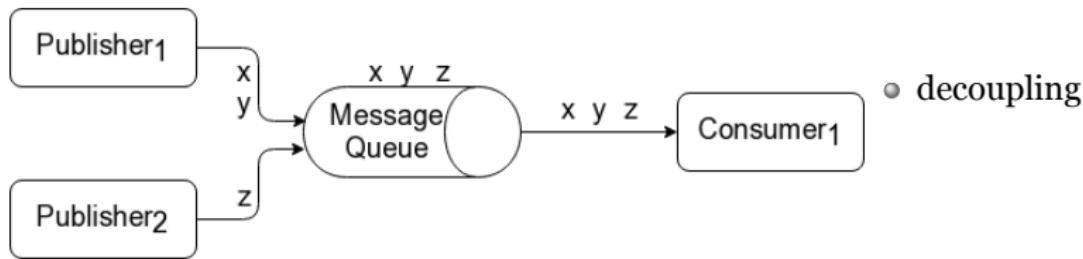
Message buses



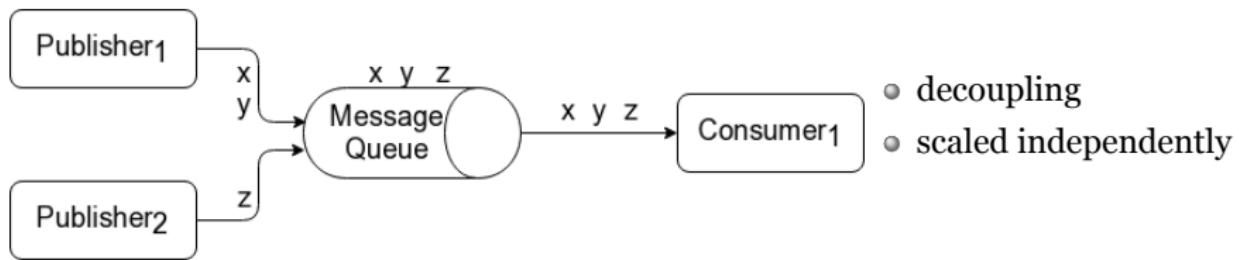
Message Queue(MQ)



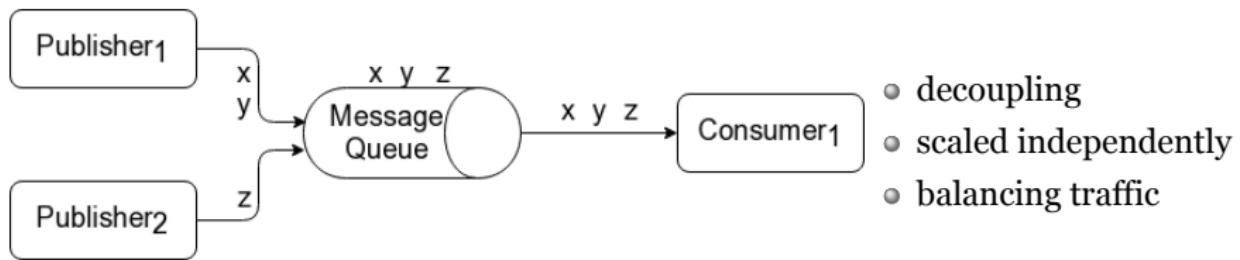
Message Queue(MQ)



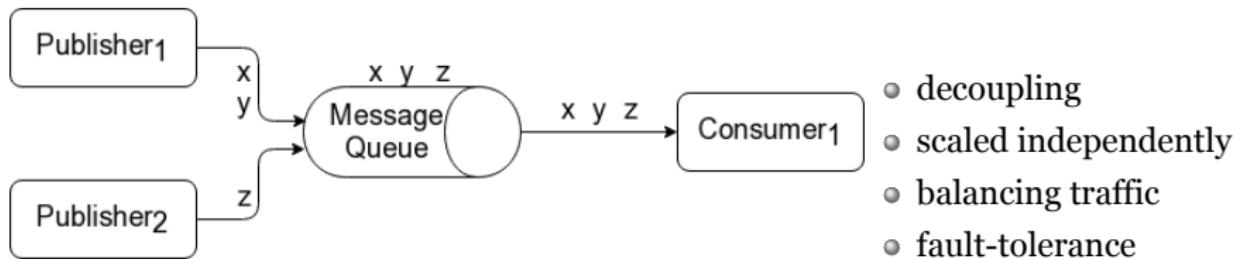
Message Queue(MQ)



Message Queue(MQ)



Message Queue(MQ)



MQ: Routing mechanisms

MQ: Routing mechanisms

- Direct Worker Queue Data Flow

MQ: Routing mechanisms

- Direct Worker Queue Data Flow
- Fanout

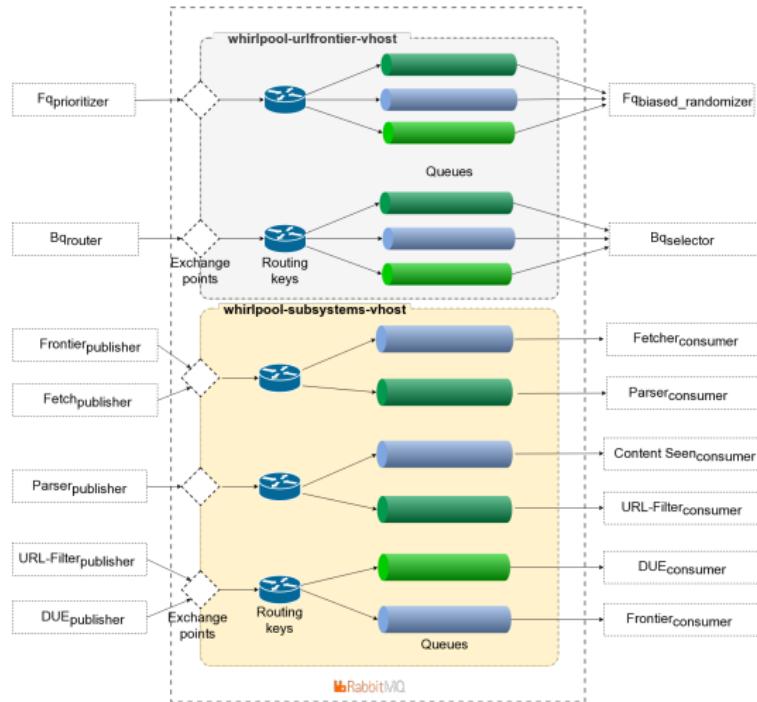
MQ: Routing mechanisms

- Direct Worker Queue Data Flow
- Fanout
- Topic

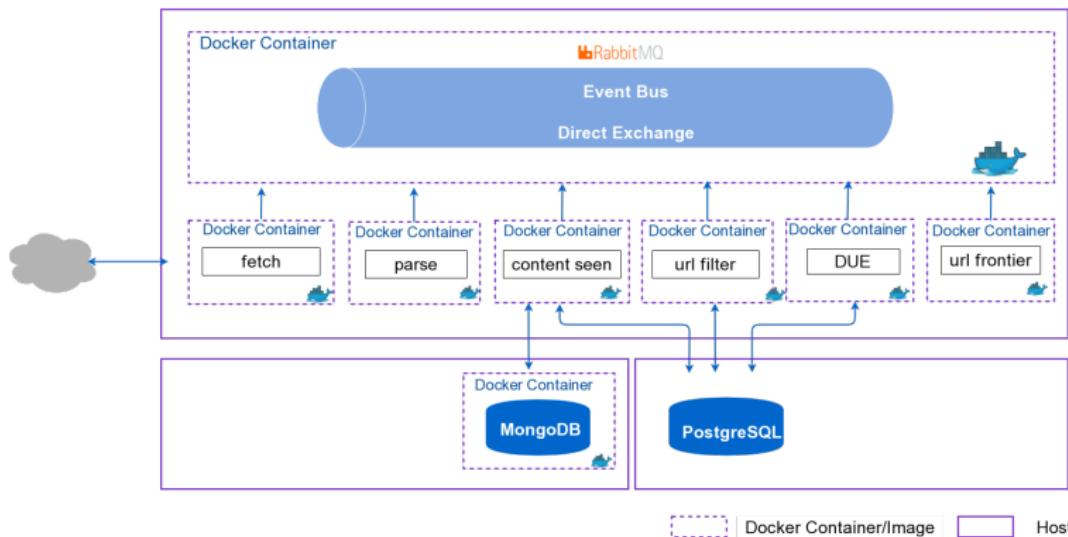
MQ: Routing mechanisms

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

Direct Worker Queue Data Flow



RabbitMQ: Message bus



development vs. production docker containers

```
1  version: '2.4'
2
3  networks:
4      default:
5          external:
6              name: whirlpool-net
7
8  services:
9
10  base:
11      image: python:3.7.4-buster
12      command: bash -c "useradd --create-home --shell /bin/bash &
13          &chown -R whirlpool:whirlpool /home/whirlpool/whirlp
14          ool-due"
15      volumes:
16          - ./:/home/whirlpool/whirlpool-due
17          - wh-due:/usr/local/lib/python3.7/site-packages
18      working_dir: /home/whirlpool/whirlpool-due
19      environment:
20          - PY_ENV=development
21      networks:
22          - default
23
24  install:
25      extends:
26          service: base
27      command: pip3 install -r requirements.txt
28
29  quick-up:
30      extends:
31          service: base
32      command: python3 ./due/main.py
33
34  volumes:
35      wh-due:
36          external: true
```

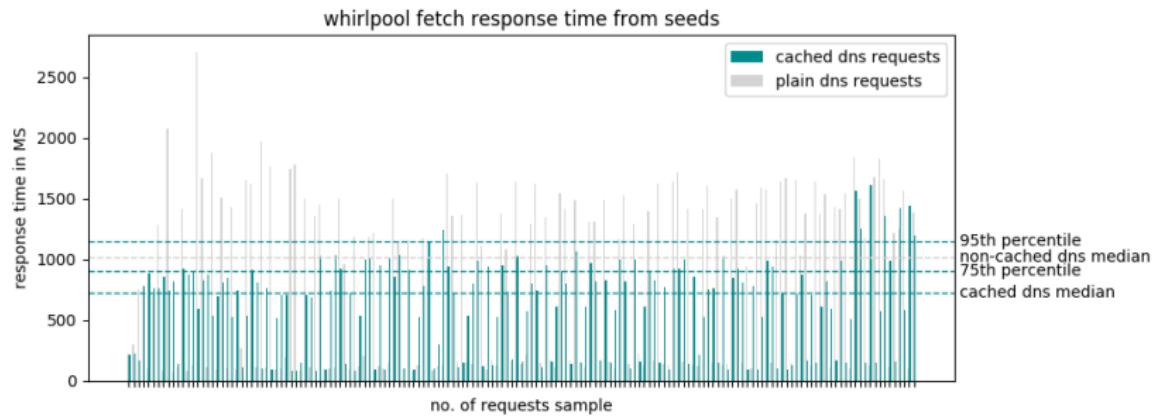
development vs. production docker containers

```
1 FROM python:3.7.4-buster as whirlpool-due-base
2
3 ENV PYTHONDONTWRITEBYTECODE=1
4 ARG WH_DUE_ROOT=/home/whirlpool/whirlpool-due
5 WORKDIR $WH_DUE_ROOT
6
7 RUN apt-get update \
8   && apt-get install -y --no-install-recommends netcat \
9   && rm -rf /var/lib/apt/lists/* \
10  && useradd --create-home --shell /bin/bash whirlpool \
11  && chown -R whirlpool:whirlpool $WH_DUE_ROOT
12
13 # files necessary to build the project
14 COPY .pylintrc ./
15 COPY requirements.txt ./
16
17 RUN mkdir logs/ \
18   && pip3 install -r requirements.txt
19
20 COPY scripts/ scripts/
21 COPY due/ due/
22
23 # docker image for dev target
24 FROM whirlpool-due-base as whirlpool-due-dev
25
26 COPY scripts/wait-for-it.sh scripts/wait-for-it.sh
27 ENTRYPOINT ["bash ./scripts/wait-for-it.sh"]
28
29 # docker image for prod target
30 FROM whirlpool-due-base as whirlpool-due-prod
31
32 COPY scripts/wait-for-it-prod.sh scripts/wait-for-it-prod.sh
33 ENTRYPOINT ["bash ./scripts/wait-for-it-prod.sh"]
```

development vs. production docker containers

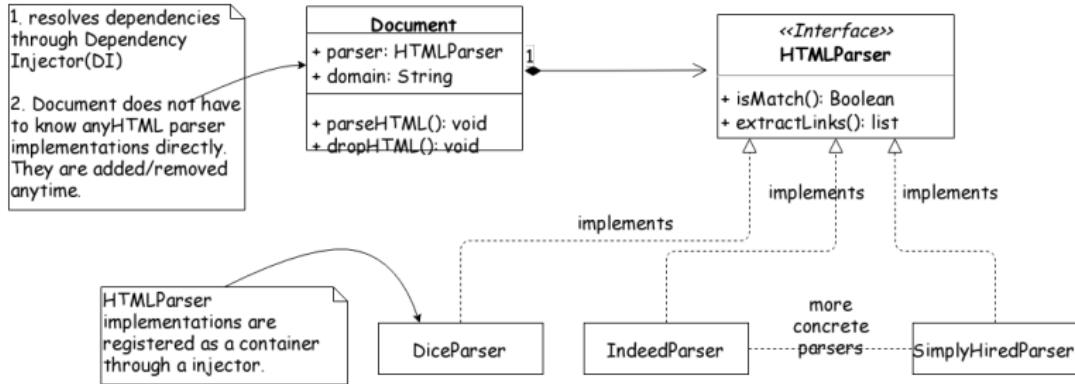
Whirlpool: Fetcher

Whirlpool: Fetcher



Whirlpool: Parser

Parser



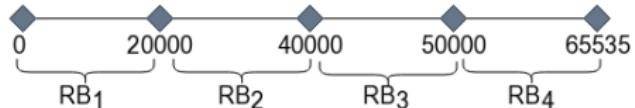
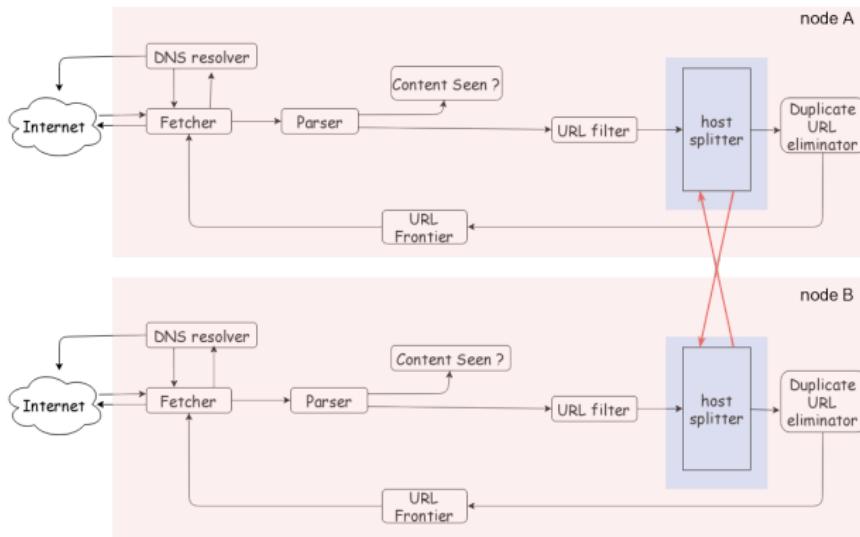
Whirlpool: Deduplication

Deduplication

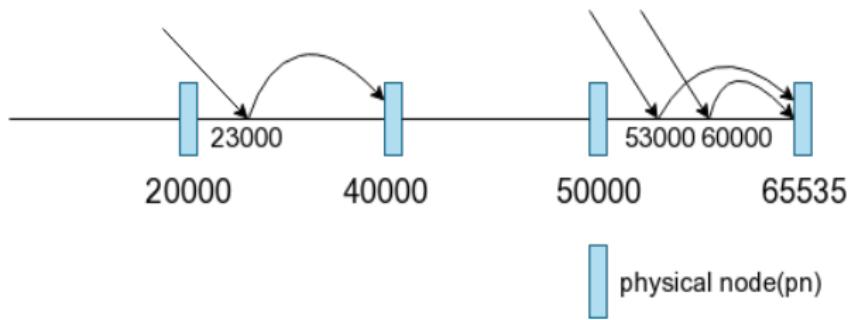
to add something

Whirlpool: Distributed Crawling

Host splitting in Mercator



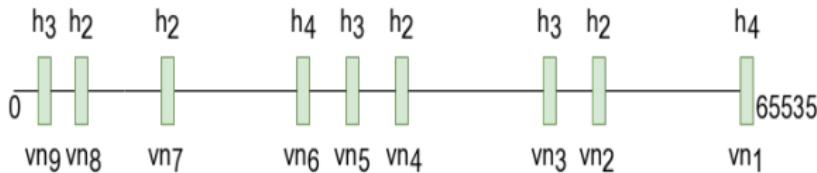
Rebalancing: mod N



Rebalancing: mod N

absolute URL	$m = \text{hash}(\text{absolute URL})$	$m \bmod N(1)$	$m \bmod N(2)$	$m \bmod N(3)$	$m \bmod N(4)$
jobs.com/browses	24	n1	n1 0	n1 0	n1 0
dice.com/browse-titles	59956	n1	n2 1	n1 1	n4 1
monster.com/categories	10090	n1	n1 0	n1 0	n3 1
simplyhired.com/find-jobs	23000	n1	n2 1	n2 0	n4 1
indeed.com/jobs/unix-admin	41103	n1	n2 1	n1 1	n4 1
glassdoor.com/hires/software-developer	35114	n1	n1 0	n3 1	n3 0
data movement		0%	50%	50%	60%

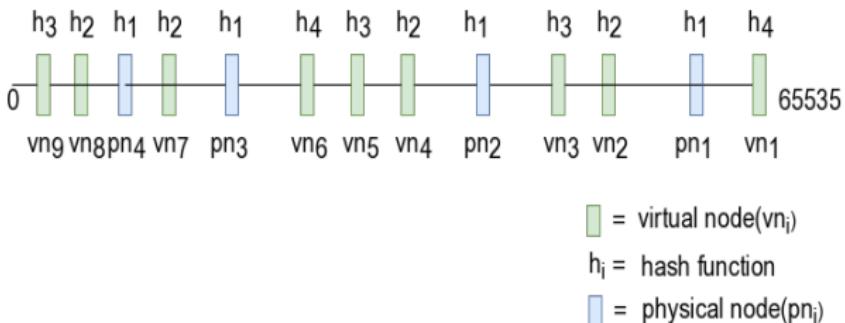
Rebalancing: virtual nodes



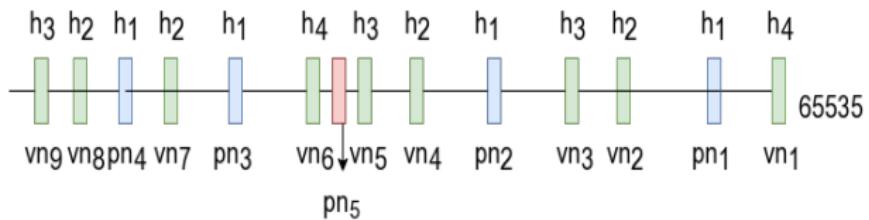
■ = virtual node(vn_i)

h_i = hash function

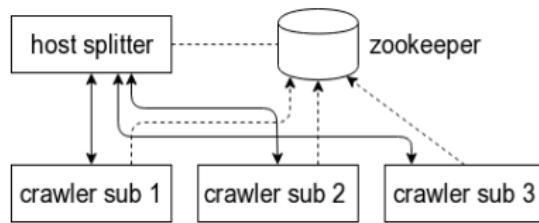
Rebalancing: virtual nodes



Rebalancing: virtual nodes



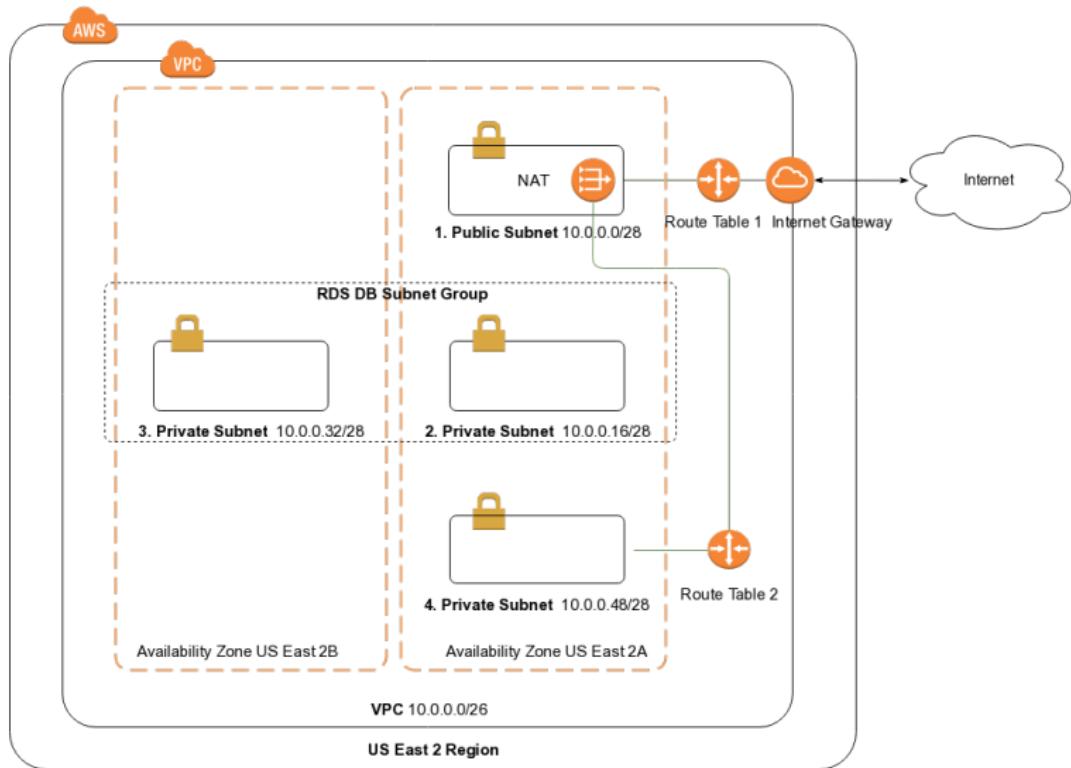
Zookeeper



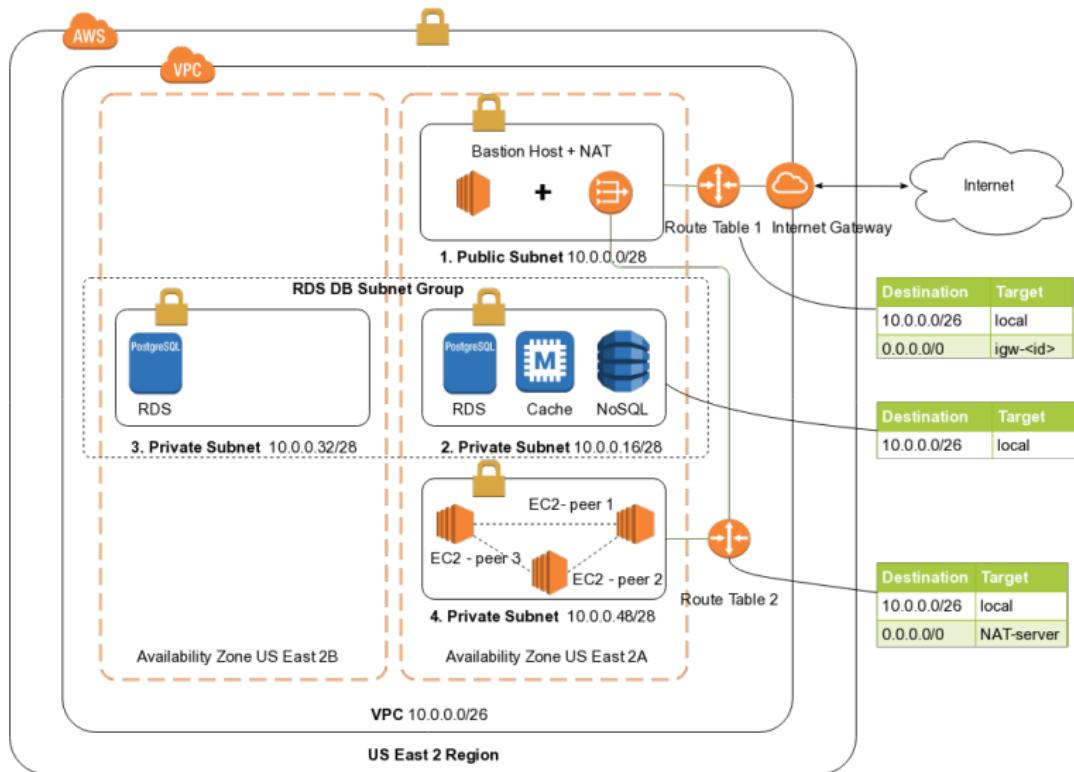
hash range	vnode(vni)	pnode(pni)	IP
0 - 3k	vn9	pn4	10.0.0.4
3k - 6k	vn8	pn4	10.0.0.4
6k - 10k	vn7	pn3	10.0.0.3
10k - 15k	vn6	pn2	10.0.0.2
15k - 20k	vn5	pn2	10.0.0.2
20k - 35k	vn4	pn2	10.0.0.2
35k - 45k	vn3	pn1	10.0.0.1
45k - 55k	vn2	pn1	10.0.0.1
55k - 65k	vn1	pn4	10.0.0.4

Whirlpool: Operations

From 10,000 ft.



From 5,000 ft.



Infrastructure as a code

Future work

future to do

to add something

Thank you! Questions ?