

# Whirlpool

## Data Acquisition using N-node Distributed Web Crawler

Rihan Pereira, MSCS

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

December 6, 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

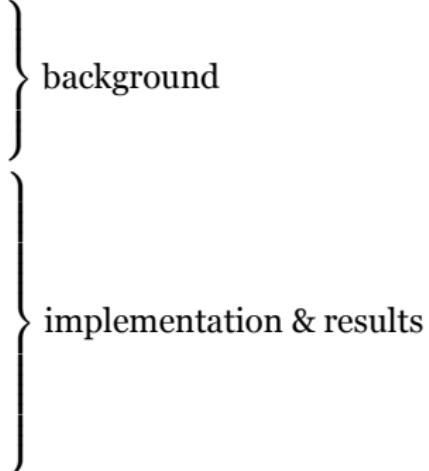
- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication

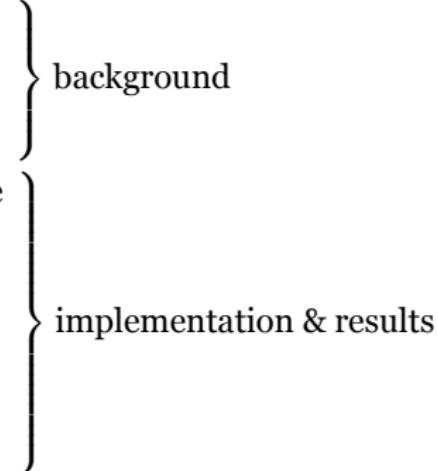


- Motivation & Contributions
- Crawler characteristics & history
- Mercator 1999 (Heydon & Najork)
- Software Design Principles
- Whirlpool: Event driven architecture
- Whirlpool: Fetcher
- Whirlpool: Parser
- Whirlpool: Deduplication
- Whirlpool: Distributed Crawling

}

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 is divided into two main sections by curly braces on the right side of the slide. The first section, containing the first five items, is labeled 'background'. The second section, containing the remaining five items, is labeled 'implementation & results'.

- 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

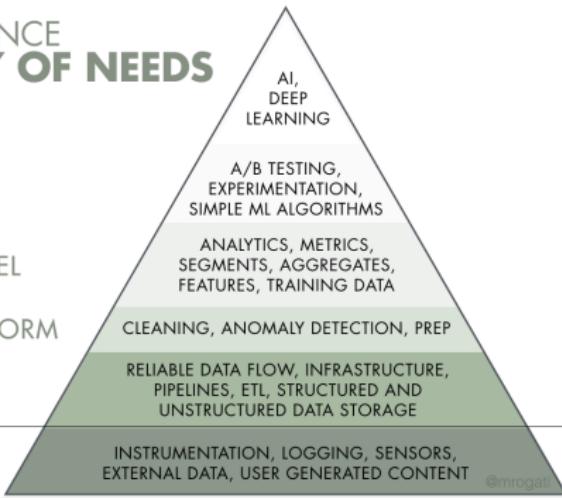
LEARN/OPTIMIZE

AGGREGATE/LABEL

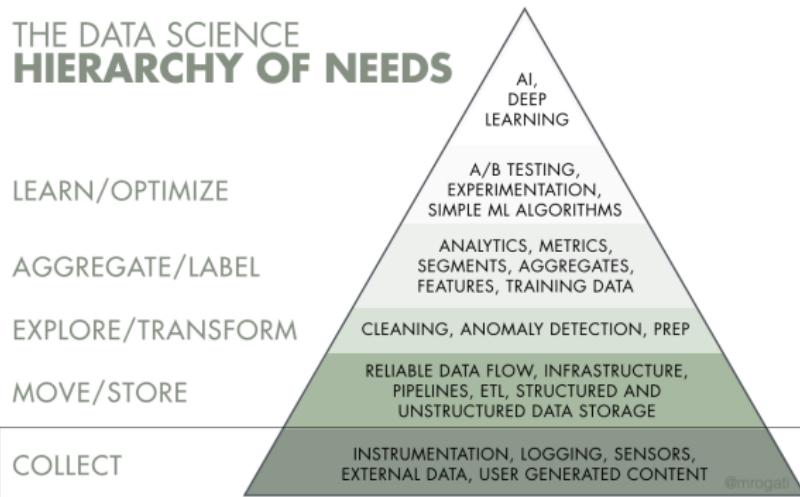
EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



# Motivation



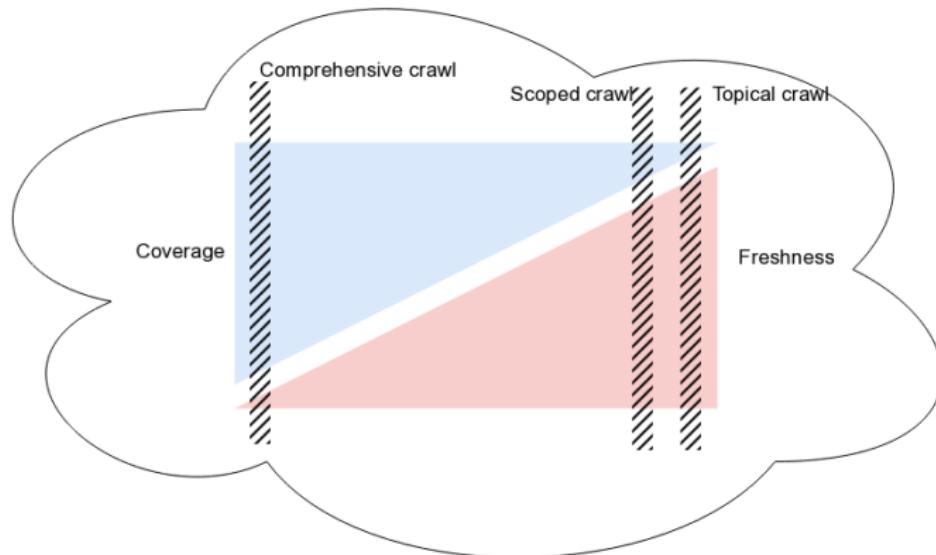
Self-actualization (AI) is great, but you first need food, water, and shelter (data literacy, collection, and infrastructure)."

# 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

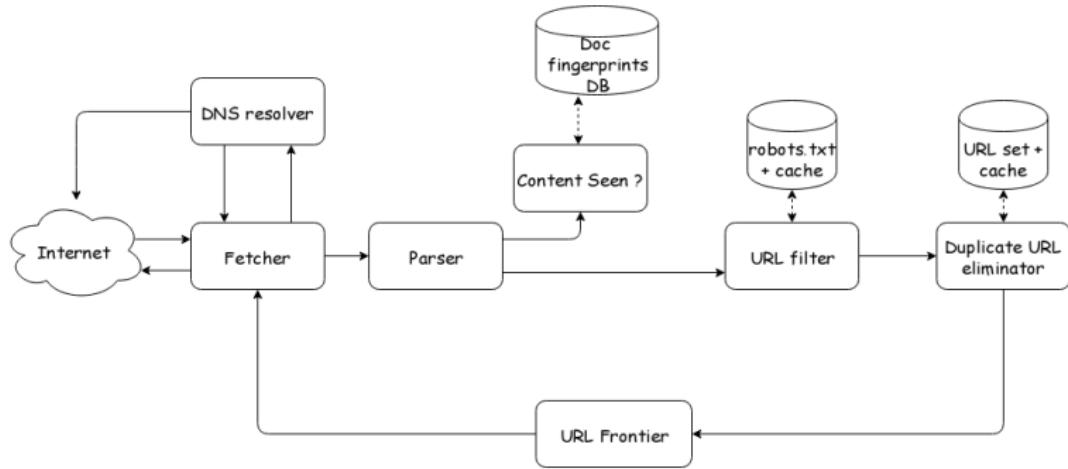


Figure: Mercator building blocks (Heydon & Najork)

# Mercator: URL filter

```
User-agent: *
Allow: /
Disallow: /unsubscribe
Disallow: /job/track-pageview
Disallow: /candidate/
Disallow: /alert-job/
Disallow: /facebook/
Disallow: /amazon/
Disallow: /job/track-apply-click
Disallow: /job/monster-apply-request
Disallow: /?search=
Disallow: /geo/job_title
Disallow: /blog/wp-admin/
Disallow: /contact/zip-resume
Disallow: /resume-database/preview-iframe
Disallow: /admin
Disallow: /login/ # want the main login page allowed, but not pages under /login/
Disallow: /mu
Disallow: /trk
Disallow: /record-event
Disallow: /events
Disallow: /apply/ # want apply allowed, but not pages under /apply/
Disallow: /eclk/
Disallow: /pixel/
Disallow: /chat/get-chat-auth-token
Disallow: /ajax/

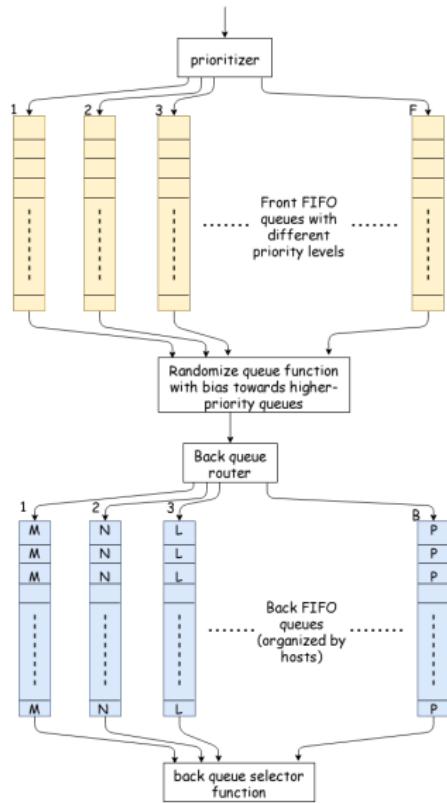
# Block URLs that are likely added by js clipboard library
Disallow: /jobs/*/closest$
Disallow: /jobs/*/clipboard-action$
Disallow: /jobs/*/is$

# Block temporary pages of the go seo app
Disallow: /svc/seo/

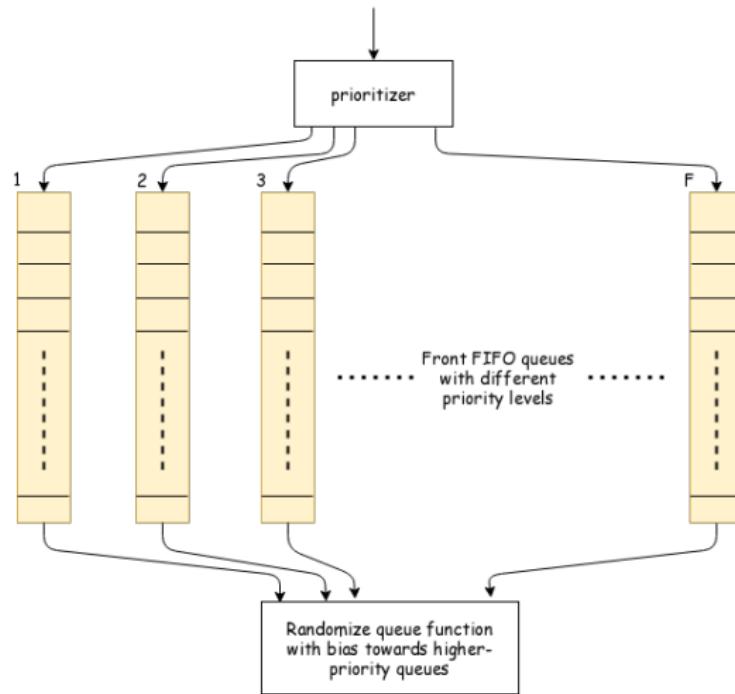
User-agent: 008
Disallow: /

User-agent: OmniExplorer_Bot
Disallow: /
```

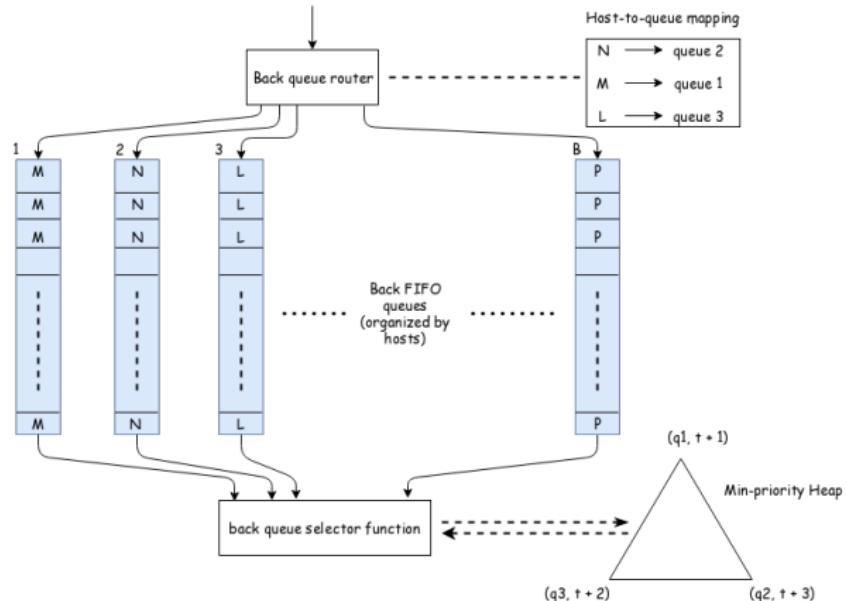
# URL Frontier Scheme



# Front queue (Frontier Queue)



# Back queue (Frontier Queue)



## Software Design Principles

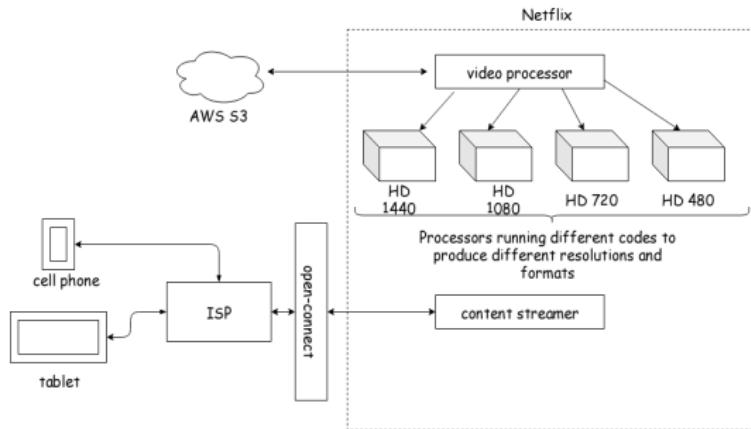
# Designing scalable systems

# Designing scalable systems

- Adding identical copies of components

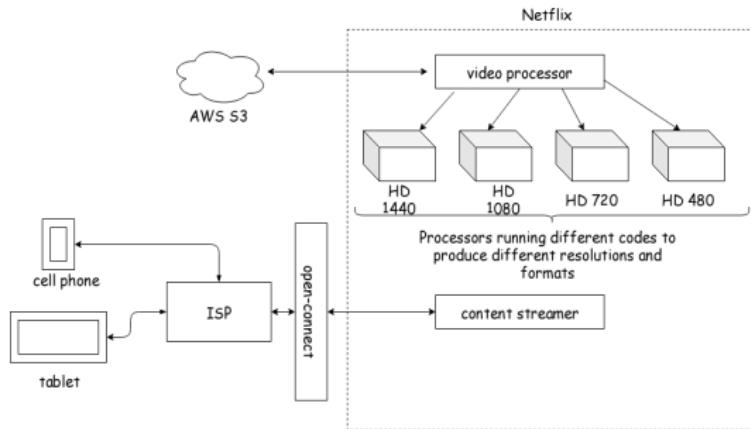
# Designing scalable systems

- Adding identical copies of components
- Functional partitioning



# Designing scalable systems

- Adding identical copies of components
- Functional partitioning



- Data partitioning

# State Management

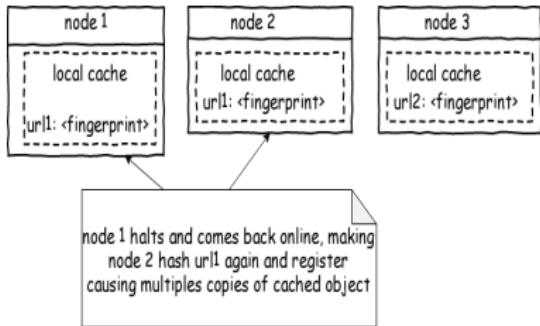


Figure: identical copies of same cached object

# State Management

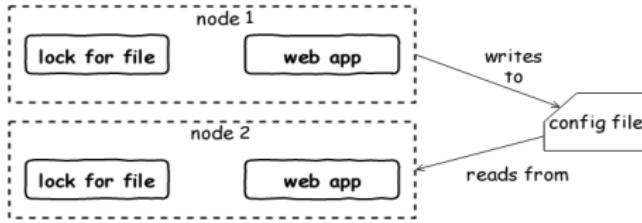


Figure: Using local locks to access shared resources

# State Management

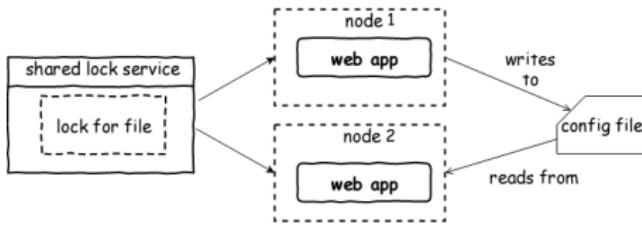


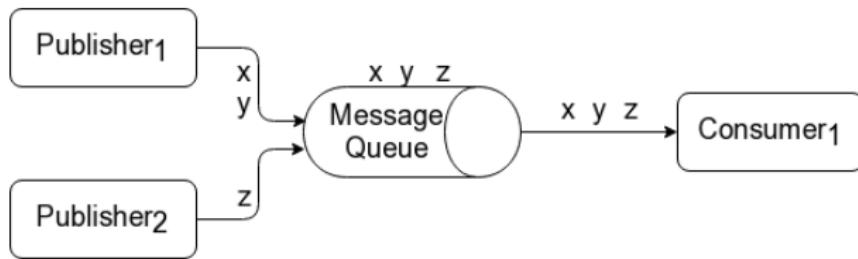
Figure: using shared locks to access shared resources

## Whirlpool: Event-driven architecture

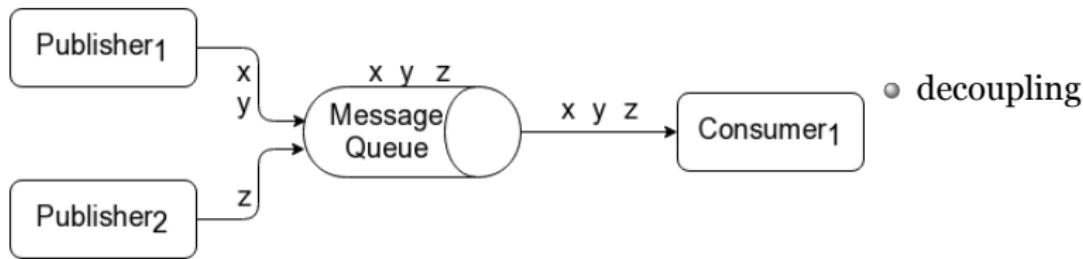
# Message buses



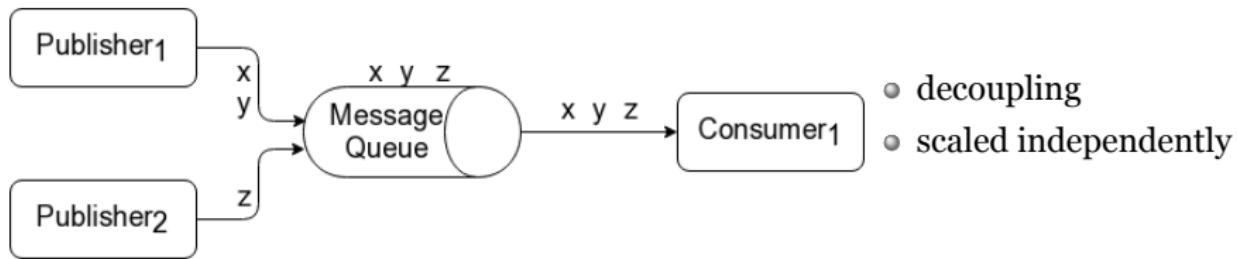
# Message Queue(MQ)



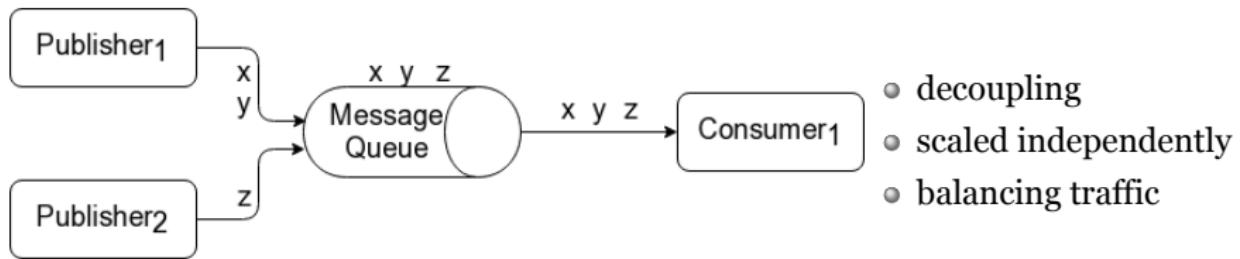
# Message Queue(MQ)



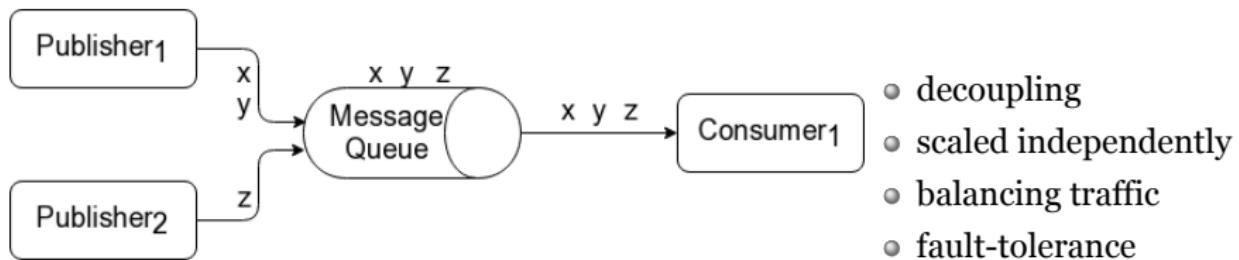
# Message Queue(MQ)



# 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

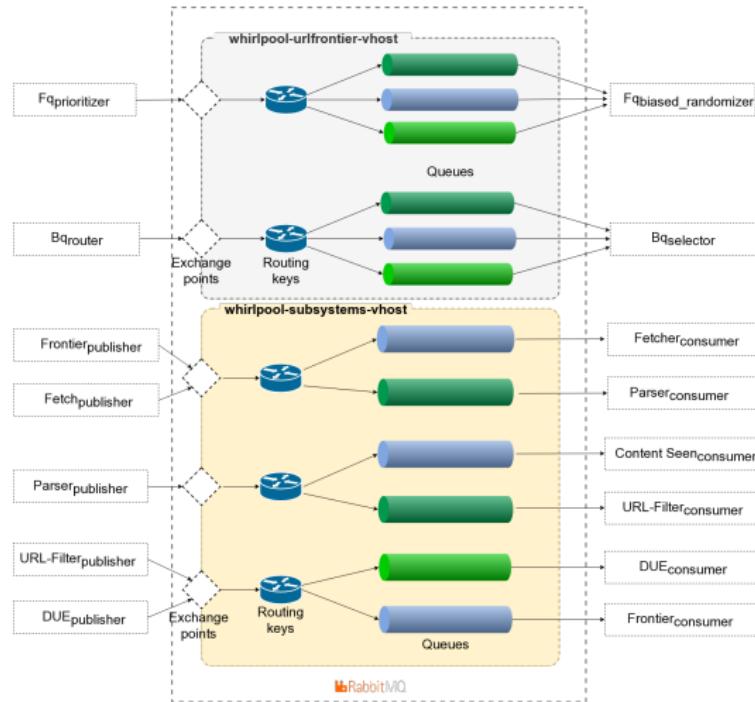
# MQ: Routing mechanisms

- Direct Worker Queue Data Flow
- Fanout
- Topic

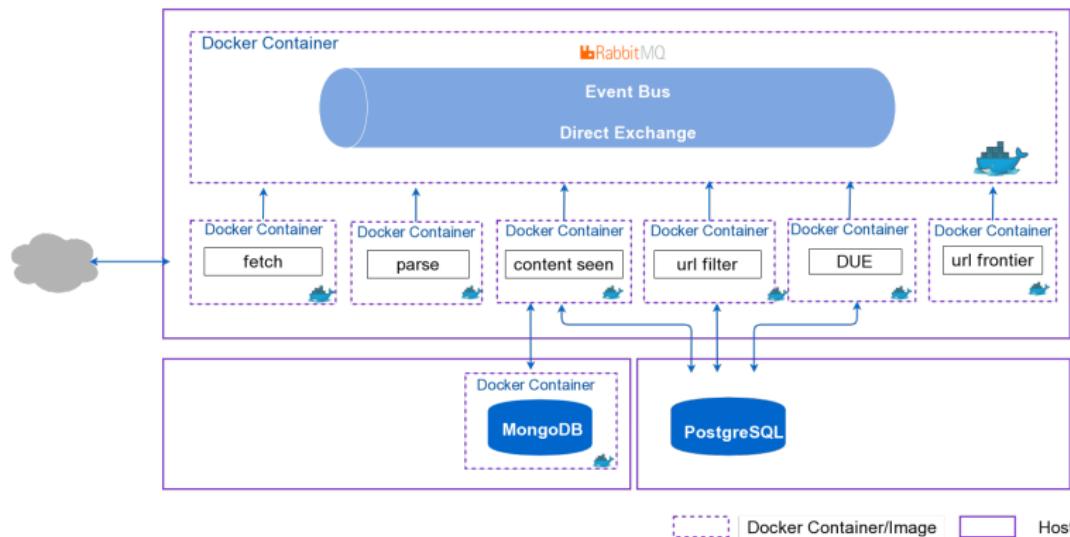
# MQ: Routing mechanisms

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

# Direct Worker Queue Data Flow



# RabbitMQ: Message bus



# Language runtime



# 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
```

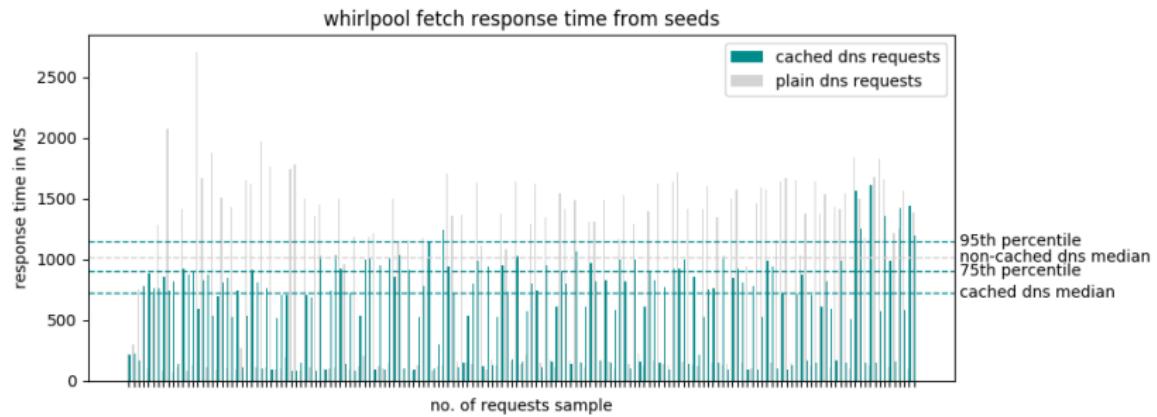
# development vs. production docker containers

```
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

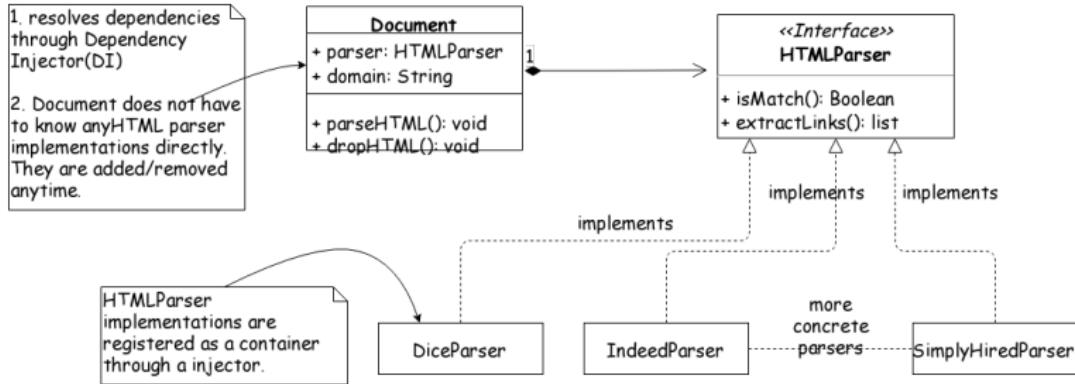
Whirlpool: Fetcher

# Whirlpool: Fetcher



## Whirlpool: Parser

# Parser



## Whirlpool: Deduplication

# Fingerprinting

<u><a href="#">Id</a></u> [PK] integer	<u><a href="#">page_fp</a></u> character varying (255)	<u><a href="#">page_type</a></u> character varying (255)	<u><a href="#">whirlpool_page_id</a></u> character varying (255)	<u><a href="#">domain</a></u> character varying (255)	<u><a href="#">fp_alg</a></u> character varying (255)	<u><a href="#">content_page_id</a></u> character varying (255)
24	b838fe385ab57d3atfe538ed...	nc	5d8d51539b96be004f5af9f4	www.simplyhired.com	sha1	5d8d5153436a1c002df172a0
25	def9e10d1c7ad7bf92465d41...	nc	5d8d57469b96be004f5af9f6	www.dice.com	sha1	5d8d5756a17e4d002da669c1
26	fa4cbe7184263bb6ee885638...	nc	5d8d57be9b96be004f5af9f8	www.dice.com	sha1	5d8d57bea17e4d002da669c3
27	d30ccf29df3fa34e6f562e428...	nc	5d8d57f89b96be004f5af9fa	www.dice.com	sha1	5d8d57f8a17e4d002da669c5
28	3dba207c69927ca4dc78d426...	nc	5d8d8d71db02b70052383520	job-openings.monster.com	sha1	5d8d8d71825679002db3cd2f
29	d4c714f73c767bcd021d8a1...	nc	5d8d8e36db02b70052383522	job-openings.monster.com	sha1	5d8d8e36825679002db3cd31
30	10ae4d73bbc028349caffac81...	nc	5d8d8e70db02b70052383524	job-openings.monster.com	sha1	5d8d8e70825679002db3cd33

# Bloom filter

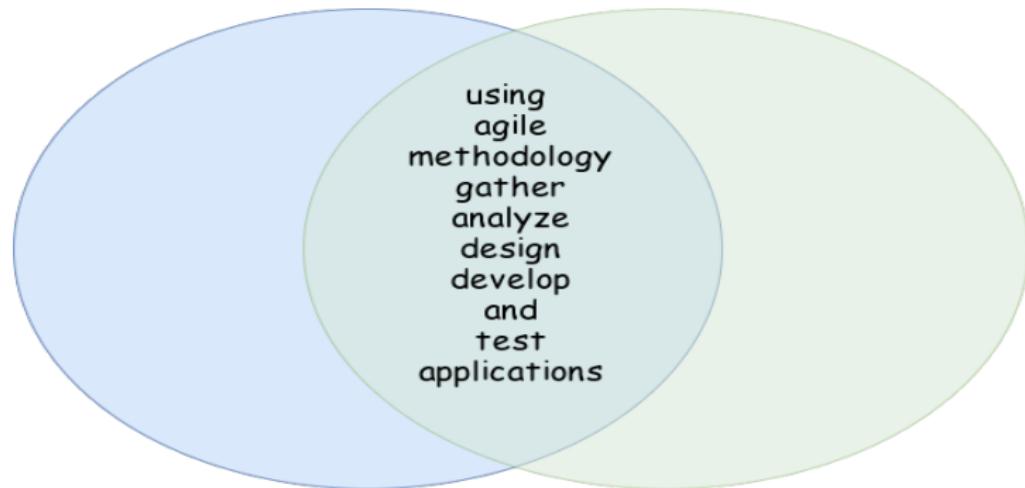
# Shingling: Near Deduplication

# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- using agile methodology gather analyze design develop and test applications

# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- using agile methodology gather analyze design develop and test applications

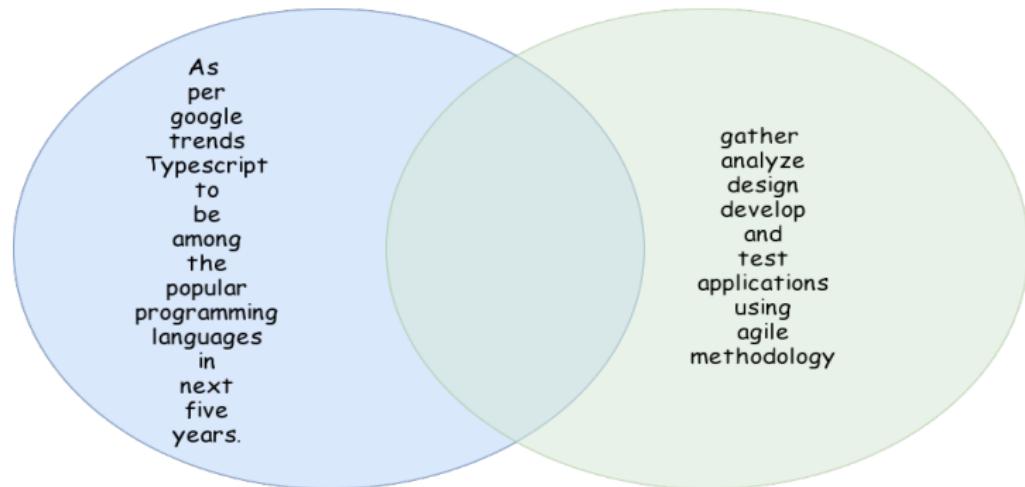


# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- As per google trends Typescript to be among the popular programming languages in the next five years

# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- As per google trends Typescript to be among the popular programming languages in the next five years

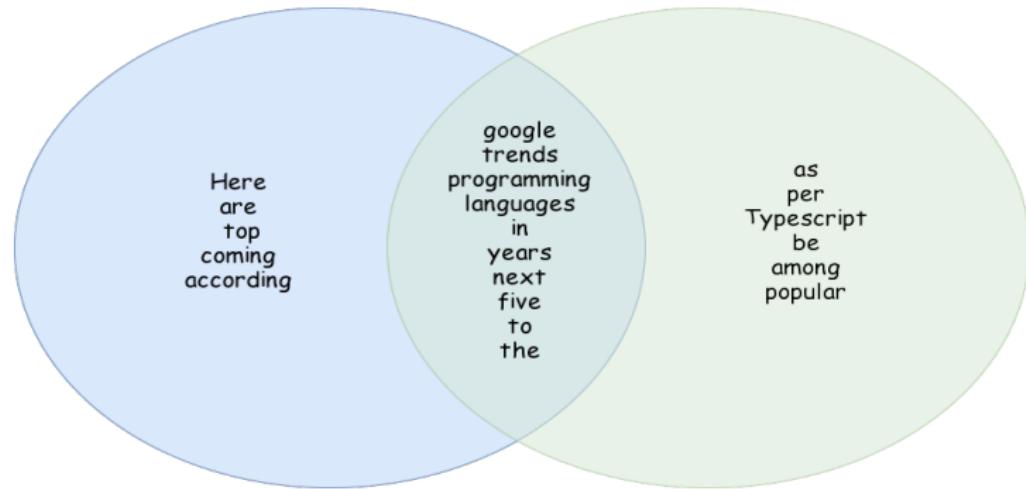


# Shingling: Near Deduplication

- Here are the next top five programming languages in coming years according to google trends
- As per google trends Typescript to be among the popular programming languages in the next five years

# Shingling: Near Deduplication

- Here are the next top five programming languages in coming years according to google trends
- As per google trends Typescript to be among the popular programming languages in the next five years

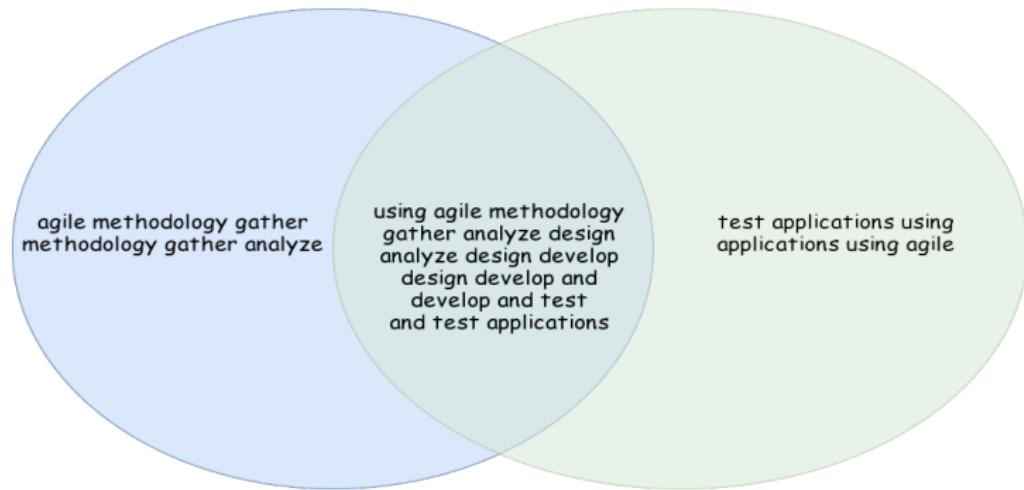


# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- using agile methodology gather analyze design develop and test applications

# Shingling: Near Deduplication

- gather analyze design develop and test applications using agile methodology
- using agile methodology gather analyze design develop and test applications

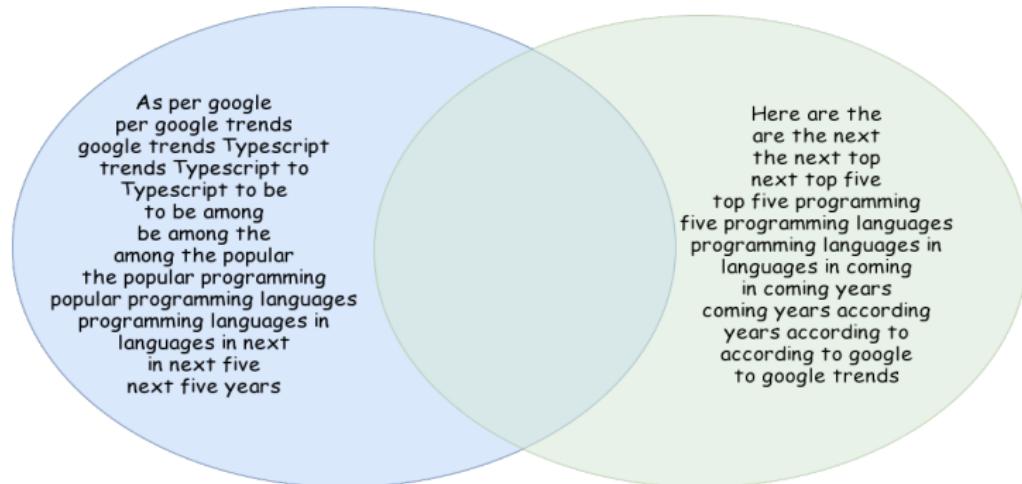


# Shingling: Near Deduplication

- Here are the next top five programming languages in coming years according to google trends
- As per google trends Typescript to be among the popular programming languages in the next five years

# Shingling: Near Deduplication

- Here are the next top five programming languages in coming years according to google trends
- As per google trends Typescript to be among the popular programming languages in the next five years



# Shingling: Near Deduplication

## monster.com

simhash_a	simhash_b	jaccard index	hamming weight
-1571308417	1314365132	0.04	0.52
1314365132	-1906275206	0.15	0.39
-1571308417	-1906275206	0.03	0.48

## simplyhired.com

simhash_a	simhash_b	jaccard index	hamming weight
25104594	-897371832	0.1	0.19
-897371832	-880463811	0.72	0.5
25104594	-880463811	0.01	0.22

## dice.com

simhash_a	simhash_b	jaccard index	hamming weight
-656343461	-660545702	0.9	0.8
-660545702	1749346824	0.1	0.35
-656343461	1749346824	0.1	0.43

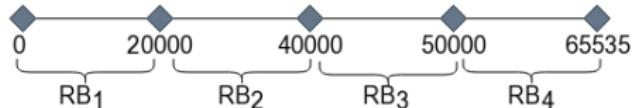
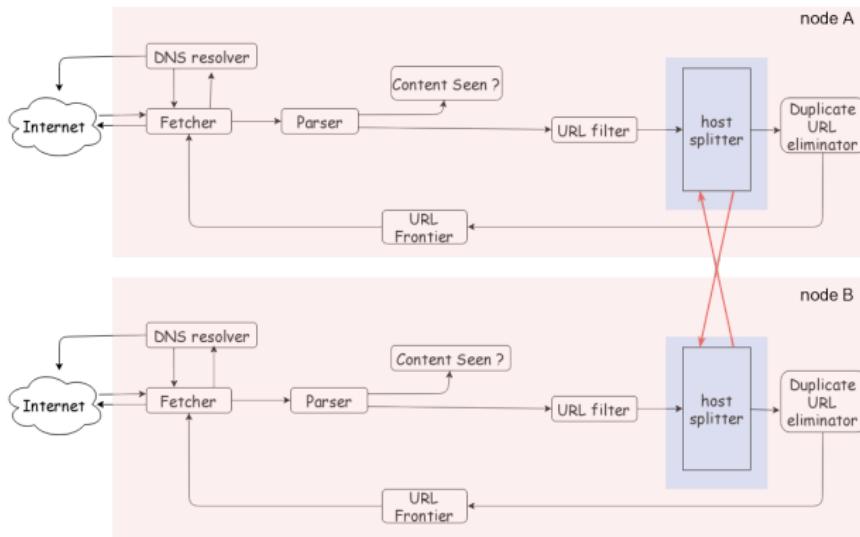
# Shingling: Near Deduplication

<b>id</b> [PK] integer	<b>page_fp</b> character varying (255)	<b>page_type</b> character varying (255)	<b>whirlpool_page_id</b> character varying (255)	<b>domain</b> character varying (255)	<b>fp_alg</b> character varying (255)	<b>content_page_id</b> character varying (255)
32	-656343461	nc	5da2785cd410370044b269a4	www.dice.com	simhash	5da27972d410370044b269a5
33	-660545702	nc	5da28b327f2d693fe92077aa	www.dice.com	simhash	[null]
34	1749346824	nc	5da28e0c1244bc30edf4026f	www.dice.com	simhash	5da28e0c1244bcd4eef40270
35	25104594	nc	5da2961ef31b193ffbb426bca	www.simplyhired.com	simhash	5da2961ef31b194674426bcb
36	-897371832	nc	5da2966b3c81d51360729b22	www.simplyhired.com	simhash	5da2966b3c81d52104729b23
37	-880463811	nc	5da296cf153a6c17cac78e8f	www.simplyhired.com	simhash	[null]
38	-1571308417	nc	5da298c6dd568b9a7fa681dd	www.jobs-openings.monster.c...	simhash	5da298c6dd568b87ba681de
39	1314365132	nc	5da299011cd38e47840c4f50	www.jobs-openings.monster.c...	simhash	5da299011cd38e27210c4f51
40	-1906275206	nc	5da29924d66cd79ec310c1da	www.jobs-openings.monster.c...	simhash	5da29924d66cd783c610c1db

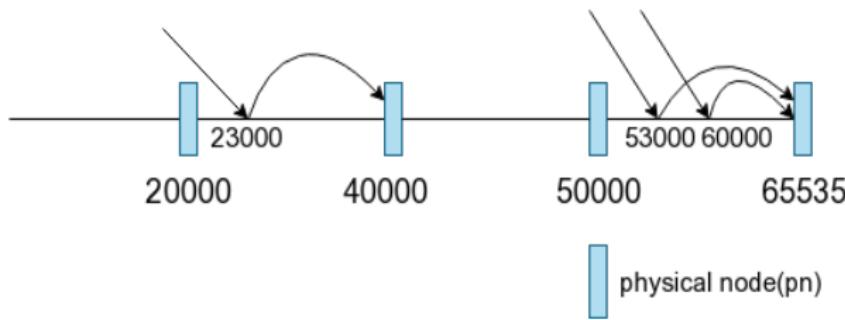
# Locality Sensitive Hashing

## Whirlpool: Distributed Crawling

# Host splitting in Mercator



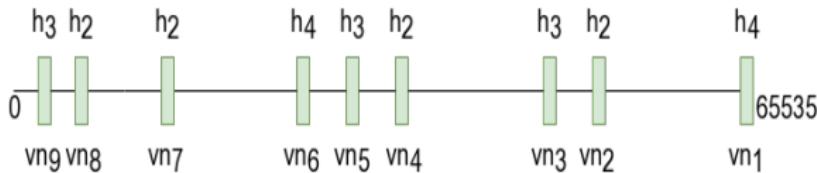
# Rebalancing: mod N



# 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%

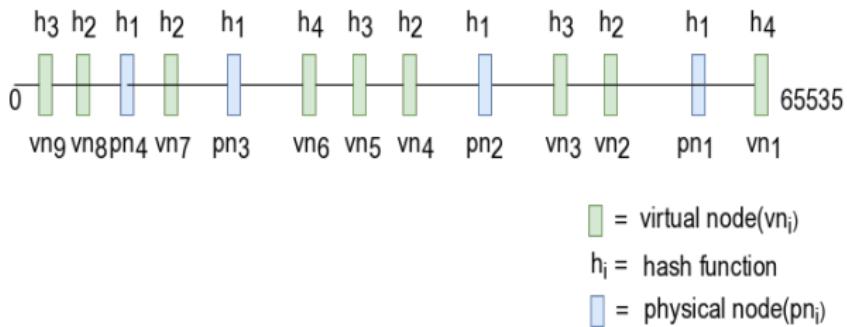
# Rebalancing: virtual nodes



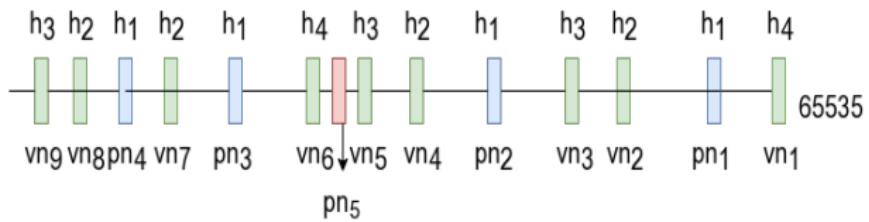
■ = virtual node( $vn_i$ )

$h_i$  = hash function

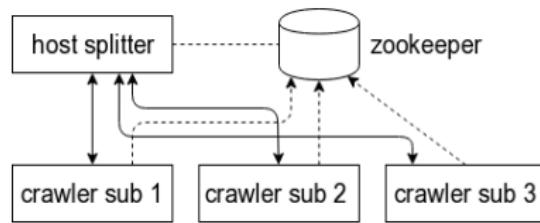
# Rebalancing: virtual nodes



# Rebalancing: virtual nodes



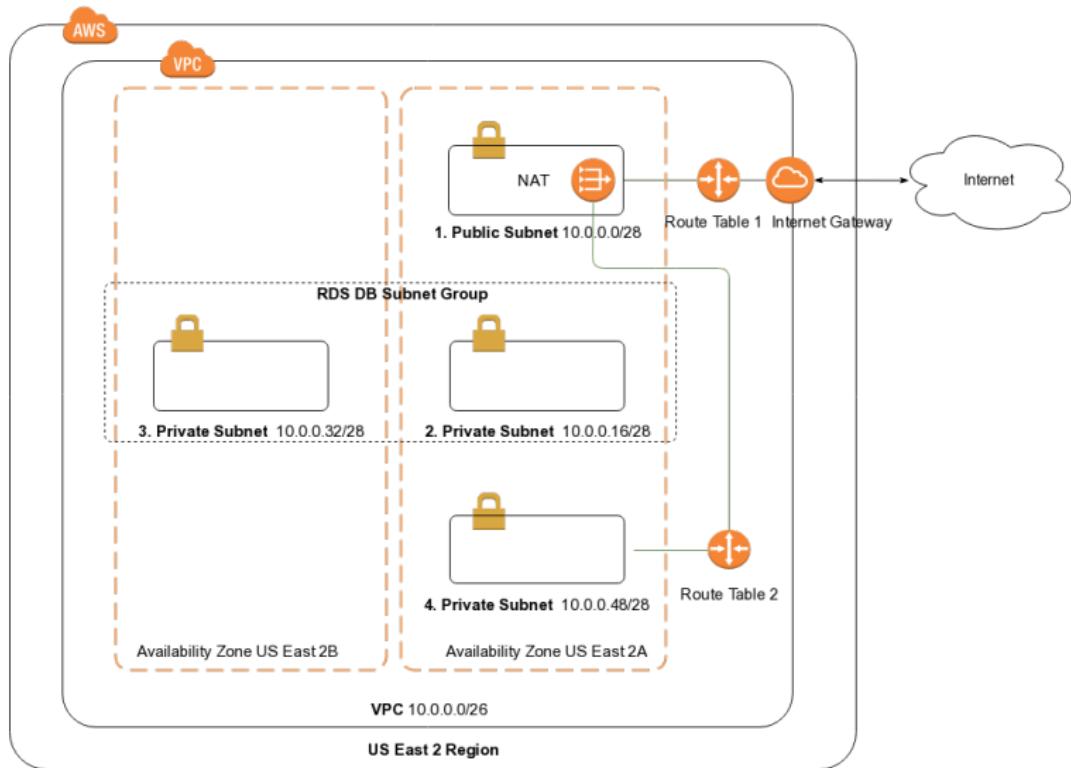
# Zookeeper



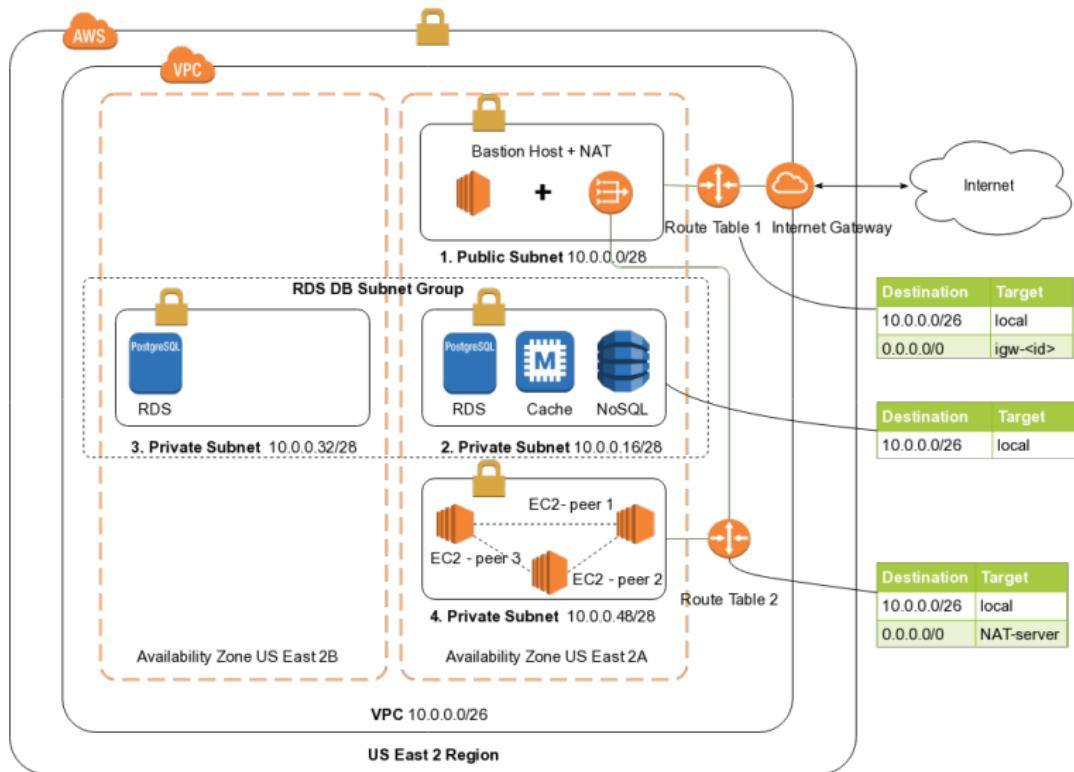
hash range	vnode(vni)	pnode(pni)	IP
0 - 3k	vn9	pn4	10.0.0.4
3k - 6k	vn8	pn4	10.0.0.4
6k - 10k	vn7	pn3	10.0.0.3
10k - 15k	vn6	pn2	10.0.0.2
15k - 20k	vn5	pn2	10.0.0.2
20k - 35k	vn4	pn2	10.0.0.2
35k - 45k	vn3	pn1	10.0.0.1
45k - 55k	vn2	pn1	10.0.0.1
55k - 65k	vn1	pn4	10.0.0.4

## Whirlpool: Operations

# From 10,000 ft.



# From 5,000 ft.



# Infrastructure as a code

# Recap

- Microservices
- Load balancing
- Decoupling
- Extensible
- Distributed Systems(Partitioning)
- Logging & metrics calculation
- Fault tolerance
- Resiliency

## Future work

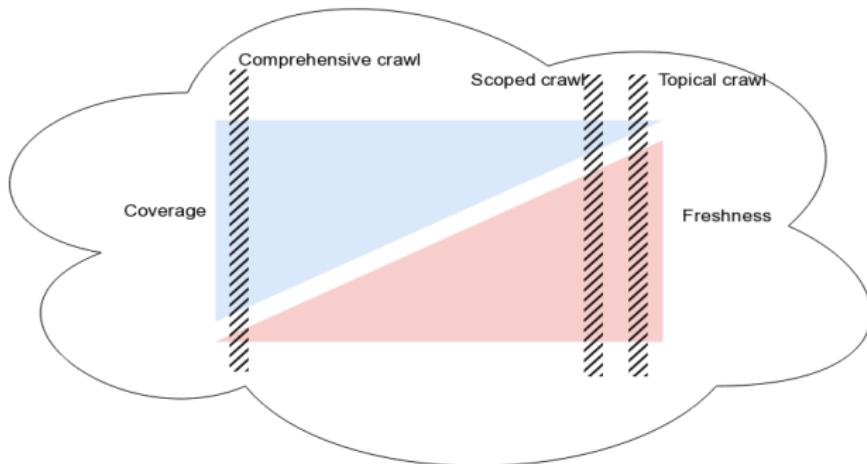
# future work

Indexing crawled data & logs using Elasticsearch



# future work

## Adding comprehensive coverage



# future work

## Content Extraction

# future work

Categorizing Jobs using ML algorithm

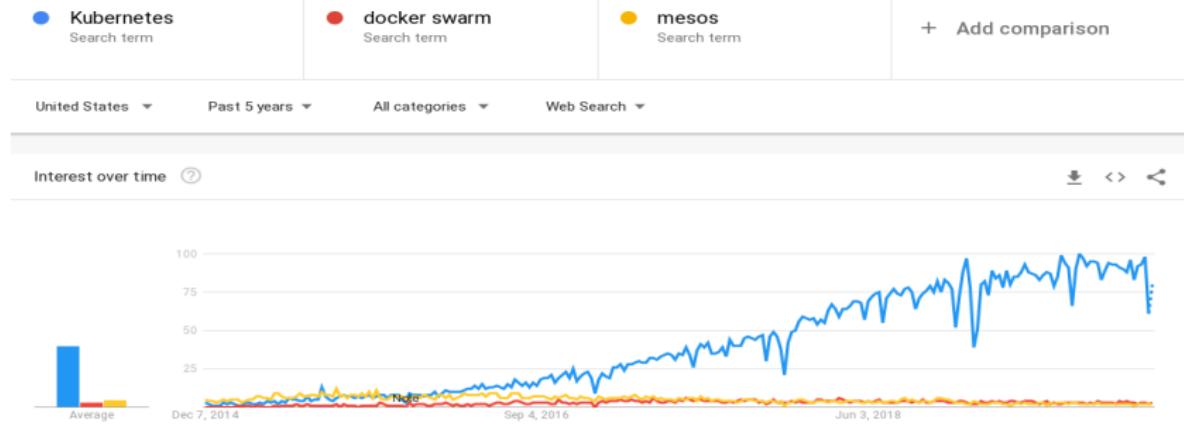
# future work

## SeenTest: Inserting and Querying Simhashes using LSH

# future work



## Deploying whirlpool using kubernetes



Thank you! Questions ?