

Introduction au **W**ebscraping

Léa Christophe
Robin Cura
Hugues Pecout
Alexandre Cebeillac
Sébastien Rey-Coyrehourcq

GT Webscraping
21 juin 2024

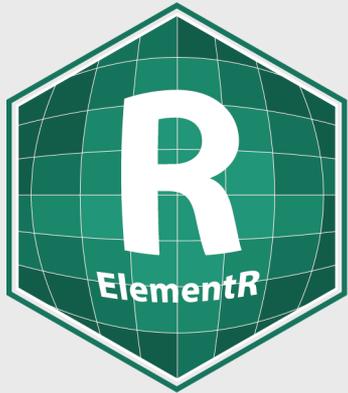


Géographie-cités
UMR 8504

ides
UMR6266 CNRS Rouen



UMR 8586 (CNRS)



W**e**bscraping

1. Introduction (5 min)
2. Exemples d'application (10-15 min)
3. Mise en pratique

Introduction

Définition(s)

Web scraping, web harvesting, or web data extraction is data scraping used for **extracting data from websites**. Web scraping **software** may directly access the World Wide Web using the Hypertext Transfer Protocol or a web browser. While web scraping can be done manually by a software user, **the term typically refers to automated processes implemented using a bot or web crawler**.

Wikipedia (en)

“ Le WebScraping est **une pratique** qui nécessite l’articulation plus ou moins complexe de méthodes, de logiciels et d’infrastructures pour **naviguer et extraire du contenu** (des informations) d’un ou de plusieurs sites Web **de façon automatisé**.”

notre GT (wip)

Introduction

De multiples utilisations pour les SHS ...

- Constitution de corpus ...
 - **en // de corpus existants**
 - **originaux**
- ... faits de de sources et de nature de sources variés (son, images, textes)
- ... collectés avec un
 - **flux continu**
 - **échantillonnage**

études comparatives
sauvegarde / patrimoine
mondialisé

suivi longitudinal
mixte quali/quantitatif

► Intérêt **croissant** au sein des laboratoires en SHS ...

Introduction

... mais des verrous à différents niveaux

- **(I) légalité** de la :

- collecte
- stockage
- de la (re)diffusion



GT Webscraping
(partie éthique & légalité)

- **(II) contrôle qualité**

- validation
- harmonisation



trop problème dépendant

- **(III) techniques** lié à l'objet à capturer :

- résilience
- intégration
- reproductibilité



GT Webscraping
(docker, sgbd, forge, etc.)

Exemple : <https://journals.openedition.org/cybergeos/36478>

Discuté ici (vidéo à venir) : <https://atelierdatapaperadn.sciencesconf.org/>

Légalité

une multiplicité de facteurs interdépendants

(... +/- loi du plus fort à l'international entre multinationales ...)

- (1) Pays → Jurisprudences & droit d'auteurs
- (2) CGU / Robots.txt → Loi Nationale & Européenne
- (3) Nature de l'information → Personnelles ? Sensibles ?
- (4) Visibilité de l'information → Publiques ? Privés ?
- (5) Qui récolte l'information → Institutions ? Chercheurs ? Entreprises ?
- (6) Volumétrie → Substantielle ? Non Substantielle ?
- (7) Finalité de la récolte → Lucratif / Non lucratif ?
Avec ou Sans modifications ?
Traitements, Aggrégation ?

Quelques exemples



```

{
  "b": {
    "b_ranking_breakdown": {},
    "b_behaviors": {},
    "b_recommended_checkout": null,
    "b_recommended_checkout": null,
    "b_accommodation_classification_rating_data": {},
    "b_preferred": 1,
    "b_id": "8792793",
    "b_url": "/hotel/fr/les-toits-de-r-rooms1-group_adults=2",
    "b_url": "/hotel/fr/les-toits-de-r-rooms1-group_adults=2",
    "b_backend_rank": 1,
    "b_latitude": "49.484881",
    "b_hotel_ways_sold": 0,
    "b_cst": "FR",
    "b_max_discount_rate_badge": 0,
    "b_class_badge": "",
    "b_image_url": "/xdata/images/hotel/270x_3976x13099fc04138c436w",
    "b_bh_quality_class": 3,
    "b_booking_ways": "Superbe",
    "b_status": {},
    "b_hotel_title": "Les toits de Rouen - Atypique et familial",
    "b_sw_template": "atlas_sw_hotel",
    "b_image_url_v2": "/xdata/images/hotel/270x_3976x13099fc04138c436w",
    "b_class": 0,
    "b_is_booking_home": 1,
    "b_review_nr": 42,
    "b_accommodation_type": "Appartement",
    "b_is_legal_exception_applied": 1,
    "b_is_b_name": 1,
    "b_class_is_estimated": 0,
    "b_score_from_text": "Note sur span=strong c_commentaires&bbp-c/span",
    "b_sas_themes": {},
    "b_is_preferred_plus": 0,
    "b_is_single_unit_url": 0,
    "b_marker_type": "hotel",
    "b_nr_reviews_text": "Note sur strong=42 commentaires/strong",
    "b_accommodation_type_id": 201,
    "b_review_score": "4.7",
    "b_is_booking_home_v2": 1,
    "b_longitude": "1.896518",
    "b_latitude": "49.484881"
  }
}
```

Quelques exemples



Pratiques touristiques

Airbnb vs Booking.com

Flux multi-échelles

Données Avion FlightRadar24 & Bateaux MarineTraffic



flightradar



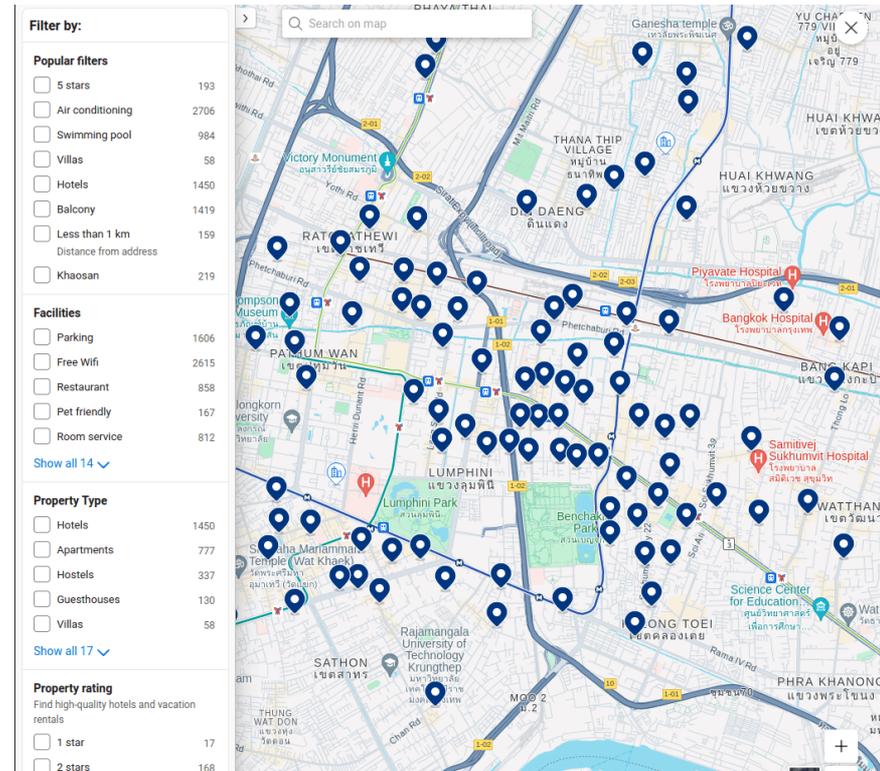
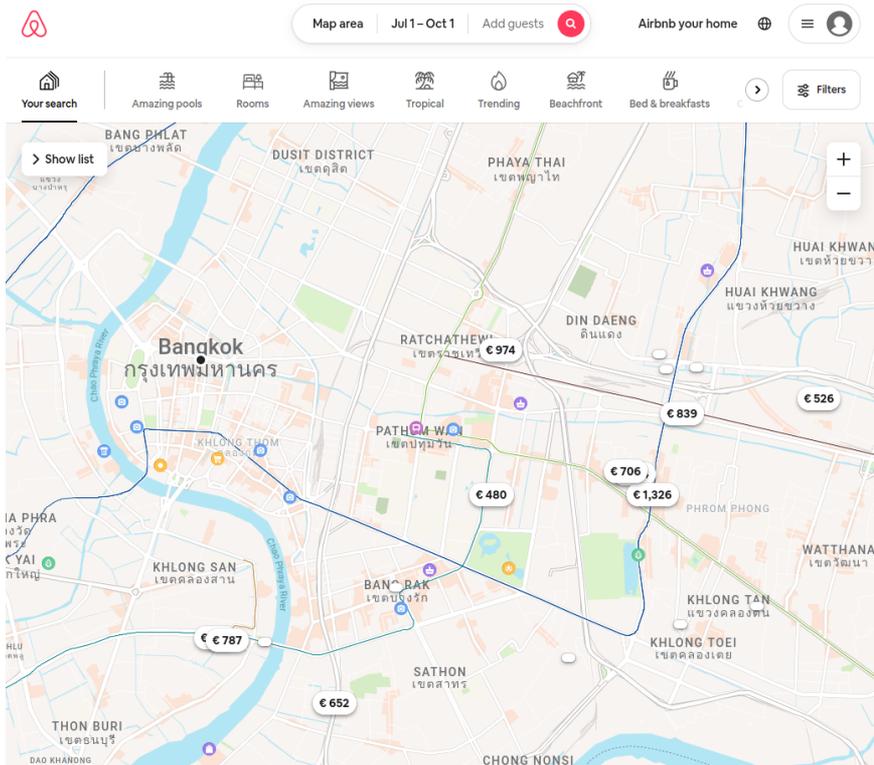
Mobilités quotidiennes & évacuation

Données trafic Bing & TomTom

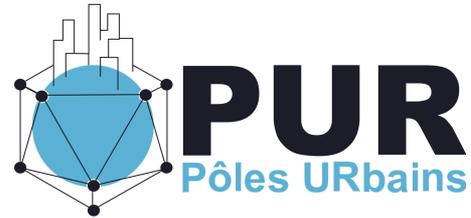
Pratiques touristiques



Comparer les localisations des types d'offres dans les villes



Arctique & système monde



ANR PUR (2016 - 2020) / Yvette Vaguet

Intégration des zones arctiques dans le système monde

Un des objectifs :

Analyser temporellement les flux affectant les zones arctiques, dont les flux touristiques.

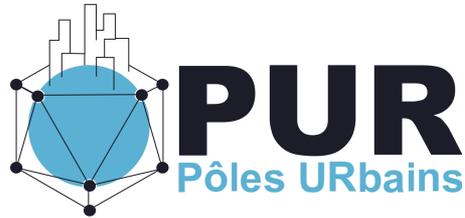
Pourquoi ?

"Tourism is a powerful global force that is turning the north into a Global North—tied into the Global South by facing the same global problem" (Veijola and Strauss-Mazzullo, 2019)

Comment ?

En créant des réseaux à partir des données les plus appropriées et accessibles

Arctique & système monde



ANR PUR (2016 - 2020) / Yvette Vaguet

Intégration des zones arctiques dans le système monde

Un des objectifs :

Analyser temporellement les flux affectant les zones arctiques, dont les flux touristiques.

Pourquoi ?

"Tourism is a powerful global force that is turning the north into a Global North—tied into the Global South by facing the same global problem" (Veijola and Strauss-Mazzullo, 2019)

Comment ?

En créant des réseaux à partir des données **7M+** les plus appropriées et accessibles

Airbnb listings worldwide

100K+

cities with Airbnb listings

191+

countries and regions with Airbnb listings



airbnb

Type

Privately held company

Industry

Lodging

Founded

August 2008;
11 years ago in San Francisco, California

Founders

Brian Chesky
Joe Gebbia
Nathan Blecharczyk

Airbnb



Collecter des données permettant de former un réseau de flux touristiques vers l'Arctique

Site web : airbnb.com

Langage : R

Framework : rvest / httr

Type de récolte : Campagne

Format : json

Caractéristiques :

- Accessible via code source
- (+) Plusieurs API officielles & pas de token privé
- (-) 80 requêtes / min / IP ==> temps de collecte long

- Des extractions de zones (insideairbnb) mais que dans certaines villes
- Des codes sont disponibles (Tom Slee) python, mais pas sous R

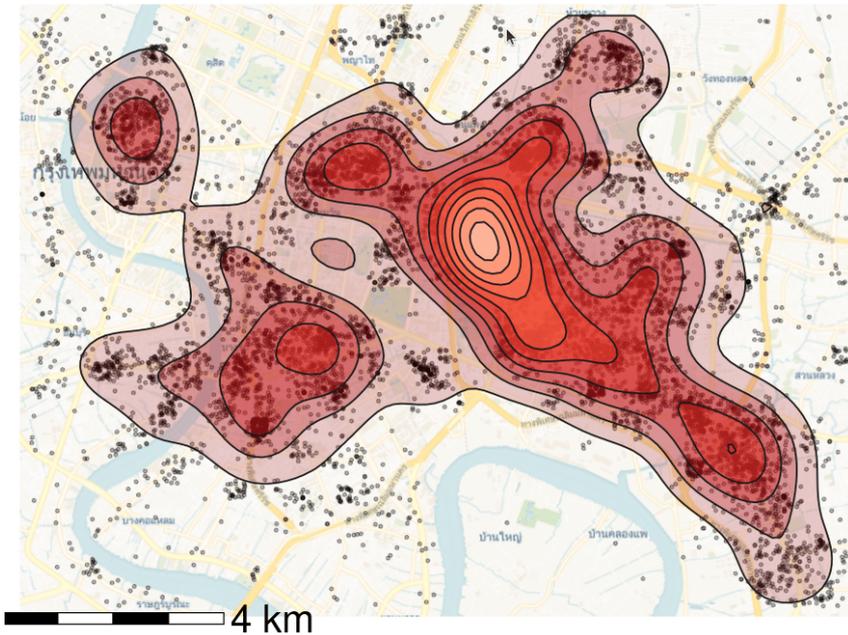
Pratiques touristiques



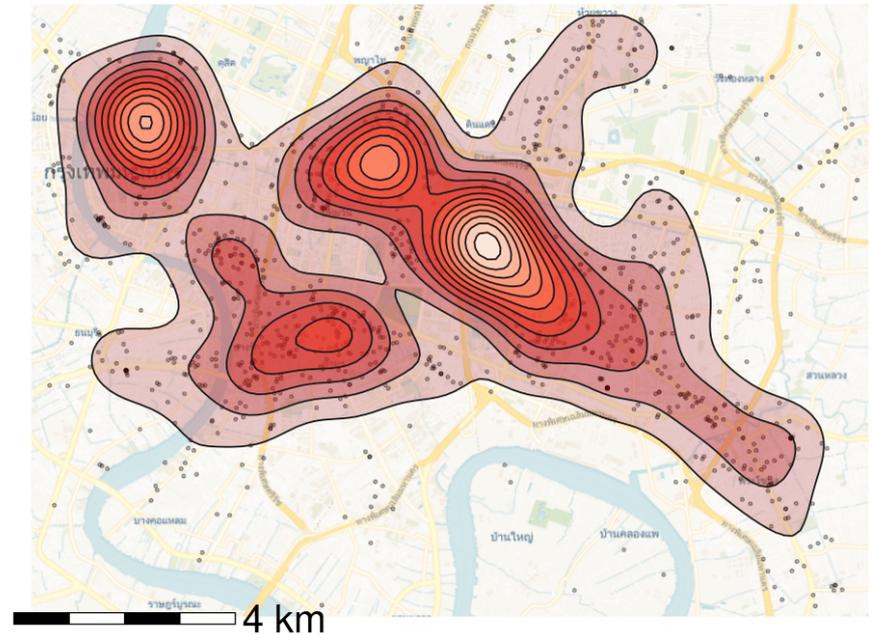
Comparer les localisations des types d'offres dans les villes



Airbnb

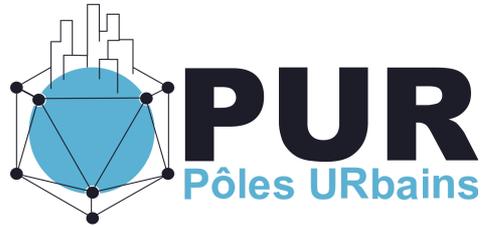


Booking



Localisation des hébergements à Bangkok

Arctique & système monde



ANR PUR (2016 - 2020) / Yvette Vaguet

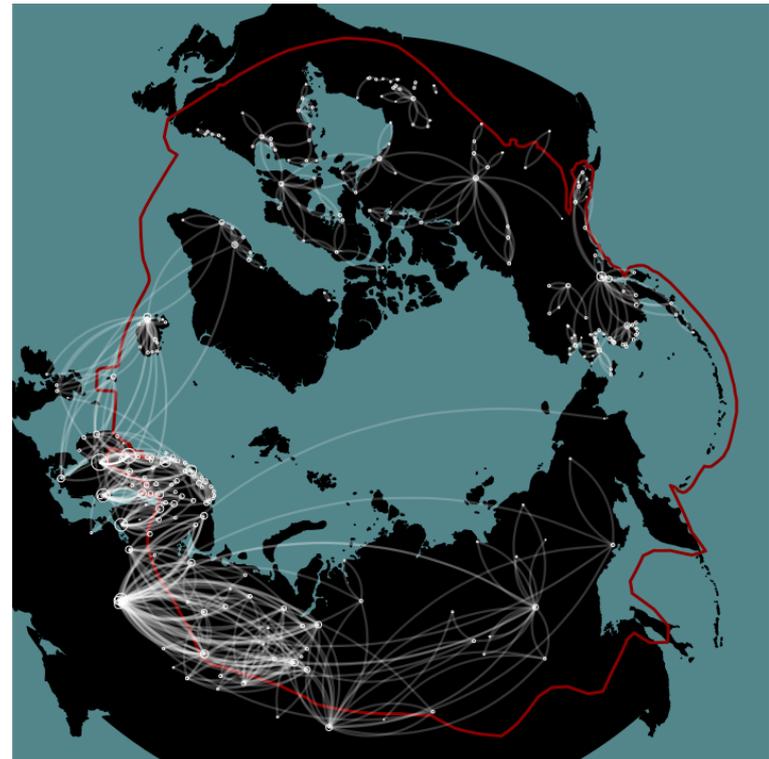
Intégration des zones urbaines arctiques dans le système monde

Un des objectifs :

Analyser temporellement les flux affectant les zones arctiques, notamment les flux aériens.

Comment ?

En créant des réseaux à partir des données les plus appropriées et accessibles



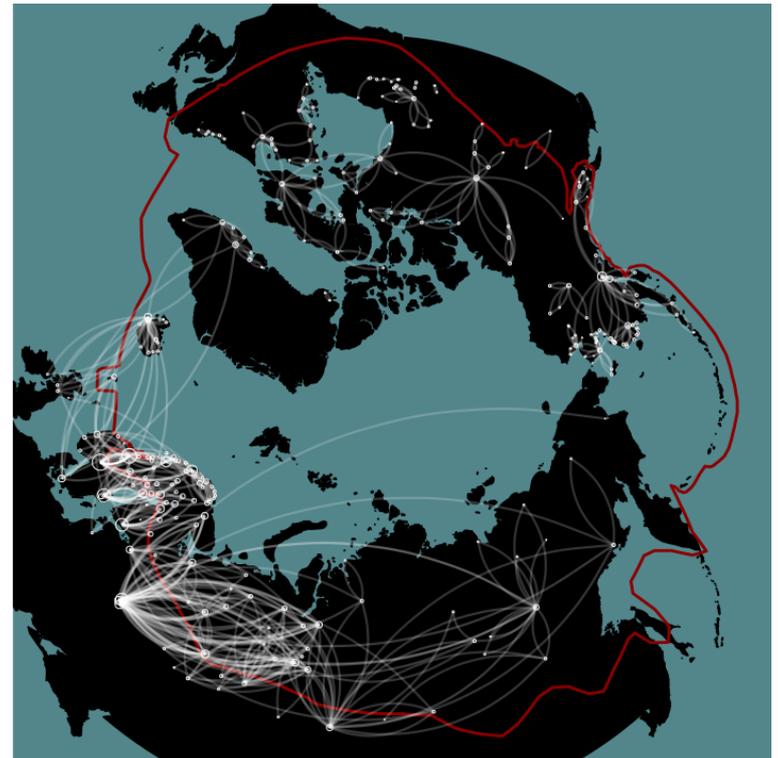
FlightRadar



flightradar 24

ANR PUR (2015 - 2019) / Yvette Vaguet

Récupérer les **départs et arrivées d'avions** pour les **aéroports mondiaux** pendant **plus d'un an**



FlightRadar



flightradar 24

ANR PUR (2015 - 2019) / Yvette Vaguet

Récupérer les **départs et arrivées d'avions** pour les **aéroports mondiaux** pendant **plus d'un an**

Site web : flightradar24.com

Langage : Python

Framework : Scrapy

Type de récolte : Continue (depuis le **2018/06/10**)

Format : mongoDB / json

Difficultés :

- architecture technique pour une récolte continue
- gestion des timezones
- intégration/consolidation des données SGBD (work in progress)



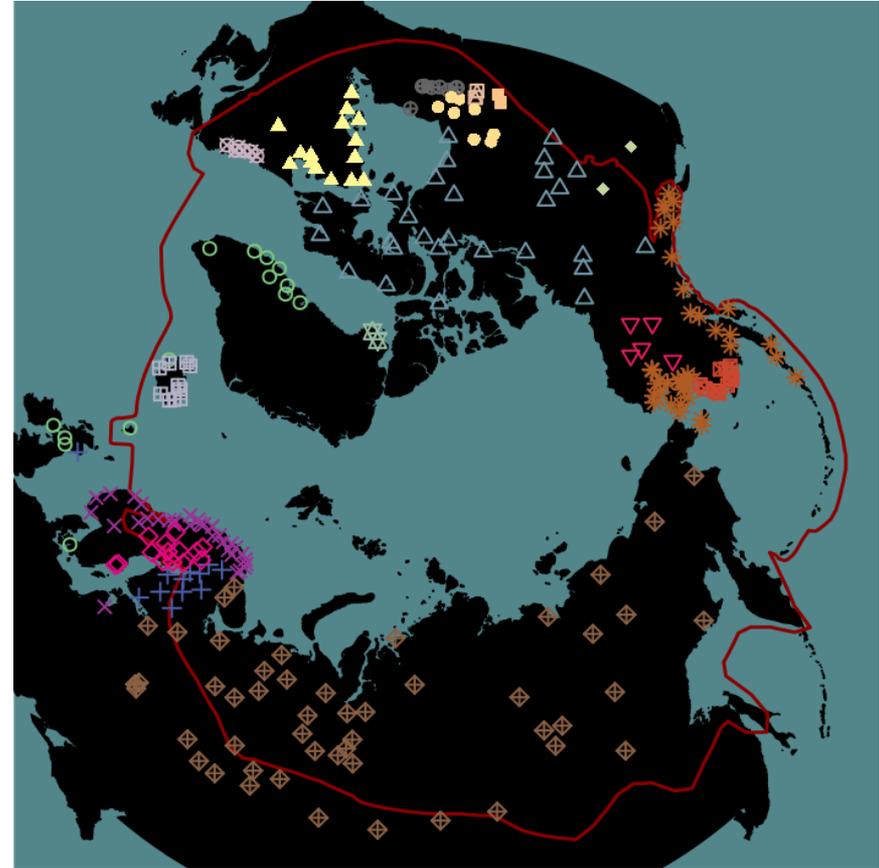
Free as in Freedom

<https://github.com/IDEES-Rouen/Flight-Scraping>

FlightRadar



Vol vers la zone arctique



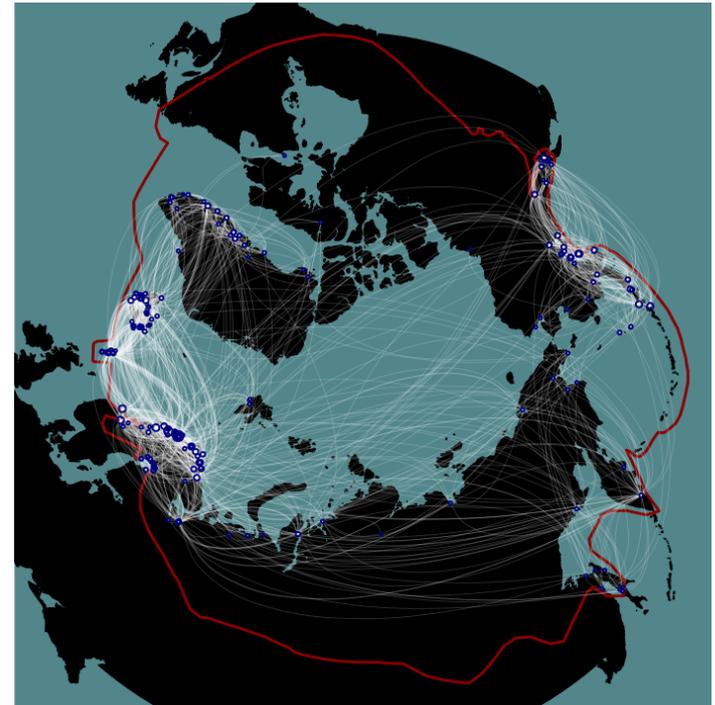
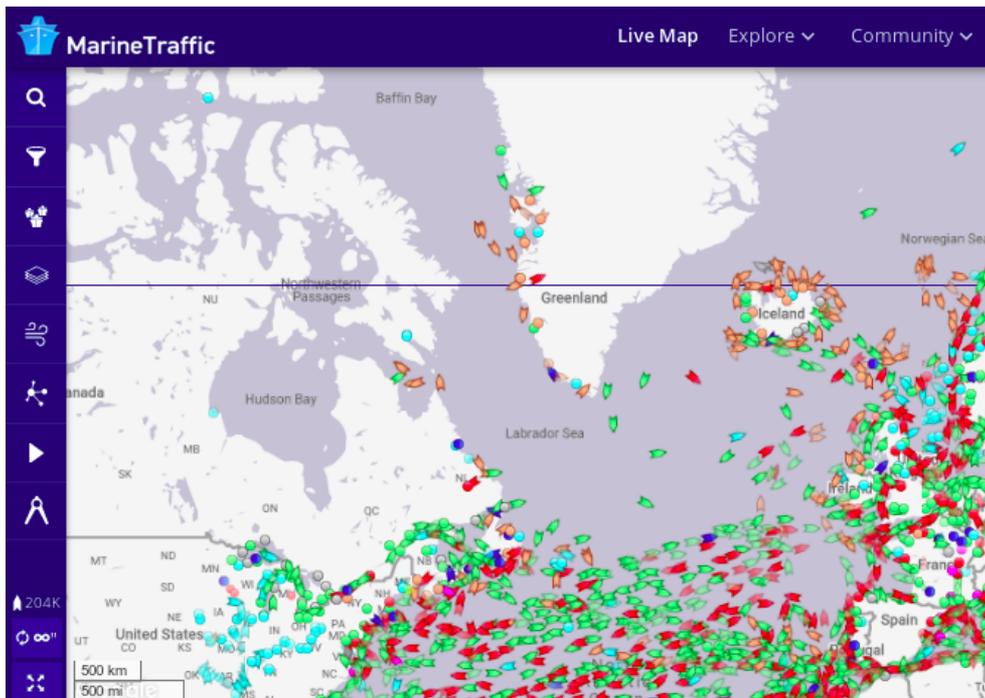
Détection de communauté

Marine Traffic

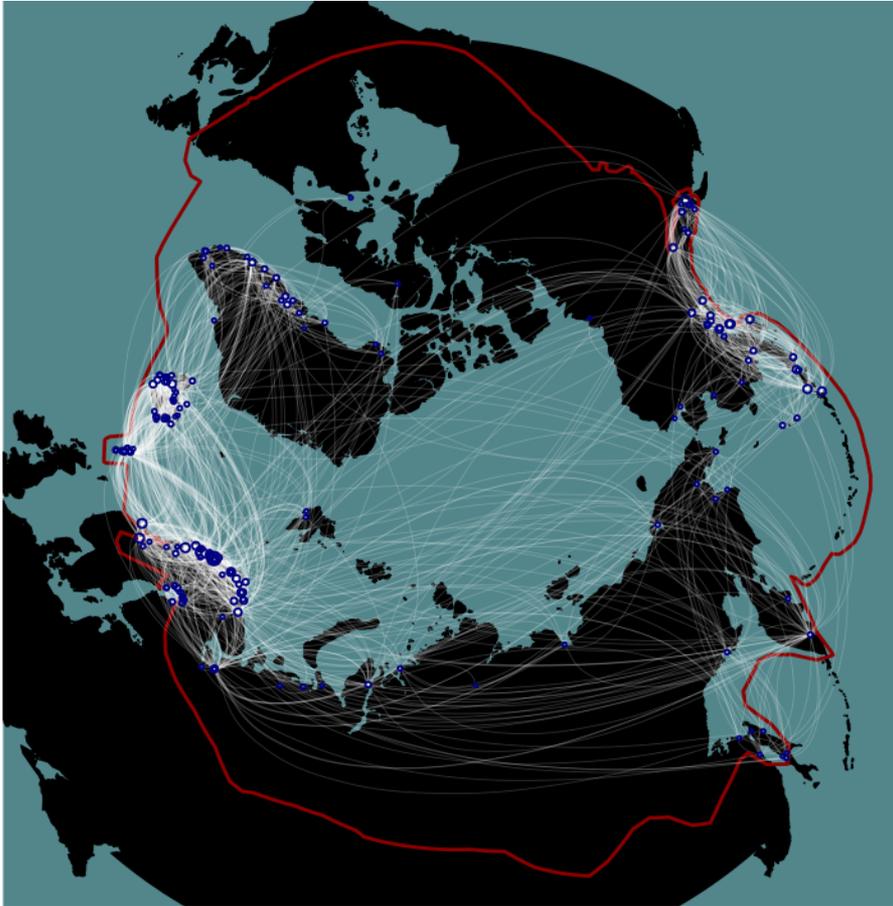


ANR PUR (2015 - 2019) / Yvette Vaguet

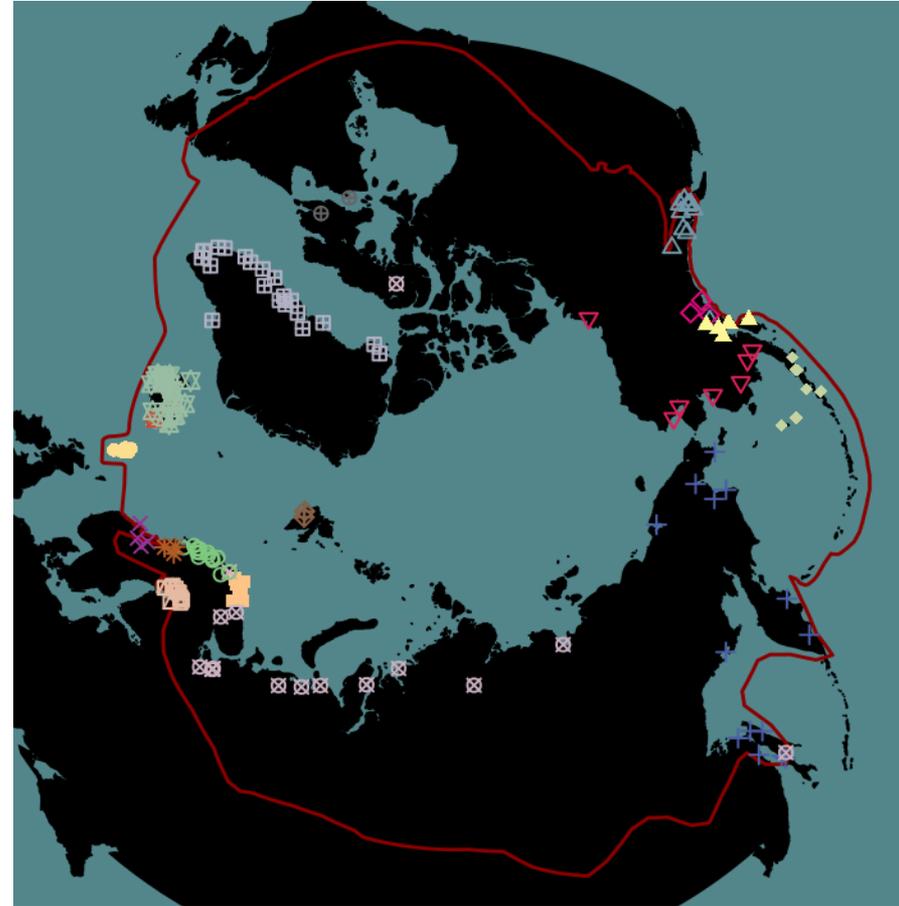
Récupérer les **départs et arrivées de bateaux** partant où allant vers la zone arctique



Marine Traffic



Flux maritime vers la zone arctique

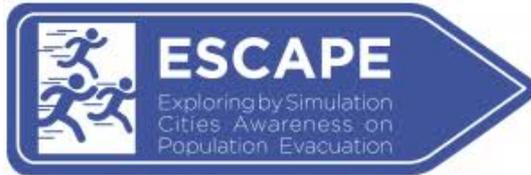


Détection de communauté

Mobilités urbaines

FP7 DENFREE & ANR MO3

Mobilité quotidiennes & propagation des épidémies de dengue



ANR Escape (2016 - 2024)

Modéliser les mobilités lors d'évacuation en cas d'urgence

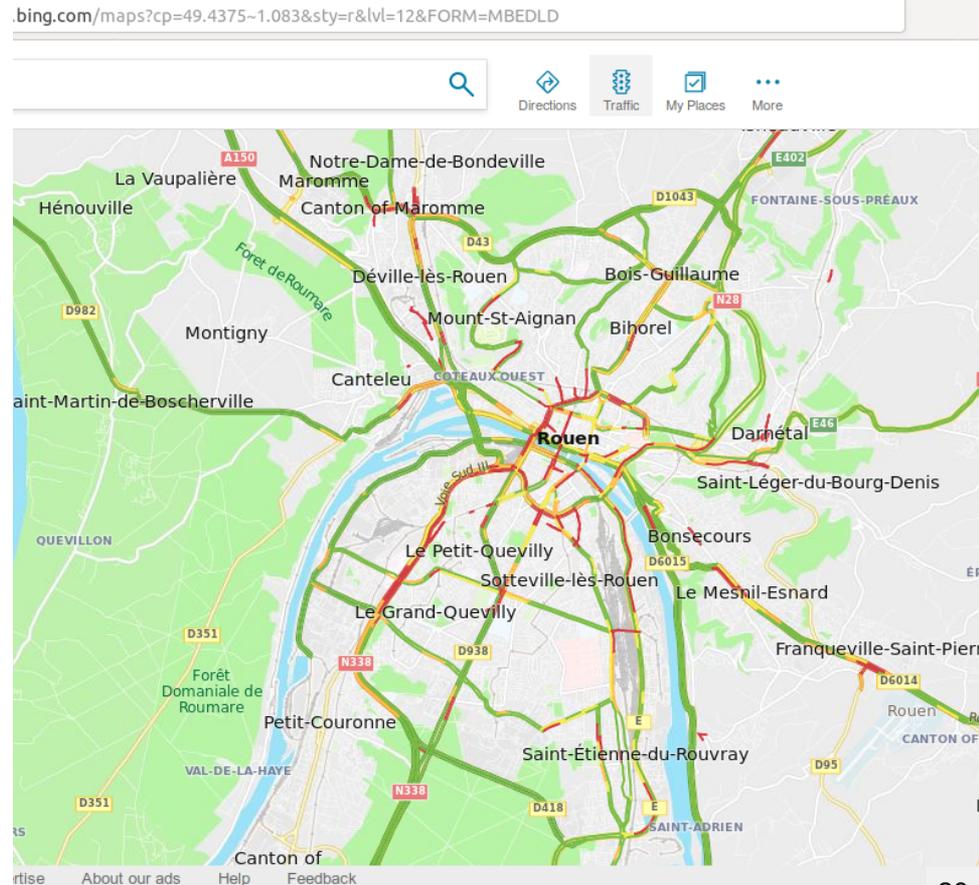
Besoin de données sur les mobilités comme les conditions de trafic pour nourrir et/ou valider les modèles

Condition de trafic



Bing Map (Microsoft)

Collecte en temps réel des conditions de circulation



Condition de trafic



Bing Map (Microsoft)

Collecte en temps réel des conditions de circulation

Site web : /www.bing.com/maps

Langage : R

Type de récolte : Campagne

Format : Raster (image)

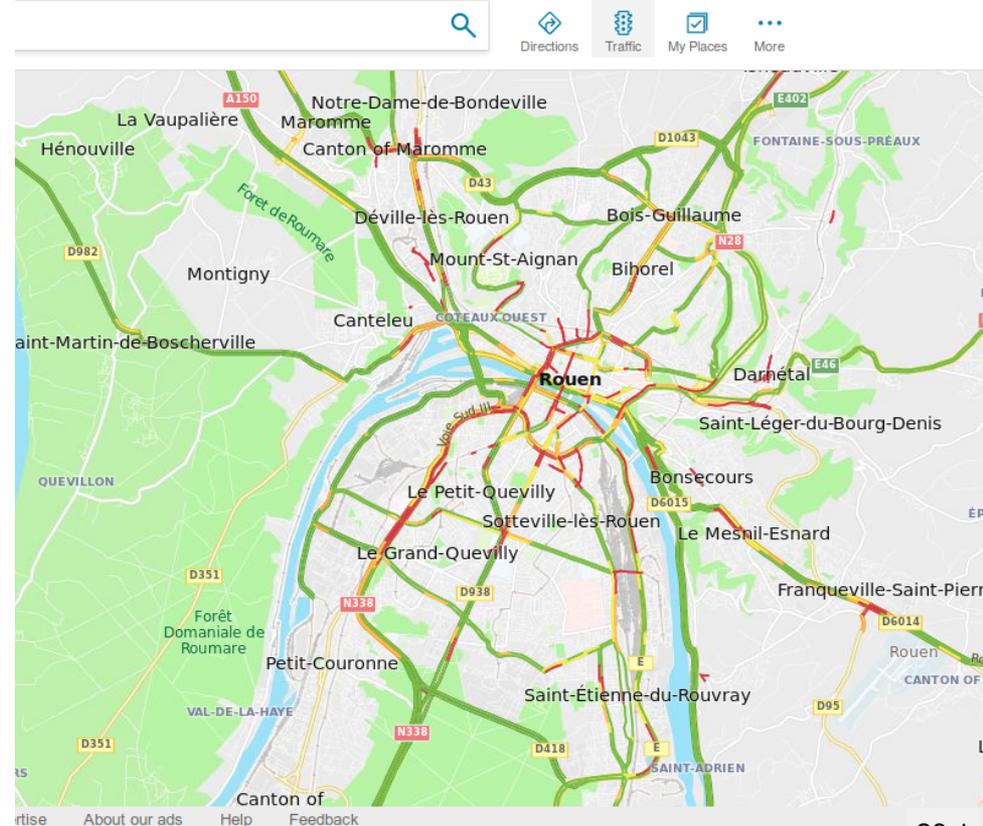
Avantages & limites :

- Trafic en temps réel
- Disponible dans 55 pays
- Méthodologie de Bing inconnue

Principe :

- Collecter et reprojeter des images
- Extraire les conditions de circulation

bing.com/maps?cp=49.4375~-1.083&sty=r&lvl=12&FORM=MBEDLD



Condition de trafic



Bing Map (Microsoft)

Collecte en temps réel des conditions de circulation

[.bing.com/maps?cp=49.4375~-1.083&sty=r&lvl=12&FORM=MBEDLD](https://www.bing.com/maps?cp=49.4375~-1.083&sty=r&lvl=12&FORM=MBEDLD)

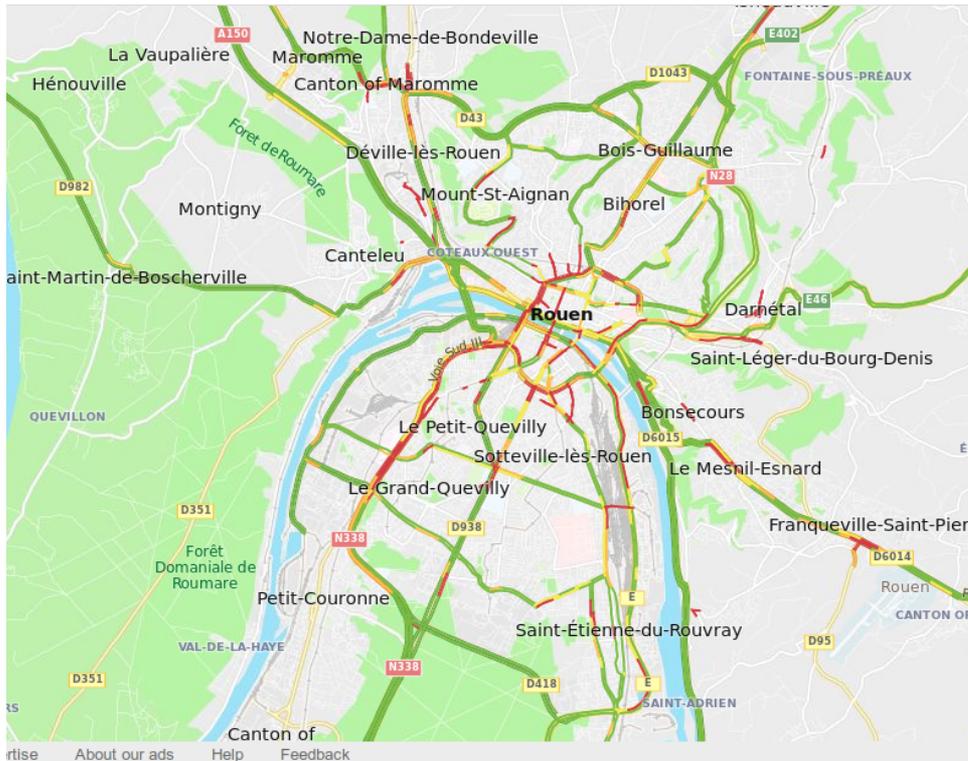


Directions

Traffic

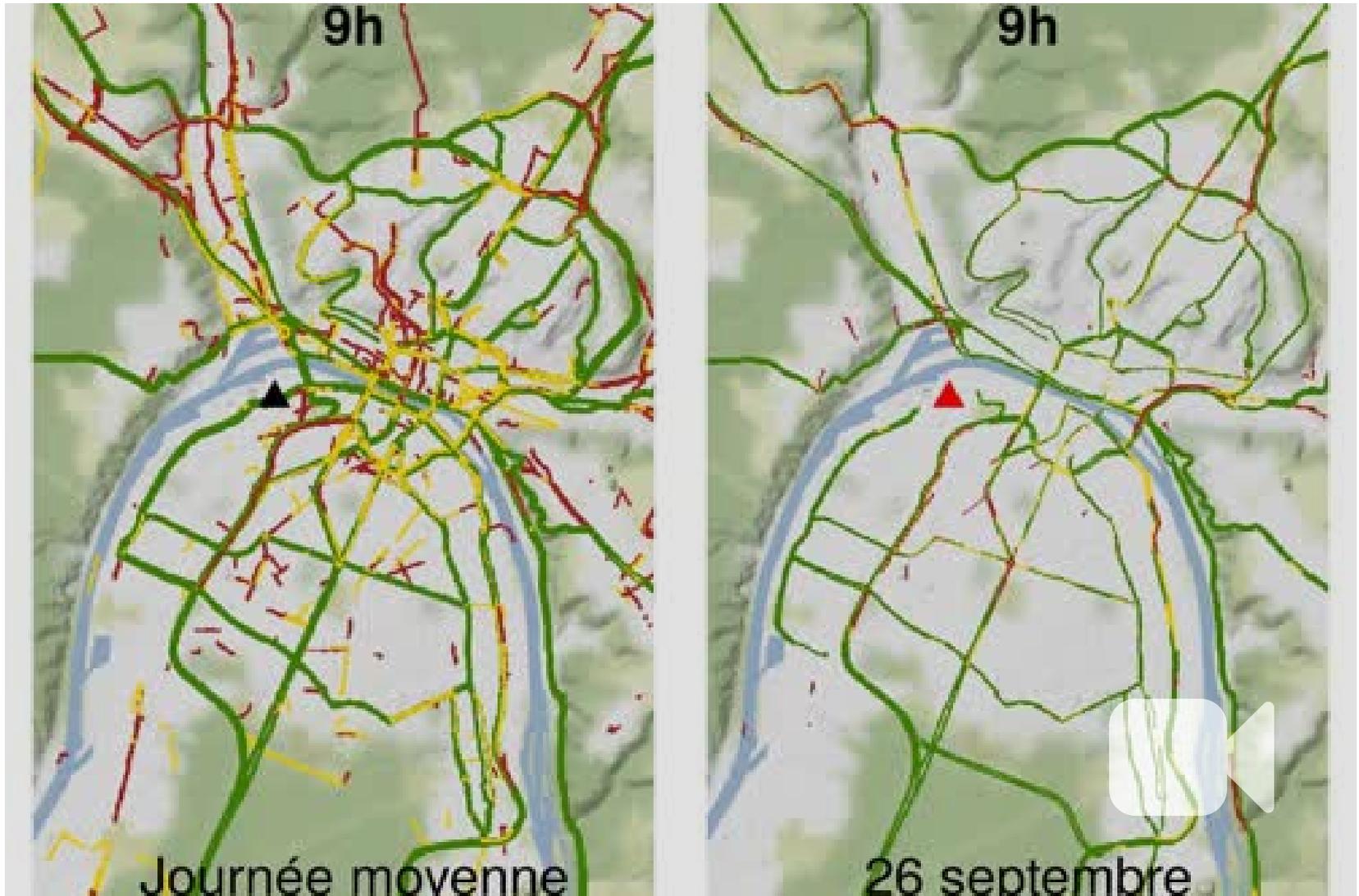
My Places

More



rtise About our ads Help Feedback

Condition de trafic



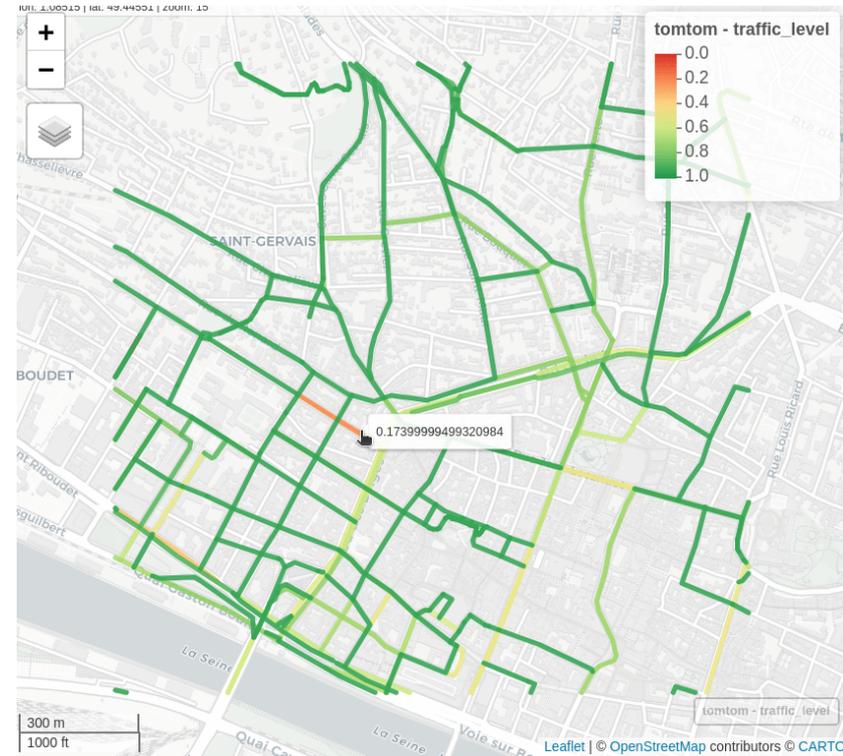
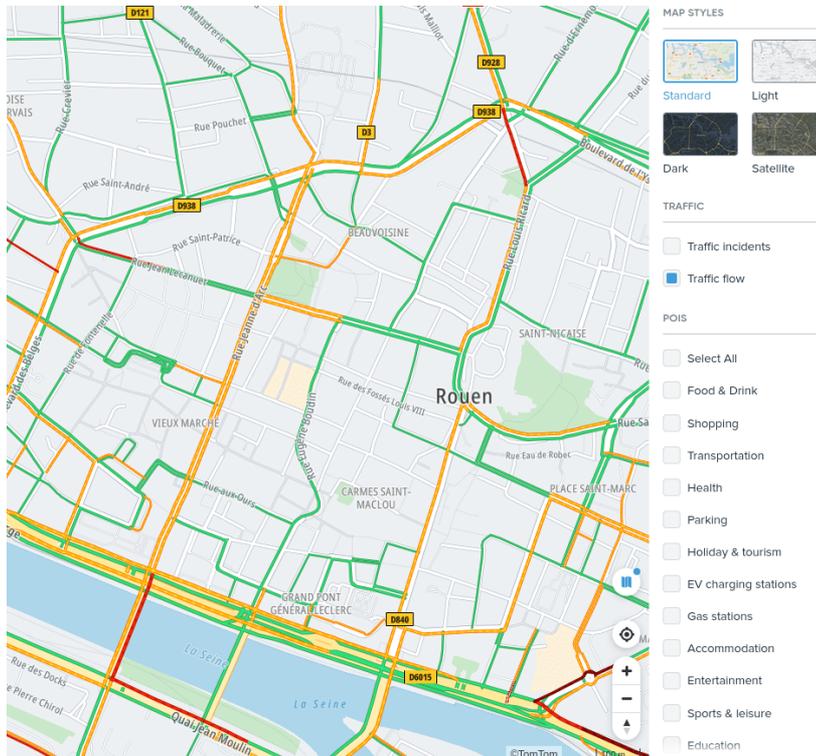
Trafic après l'explosion de Lubrizol

Condition de trafic



TomTom

Collecte en temps réel des conditions de circulation



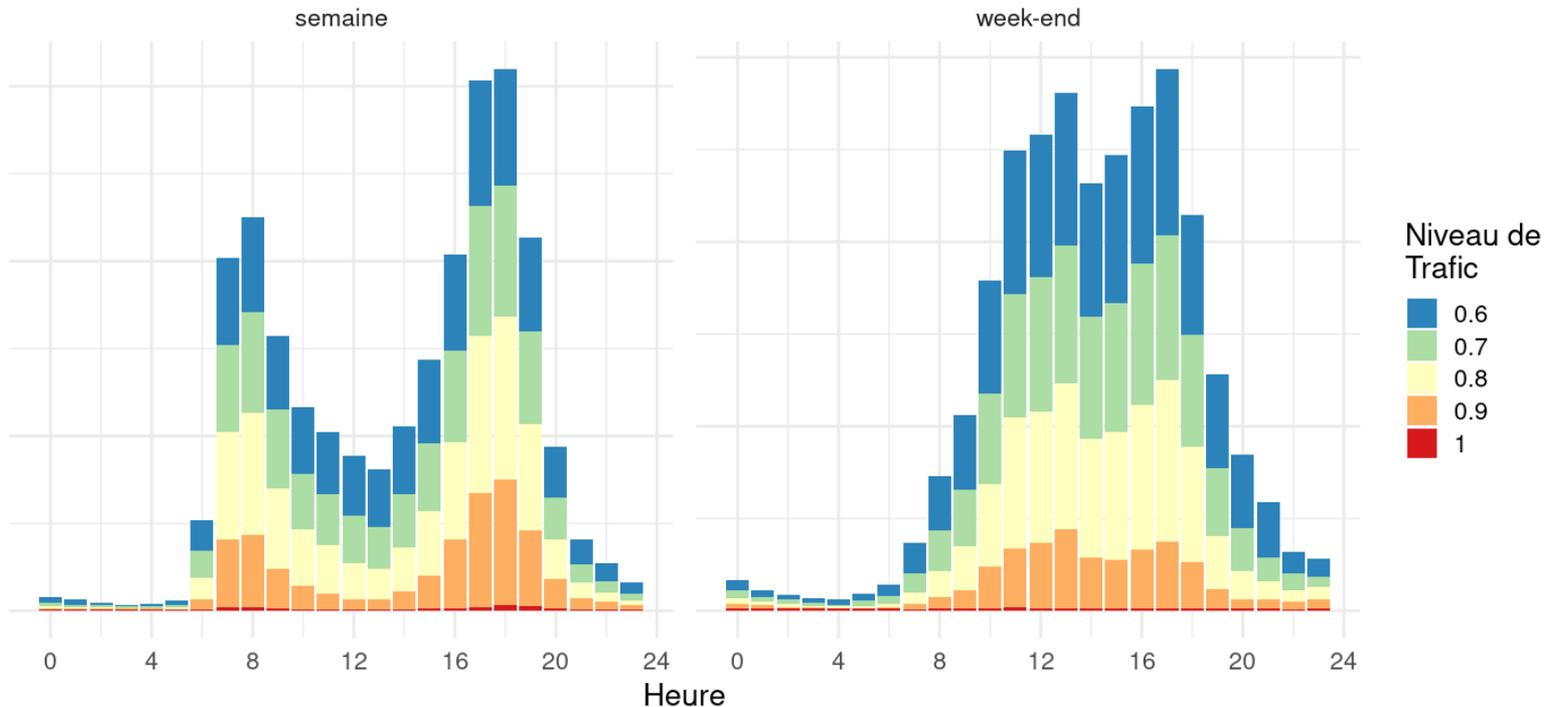
Données vectorielles scrappées plus complètes (niveau de trafic)

Condition de trafic



TomTom

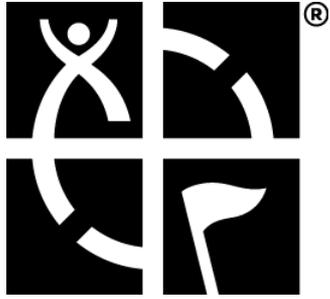
Collecte en temps réel des conditions de circulation



Permet la validation de modèles de mobilités

Écart entre mobilité simulées et trafic observé

GeoCaching



GRR Tenum (2015 - 2019) / Philippe Vidal

Récupérer **toutes les informations** sur les ~ **13000 géocaches** (1/page) de Normandie à un instant **t**

Site web : geocaching.com

Langage : Python

Framework : Scrapy

Type de récolte : Campagne

Format : json

Difficultés :

- Pas d'API & limitation à 1000 résultats par recherche
- Mon premier scraping
- Login/Pwd pour accès à la donnée complète
- ASPX avec gestion des états précédents



DOI [10.5281/zenodo.3542261](https://doi.org/10.5281/zenodo.3542261)

Téléchargement automatique

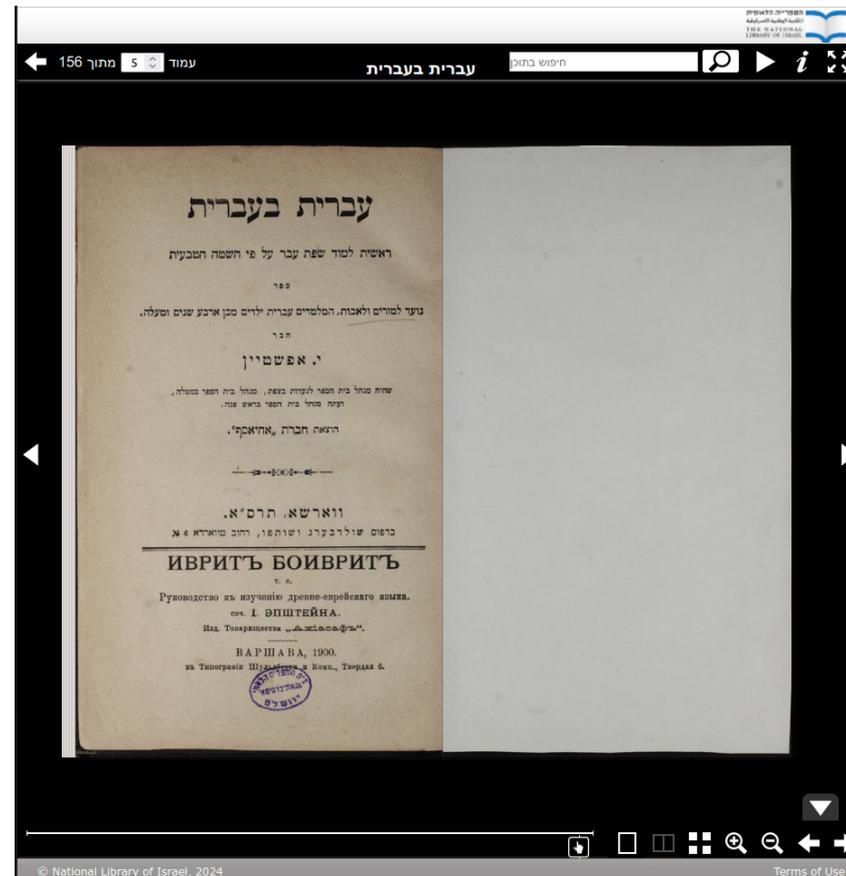


Projet RachelNet / Denis Eckert
Récupération d'ouvrages numérisés en ligne

**Téléchargement semi-automatisé
des pages d'ouvrages numérisés
accessibles uniquement en
navigation.**

- Permettre une consultation et un
archivage hors-ligne et plus rapide des
documents

```
1 library(tidyverse)
2 library(curl)
3 library(glue)
4 library(magick)
5
6 ## Livre Bibliothèque Nationale d'Israel
7 prenter <- "https://rosetta.nli.org.il/delivery/DeliveryManagerServlet?dps_func=stream&dps_pid=FL25784474"
8 dernier <- "https://rosetta.nli.org.il/delivery/DeliveryManagerServlet?dps_func=stream&dps_pid=FL25785206"
9
10 tous <- 25784474:25785206
11 errors <- numeric()
12 for (i in tous) {
13   thisAdress <- glue("https://rosetta.nli.org.il/delivery/DeliveryManagerServlet?dps_func=stream&dps_pid=FL{i}")
14   tryCatch({
15     curl_download(url = thisAdress, destfile = glue("scraping_denis/bouquin_bibli_israel/{i}.jpg"))
16   },
17   error=function(e){cat("ERROR :",conditionMessage(e), "\n")}
18 }
19
20 system('convert "*.jpg" -quality 100 ouvrage.pdf')
```



Marché immobilier mexicain



PrimoMex / A. Ribardière, J.-F. Valette & L. Salinas
Caractériser le marché immobilier et son évolution dans
l'agglomération de Mexico (ZMVM)

Un des objectifs :

Étudier l'évolution des prix à
différentes échelles depuis 2015

Moyen : Base 2015 exhaustive,
scraping de l'ensemble des
annonces en 06/2024

Un des enjeux :

Détecter les annonces
incomplètes, mal renseignées, ou
trompeuses.

The screenshot displays the Mercado Libre website interface. At the top, there is a yellow navigation bar with the Mercado Libre logo, a search bar, and various utility links. Below the navigation bar, the main content area is titled 'Departamentos' and shows a list of real estate listings on the left and a map of Mexico City on the right. The listings include details such as price, area, and number of rooms. The map shows various neighborhoods in Mexico City, with blue dots indicating the locations of the listed properties. The bottom of the page features a footer with links for 'Trabaja con nosotros', 'Términos y condiciones', 'Promociones', 'Cómo cuidamos tu privacidad', 'Accesibilidad', 'Ayuda', and 'Hot Sale'.

Mobilités urbaines



FP7 DENFREE & ANR MO3 / Eric Daudé
Modéliser et simuler les épidémies de dengue

Un des objectifs :

Comprendre le rôle des mobilités quotidiennes dans la propagation des épidémies à Bangkok & Delhi

Moyen :

Créer un modèle de mobilité, individu centré, à base d'agents (couplé à d'autres modèles)

Concept d'espace d'activité

Un des enjeux :

Collecte de données de mobilité

Mobilités urbaines



FP7 DENFREE & ANR MO3 / Eric Daudé
Modéliser et simuler les épidémies de dengue

Un des objectifs :

Comprendre le rôle des mobilités quotidiennes dans la propagation des épidémies à Bangkok & Delhi

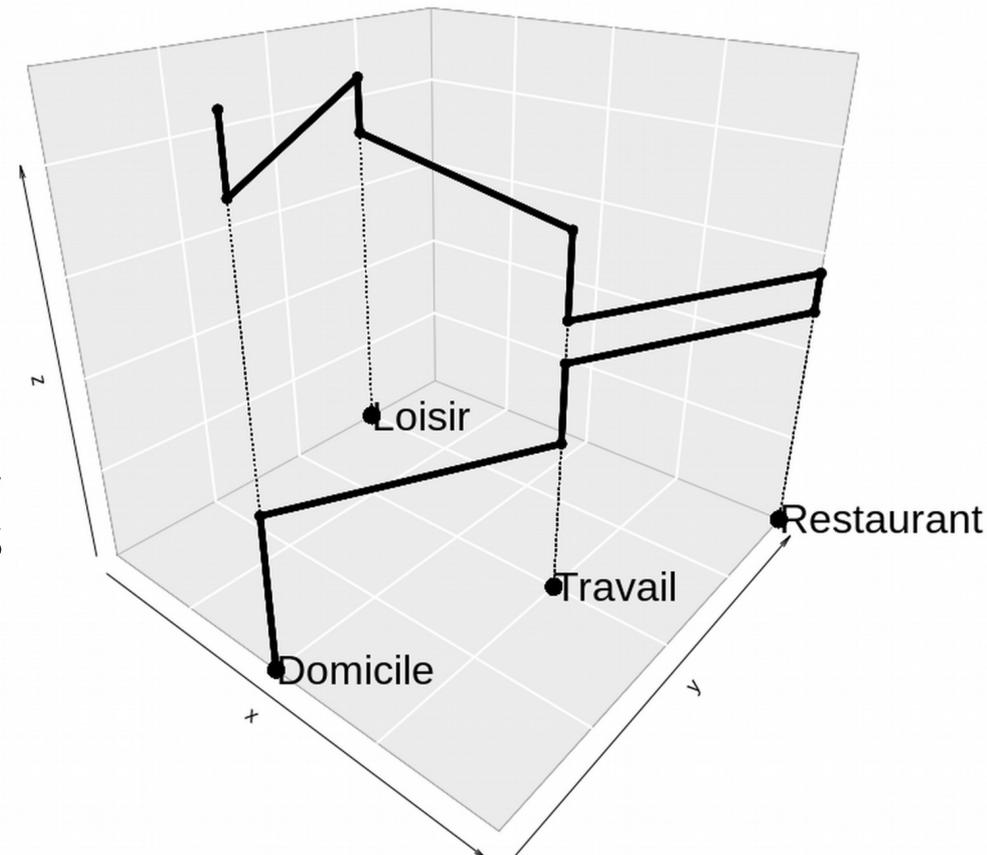
Moyen :

Créer un modèle de mobilité, individu centré, à base d'agents (couplé à d'autres modèles)

Concept d'espace d'activité

Un des enjeux :

Collecte de données de mobilité



Mobilités urbaines



FP7 DENFREE & ANR MO3 / Eric Daudé
Modéliser et simuler les épidémies de dengue

Un des objectifs :

Comprendre le rôle des mobilités quotidiennes dans la propagation des épidémies à Bangkok & Delhi

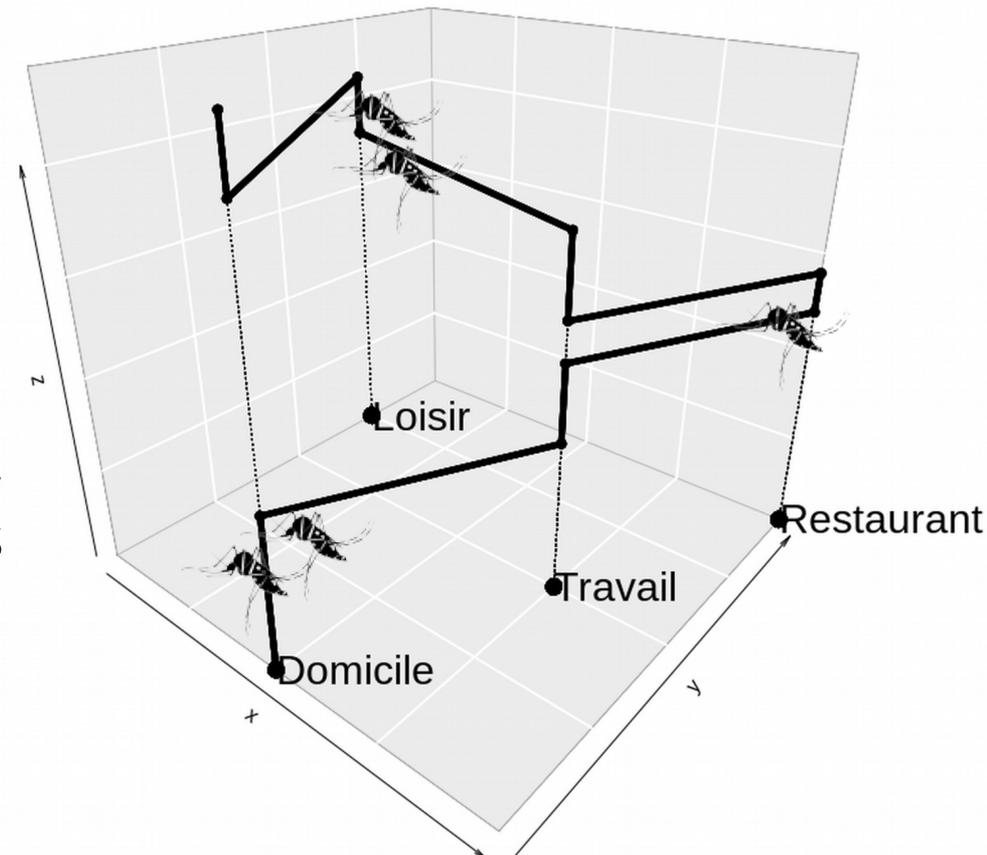
Moyen :

Créer un modèle de mobilité, individu centré, à base d'agents (couplé à d'autres modèles)

Concept d'espace d'activité

Un des enjeux :

Collecte de données de mobilité



Mobilités urbaines



FP7 DENFREE & ANR MO3 / Eric Daudé
Modéliser et simuler les épidémies de dengue

Un des objectifs :

Comprendre le rôle des mobilités quotidiennes dans la propagation des épidémies à Bangkok & Delhi

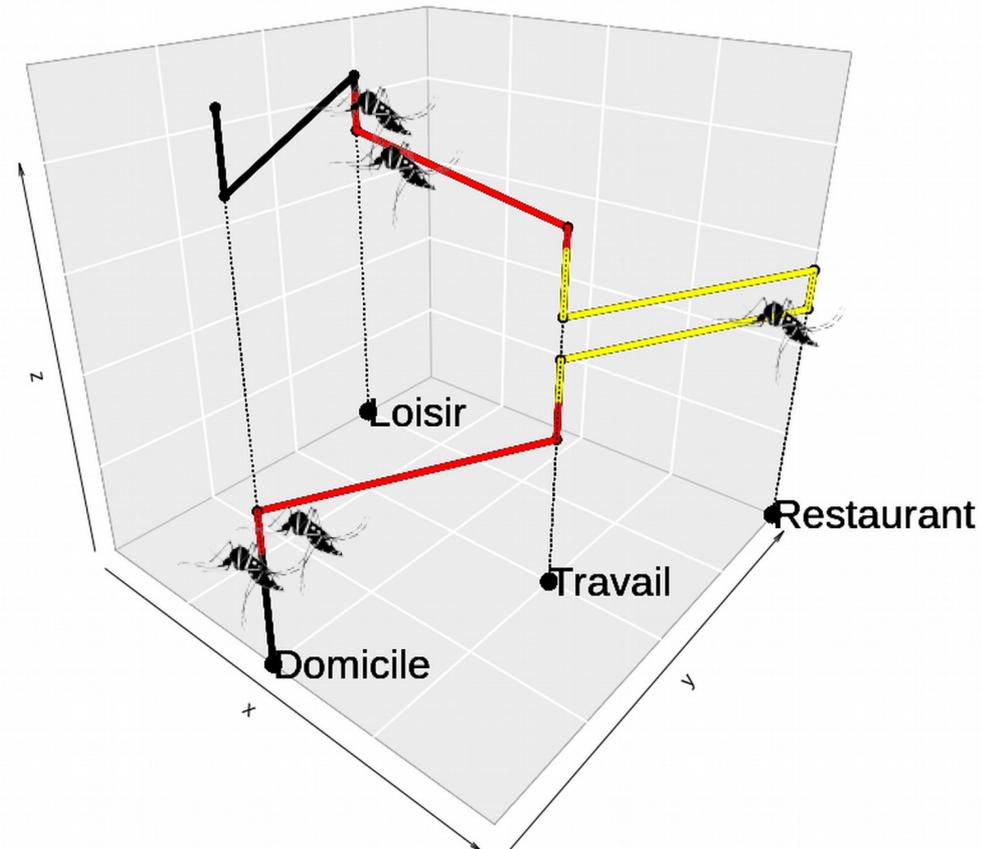
Moyen :

Créer un modèle de mobilité, individu centré, à base d'agents (couplé à d'autres modèles)

Concept d'espace d'activité

Un des enjeux :

Collecte de données de mobilité



Évacuations massives



ANR Escape (2016 - 2024) / Eric Daudé

"Contribuer à la conception de systèmes d'aide à la décision dans le cas d'évacuations massives"

Un des objectifs :

Simuler les mobilités en cas de catastrophes (industrielles ou naturelles). Multisites

Comment ?

Simulation à base d'agents

Préalable :

Connaître les comportements de mobilités en condition "normale"

==> Intérêt des conditions de trafic
(en plus des EMD)

Évacuations massives



ANR Escape (2016 - 2024) / Eric Daudé

"Contribuer à la conception de systèmes d'aide à la décision dans le cas d'évacuations massives"

Un des objectifs :

Simuler les mobilités en cas de catastrophes (industrielles ou naturelles). Multisites

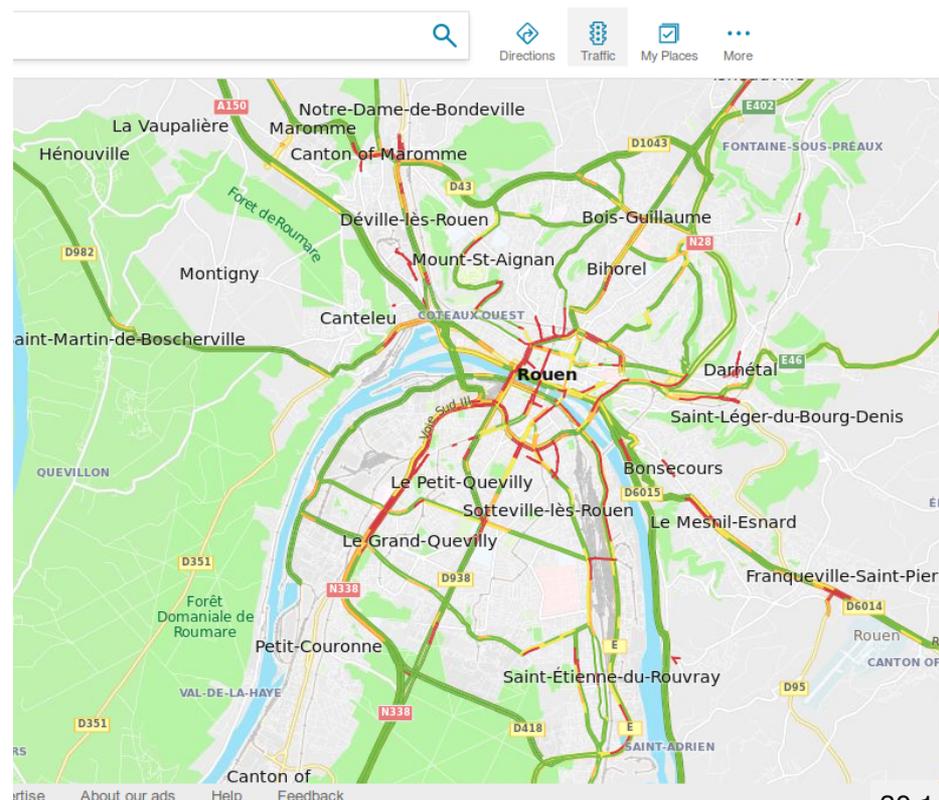
Comment ?

Simulation à base d'agents

Préalable :

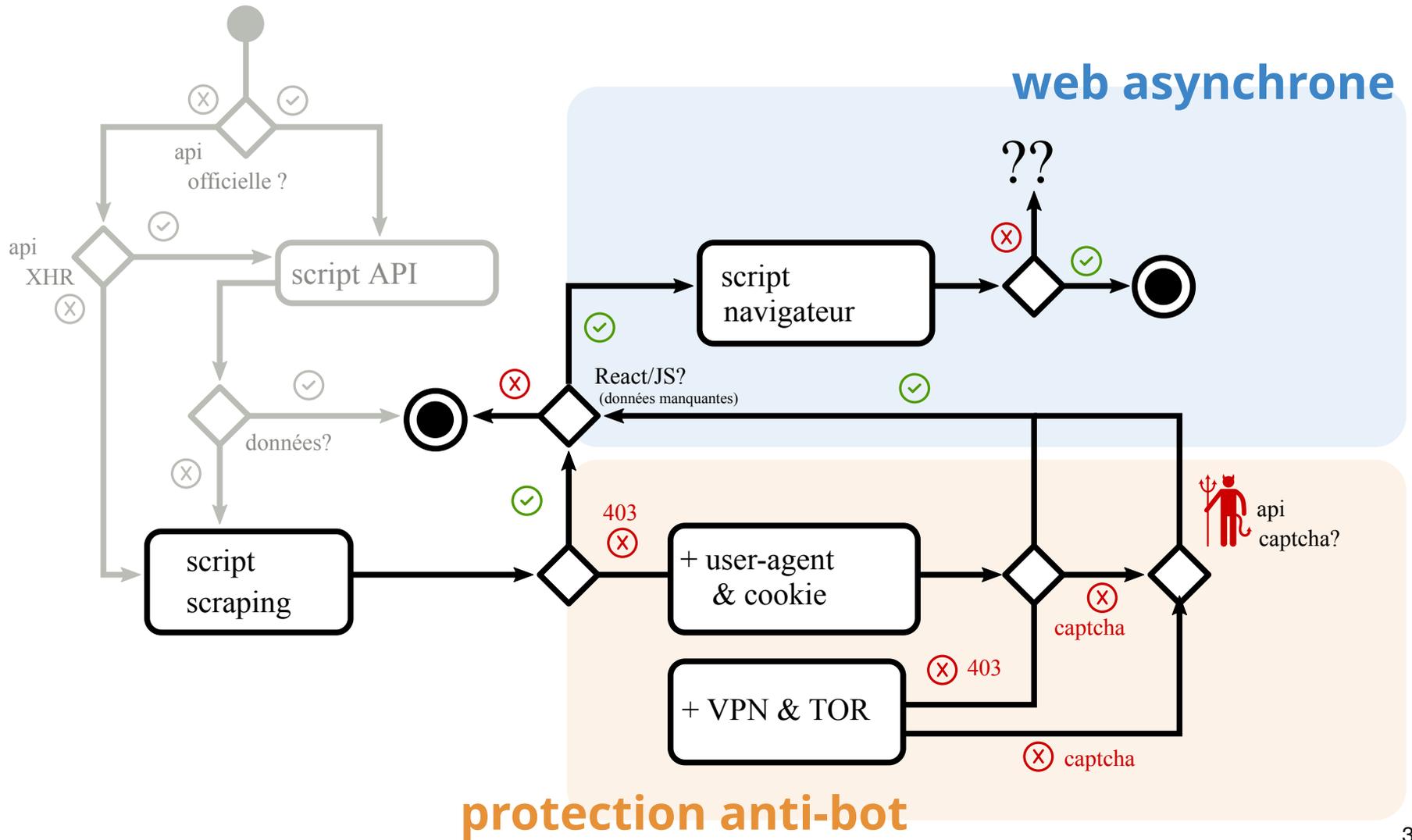
Connaître les comportements de mobilités en condition "normale"
==> Intérêt des conditions de trafic
(en plus des EMD)

bing.com/maps?cp=49.4375-1.083&sty=r&lvl=12&FORM=MBEDLD



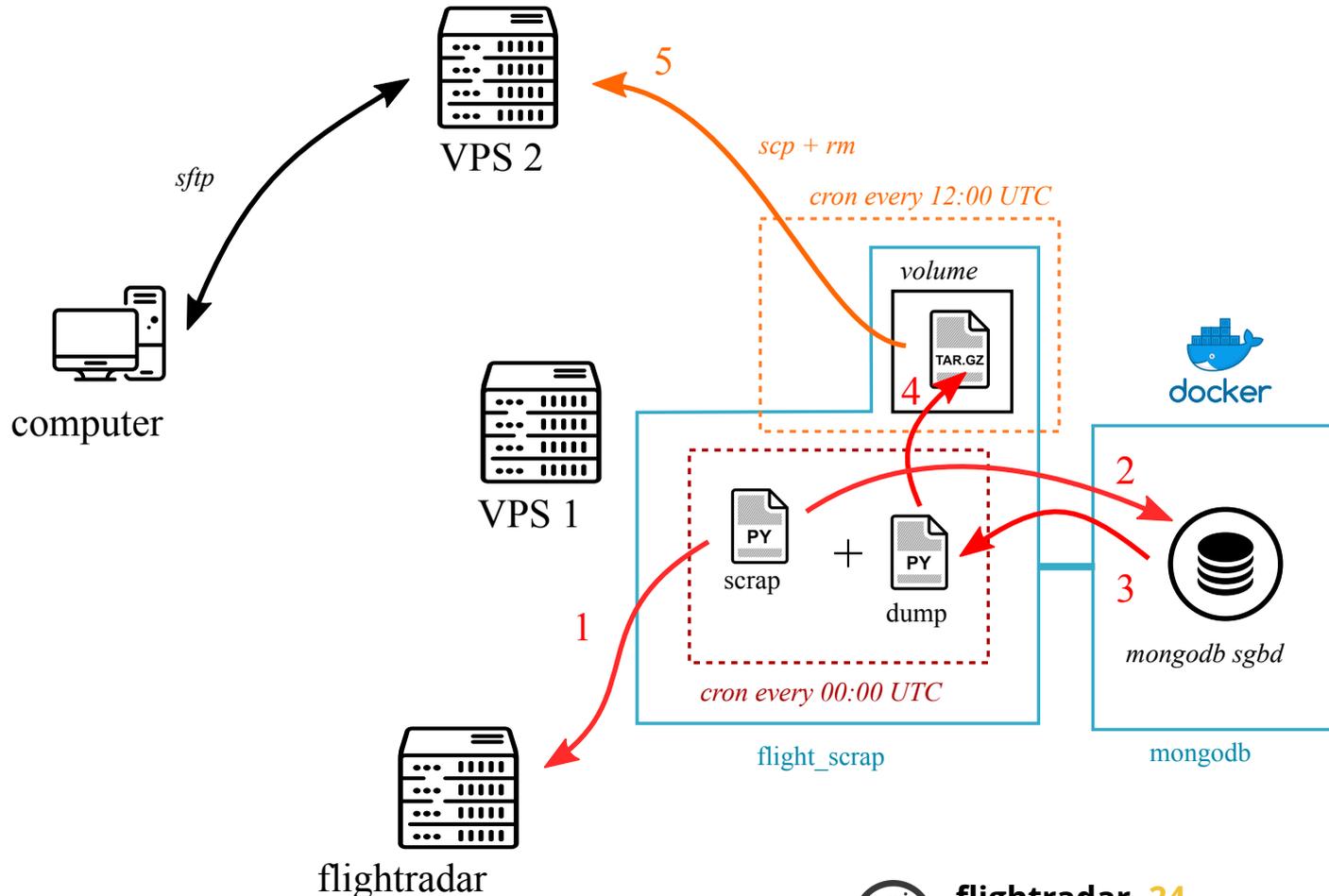
(III) Verrous techniques

(a) définir un logigramme pour mieux s'y retrouver



(III) Verrous techniques

(b) mettre une campagne de collecte h24/7J en place, une problématique complexe rapidement ...



flightradar 24

<https://www.flightradar24.com>

<https://github.com/IDEES-Rouen/Flight-Scraping>