

Extraction d'information

Philippe Langlais

`felipe@iro.umontreal.ca`

27 novembre 2018

École **IVADO/RALI**

`https://ivado.ca/formations/ecoles/ecole-analyse-du-langage-naturel-2018/`

Extraction d'information

Extraction d'information ouverte

Intermède

Structuration

Applications

Conclusion

Extraction d'information

What ?

https://en.wikipedia.org/wiki/Information_extraction

*Information extraction. Information extraction (IE) is the task of automatically **extracting structured information from unstructured** and/or semi-structured machine-readable **documents**. In most of the cases this activity concerns processing human language texts by means of natural language processing (NLP).*

Extraire quoi ?

https://fr.wikipedia.org/wiki/Tigran_Hamasyan

[En 2009]_{TEMP}, il enregistre [Red Hail]_{ALBUM}, un album au carrefour du jazz, du metal et du folklore arménien, avec son nouveau quintet de jeunes musiciens [Aratta Rebirth]_{ORG} : [Areni Agbabian]_{PER} (voc), [Ben Wendel]_{PER} (ts), [Charles Altura]_{PER} (g), [Sam Minaie]_{PER} (b) et [Nate Wood]_{PER} (d). Ils se produisent dans plusieurs grands festivals internationaux, de [Montréal]_{LOC} à [Nice]_{LOC} en passant par [Vienne]_{LOC} ou [Rotterdam]_{LOC} ([North Sea Jazz Festival]_{EVENT}).

- entités nommées, expressions temporelles, valeurs numériques
- relations entre entités

ex : ⟨Aratta Rebirth, se produit, dans plusieurs grands festivals internationaux⟩

Quelles relations ?

https://fr.wikipedia.org/wiki/Tigran_Hamasyan

[En 2009]_{TEMP}, il enregistre [Red Hail]_{ALBUM}, un album au carrefour du jazz, du metal et du folklore arménien, avec son nouveau quintet de jeunes musiciens [Aratta Rebirth]_{ORG} : [Areni Agbabian]_{PER} (voc), [Ben Wendel]_{PER} (ts), [Charles Altura]_{PER} (g), [Sam Minaie]_{PER} (b) et [Nate Wood]_{PER} (d). Ils se produisent dans plusieurs grands festivals internationaux, de [Montréal]_{LOC} à [Nice]_{LOC} en passant par [Vienne]_{LOC} ou [Rotterdam]_{LOC} ([North Sea Jazz Festival]_{EVENT}).

- relations particulières : $\langle \textit{artiste}, a_produit, \langle \textit{album} \rangle$
(Tigran_Hamasyan, Red_Hail)
- ou toutes les relations exprimées ... ou le plus possible !

Dans quel but ?

Peupler une base de connaissances pour :

- une meilleure recherche d'information

Pianistes comme Tigran Hamasyan

Google 🔍

Tous Vidéos Images Actualités Shopping Plus Paramètres Outils




Environ 19 200 résultats (0,28 secondes)

Tigran Hamasyan — Wikipédia
https://fr.wikipedia.org/wiki/Tigran_Hamasyan ▼
Tigran Hamasyan est un pianiste de jazz arménien né le 17 juillet 1987 à Gyumri (Arménie). ... Le jeune musicien fait alors la connaissance de légendes comme Wayne Shorter, Herbie Hancock, John McLaughlin ou Joe Zawinul et de ...
Biographie · Récompenses · Discographie · En tant que leader

INTERVIEW. Tigran Hamasyan, pianiste contemporain - Culturebox
<https://culturebox.francetvinfo.fr> · Musique · Jazz / Blues ▼
24 mars 2017 - Le 31 mars, le pianiste arménien Tigran Hamasyan a sorti un nouvel album ... Comme il aime à le faire dans ses disques, Tigran Hamasyan a ...

Tigran Hamasyan, pianiste de jazz virtuose - Culturebox
<https://culturebox.francetvinfo.fr> · Musique · Jazz / Blues ▼
29 mars 2010 - À 22 ans, le pianiste de jazz Tigran Hamasyan est un prodige qui a déjà ... récompensé dans les plus grands festivals de jazz, comme celui de ...

Vidéos

		
Le surprenant clip du pianiste Tigran L'Express - 12 juin 2013	Tigran Hamasyan "The Spinners" @Jazz_in_Marcillac 8 Août 2011 Jazz In Marcillac YouTube - 11 août 2011	Tigran Hamasyan piano solo 2011 Piano Web YouTube - 30 nov. 2016

Tigran Hamasyan



Tous Vidéos Images Actualités Maps Plus Paramètres Outils

Environ 591 000 résultats (0,46 secondes)

Résultats pour **tigran hamasyan**
Rechercher plutôt [tigran hamasyan](#)

Tigran Hamasyan – Official website

www.tigranhamasyan.com/ Traduire cette page

Tigran Hamasyan, new solo album 'An Ancient Observer' coming out march 2017 .

Tigran Hamasyan — Wikipédia

https://fr.wikipedia.org/wiki/Tigran_Hamasyan

Tigran Hamasyan est un pianiste de jazz arménien né le 17 juillet 1987 à Gyumri (Arménie).

Sommaire · 1 Biographie · 2 Récompenses · 3 Discographie · 3.1 En ...

[Biographie](#) · [Récompenses](#) · [Discographie](#) · [En tant que leader](#)

Vidéos

<p>Tigran Hamasyan - Rays of Light (Official Video)</p> <p>TigranHamasyan YouTube - 10 janv. 2018</p>	<p>Tigran Hamasyan - Live @ Montreal Jazz Festival</p> <p>Culturebox YouTube - 15 déc. 2017</p>	<p>Tigran Hamasyan - Leninagone (Official Video)</p> <p>TigranHamasyan YouTube - 7 déc. 2017</p>



Tigran Hamasyan

Pianiste

Proposé sur

YouTube

Spotify

Deezer

Tigran Hamasyan est un pianiste de jazz arménien né le 17 juillet 1987 à Gyumri. [Wikipédia](#)

Date et lieu de naissance : 17 juillet 1987 (Âge: 31 ans), Gyumri, Arménie

Films : Bravo virtuose

Enseignement : Université de Californie du Sud, Tchaikovsky Music School, Yerevan

Nominations : Edison Jazz/World - Jazz International

Titres

Fides Tua

An Ancient Observer - 2017





Plus d'images

Tigran Hamasyan

Pianiste



Proposé sur

YouTube

Spotify

Deezer

Tigran Hamasyan est un pianiste de jazz arménien né le 17 juillet 1987 à Gyumri. [Wikipédia](#)

Date et lieu de naissance : 17 juillet 1987 (Âge: 31 ans), Gyumri, Arménie

Films : [Bravo virtuose](#)

Enseignement : Université de Californie du Sud, Tchaikovsky Music School, Yerevan

Nominations : Edison Jazz/World - Jazz Internationaal

Titres

Fides Tua

An Ancient Observer - 2017

The Cave of Rebirth

An Ancient Observer - 2017

Ancient Observer

An Ancient Observer - 2017

Afficher 25 autres éléments

Albums

Afficher 3 autres éléments



Mockroot

2015



An Ancient Observer

2017



Shadow Theater

2013



Luys I Luso

2015



For Gyumri

2018

Recherches associées

Afficher 10 autres éléments



Lars Danielsson



Eivind Aarset



Arve Henriksen



Jan Bang



Avishai Cohen

Dans quel but ?

Peupler une base de connaissances pour :

- une meilleure recherche d'information
- outils de fouille



Open Information Extraction

Argument 1:

Relation:

kills

Argument 2:

bacteria

All

Search

198 answers from 1523 sentences (cached)

all

drug ingredient (30)

drug (24)

medical treatment (19)

ingredient (15)

chemical compound (13)

misc.

more types

Antibiotic (165)

Chlorine (76)

Water (59)

Benzoyl peroxide (45)

Heat (40)

Antiseptic (38)

Pasteurization (35)

Cooking (34)

Vinegar (34)

Honey (28)

Tea tree oil (24)

Ultraviolet (24)

Antibiotic

URI:

<http://www.freebase.com/view/m/0tbr>

Types:

[medical procedure](#) (Nell)
[/medicine/medical_treatment](#) (FreeBase)
[/medicine/risk_factor](#) (FreeBase)
[/medicine/drug_class](#) (FreeBase)

Extracted Synonyms:

Antibiotics
the antibiotics
these antibiotics
the antibiotic
an antibiotic
most antibiotics
Many antibiotics

Extracted from these sentences:

kill **Antibiotics** kill your friendly **bacteria** , and when these have been killed the Candida yeast can mutate in
(via ClueWeb12) felipe@iro.umontreal.ca 8
Antibiotics kill **most bacteria** , but can leave resistant individual strains free to reproduce without compete

Meilleure compréhension

Un sous-ensemble de *benchmarks* de question-réponse :

- SQUAD
- MS Marco
- NewsQA
- MC Test
- Cloze-style
- **RACE**
- CNN/Daily
- WDW (What did what)
- **AI2 Elementary school Science questions**
- triviaQA
- WikiHop

- 7787 questions à choix multiples en sciences naturelles :
 - Easy set** : 5197 questions
 - Challenge set** : 2590 questions qui résistent à des baselines basés sur la RI ou la co-occurrence
- questions pour étudiants de 8 à 13 ans

Which property of a mineral can be determined just by looking at it?

(A) luster (B) mass (C) weight (D) hardness

- 7787 questions à choix multiples en sciences naturelles :
 - **Easy set** : 5197 questions
 - **Challenge set** : 2590 questions qui résistent à des baselines basés sur la RI ou la co-occurrence
- questions pour étudiants de 8 à 13 ans

Which property of a mineral can be determined just by looking at it?

(A) luster (B) mass (C) weight (D) hardness

- 7787 questions à choix multiples en sciences naturelles :
 - Easy set** : 5197 questions
 - Challenge set** : 2590 questions qui résistent à des baselines basés sur la RI ou la co-occurrence
- questions pour étudiants de 8 à 13 ans

A student riding a bicycle observes that it moves faster on a smooth road than on a rough road. This happens because the smooth road has

(A) less gravity (B) more gravity (C) less friction (D) more friction

ARC (AI2 Reasoning Challenge) [Clark et al., 2018]

- 7787 questions à choix multiples en sciences naturelles :
 - Easy set** : 5197 questions
 - Challenge set** : 2590 questions qui résistent à des baselines basés sur la RI ou la co-occurrence
- questions pour étudiants de 8 à 13 ans

A student riding a bicycle observes that it moves faster on a smooth road than on a rough road. This happens because the smooth road has

(A) less gravity (B) more gravity (C) **less friction** (D) more friction

RACE [Lai et al., 2017]

- 28k passages, 100k questions posées par des enseignants de l'anglais langue seconde à des étudiants Chinois âgés de 12 à 18 ans
- deux niveaux : RACE-M (12-15 ans), RACE-H (15-18)
- vocabulaire restreint (apprenants de l'anglais)

- dataset : <http://www.cs.cmu.edu/~glai1/data/race/>
- code : https://github.com/qizhex/RACE_AR_baselines

Passage:

In a small village in England about 150 years ago, a mail coach was standing on the street. It didn't come to that village often. People had to pay a lot to get a letter. The person who sent the letter didn't have to pay the postage, while the receiver had to. "Here's a letter for Miss Alice Brown," said the mailman.

"I'm Alice Brown," a girl of about 18 said in a low voice.

Alice looked at the envelope for a minute, and then handed it back to the mailman.

"I'm sorry I can't take it, I don't have enough money to pay it", she said.

A gentleman standing around were very sorry for her. Then he came up and paid the postage for her.

When the gentleman gave the letter to her, she said with a smile, " Thank you very much, This letter is from Tom. I'm going to marry him. He went to London to look for work. I've waited a long time for this letter, but now I don't need it, there is nothing in it."

"Really? How do you know that?" the gentleman said in surprise.

"He told me that he would put some signs on the envelope. Look, sir, this cross in the corner means that he is well and this circle means he has found work. That's good news."

The gentleman was Sir Rowland Hill. He didn't forgot Alice and her letter.

"The postage to be paid by the receiver has to be changed," he said to himself and had a good plan.

"The postage has to be much lower, what about a penny? And the person who sends the letter pays the postage. He has to buy a stamp and put it on the envelope." he said . The government accepted his plan. Then the first stamp was put out in 1840. It was called the "Penny Black". It had a picture of the Queen on it.

Questions:

- | | |
|---|--|
| 1): The first postage stamp was made ...
A. in England B. in America C. by Alice D. in 1910 | 4): The idea of using stamps was thought of by ...
A. the government
B. Sir Rowland Hill
C. Alice Brown
D. Tom |
| 2): The girl handed the letter back to the mailman because ...
A. she didn't know whose letter it was
B. she had no money to pay the postage
C. she received the letter but she didn't want to open it
D. she had already known what was written in the letter | 5): From the passage we know the high postage made ...
A. people never send each other letters
B. lovers almost lose every touch with each other
C. people try their best to avoid paying it
D. receivers refuse to pay the coming letters |
| 3): We can know from Alice's words that ...
A. Tom had told her what the signs meant before leaving
B. Alice was clever and could guess the meaning of the signs
C. Alice had put the signs on the envelope herself
D. Tom had put the signs as Alice had told him to | Answer: ADABC |

Table 1: Sample reading comprehension problems from our dataset.

- **JASPER** (1992) système Reuters pour les nouvelles financières
- **MUC** (1987-1998) Message Understanding Conference
 - Naval operations messages
 - Terrorism in Latin American countries
 - Joint ventures and microelectronics domain.
 - News articles on management changes.
 - Satellite launch reports.
- **ACE** Automatic Content Extraction Conference
- **GATE** (1995) Université de Sheffield
- **OpenCalais** (2008) Reuters (service web)
- **TAC KBP** (2009-2017). Text Analysis Conference — Knowledge Base Population

Extraction d'information

Par règles

- des règles faites à la main pour capturer une information bien précise
- utilisée par exemple pour élaborer des **wrappers**
 - in texte html
 - out des entités (ou relations) d'intérêt
- possibilité d'apprendre les patrons (pour mieux généraliser)

Pro bonne précision


Cons très spécifique (à un domaine, à un type de marquage) \Rightarrow faible rappel

Exemple (simpliste)

The screenshot shows a web browser window with the URL www.dfa.ie/travel/travel-advice/a-z-list-of-countries/. A cookie notice is displayed at the top, followed by the DFA logo and navigation links. A dark green navigation bar contains links for 'Our Role & Policies', 'Passports & Citizenship', 'Travel', 'Embassies', 'News & Media', 'About Us', and 'Brexit'. Below this, a breadcrumb trail shows 'Travel Advice' > 'A-Z list of countries'. The main content area features a list of countries under the heading 'In this section:'. The list is split into two columns: the left column contains 'Afghanistan', 'Albania', 'Algeria', 'Andorra', and 'Angola'; the right column contains 'Afghanistan', 'Albania', 'Algeria', 'Andorra', 'Angola', 'Anguilla', 'Antigua & Barbuda', 'Argentina', and 'Armenia'.

Cookies on the DFA website

We use cookies to give the best experience on our site while also complying with Data Protection requirements. Continue without changing your settings, and you'll receive cookies, or change your cookie settings at any time. [Continue](#)

 **An Roinn Gnóthaí
Eachtracha agus Trádála**
Department of
Foreign Affairs and Trade

[Gaeilge](#) | [Our Ministers](#) | [Embassies](#)

[Our Role & Policies](#) | [Passports & Citizenship](#) | [Travel](#) | [Embassies](#) | [News & Media](#) | [About Us](#) | [Brexit](#)

[Travel Advice](#) > [A-Z list of countries](#)

In this section:

- Afghanistan
- Albania
- Algeria
- Andorra
- Angola

- Afghanistan
- Albania
- Algeria
- Andorra
- Angola
- Anguilla
- Antigua & Barbuda
- Argentina
- Armenia

Exemple (simpliste)

```
854         <!-- navigation object : 2017 Breadcrumbs --><li><a href="/travel/travel-advice/">Travel Advice</a></li><span class="breadcrumb_separat
855     </ul>
856     </nav>
857 <!--** // breadcrumbs ** -->
858 </div>
859 <!--** // content-wrap **-->
860 </section>
861 <!--** // inner-page-breadcrumbs row -->
862 <!--** Main body row **-->
863 <section class="main-body main-body--general-content">
864     <div class="content-wrap">
865         <aside class="main-body__side-bar">
866             <!--^^ navigation object : 2017 Inner navigation ^^-->
867             <h3 class="heading--side-bar">In this section:
868 <span class="inner-navigation__show-more__icon">
869     <svg class="dfa-svg-icon">
870         <use xmlns:xlink="http://www.w3.org/1999/xlink" xlink:href="#next_arrow"></use>
871     </svg>
872 </span>
873 <span class="inner-navigation__show-less__icon">
874     <svg class="dfa-svg-icon">
875         <use xmlns:xlink="http://www.w3.org/1999/xlink" xlink:href="#down_arrow"></use>
876     </svg>
877 </span>
878 </h3>
879 <ul class="inner-navigation"><li><a href="/travel/travel-advice/a-z-list-of-countries/afghanistan/">Afghanistan</a></li><li><a href="/travel/travel-advice/a-z-list-
880 <!-- Script adds show more/less links to inner navs with dropdowns -->
881 <script>
882     (function addShowMoreLessLinks() {
883         if ($('[class="multilevel-linkul"]').length > 0) {
884             var $mLUL = $('[class="multilevel-linkul"]');
885             /* Add show more. show less icons if there are nested links */
```

Exemple (simpliste)

more www.dfa.ie.html | grep Argentina



```
1. csh
X      csh      #1      X      csh      #2
/1i>li>ca href="/travel/travel-advice/a-z-list-of-countries/myanmar-burma/">Myanmar/Burma</li><li>ca href="/travel/travel-advice/a-z-list
-of-countries/namibia/">Namibia</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/nepal/">Nepal</li><li>ca href="/travel/tr
avel-advice/a-z-list-of-countries/new-zealand/">New Zealand</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/nicaragua/">Nicar
agua</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/niger/">Niger</li><li>ca href="/travel/travel-advice/a-z-list-of-coun
tries/nigeria/">Nigeria</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/norway/">Norway</li><li>ca href="/travel/travel-
advice/a-z-list-of-countries/oman/">Oman</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/pacific-islands/">Pacific Islands<a
></li><li>ca href="/travel/travel-advice/a-z-list-of-countries/pakistan/">Pakistan</li><li>ca href="/travel/travel-advice/a-z-list-of-coun
tries/panama/">Panama</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/papua-new-guinea/">Papua New Guinea</li><li>ca href
="/travel/travel-advice/a-z-list-of-countries/paraguay/">Paraguay</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/peru/">Peru
</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/the-philippines/">Philippines</li><li>ca href="/travel/travel-advice/a-z
-list-of-countries/poland/">Poland</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/portugal/">Portugal</li><li>ca href="/
travel/travel-advice/a-z-list-of-countries/puerto-rico/">Puerto Rico</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/qatar/">
Qatar</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/reunion/">Reunion</li><li>ca href="/travel/travel-advice/a-z-list-o
f-countries/romania/">Romania</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/russian-federation/">Russian Federation</li
><li>ca href="/travel/travel-advice/a-z-list-of-countries/rwanda/">Rwanda</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/sai
nt-kitts-and-nevis/">Saint Kitts and Nevis</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/saint-lucia/">Saint Lucia</li>
<li>ca href="/travel/travel-advice/a-z-list-of-countries/saint-vincents-and-grenadines/">Saint Vincent and the Grenadines</li><li>ca href
="/travel/travel-advice/a-z-list-of-countries/samoa/">Samoa</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/sao-tome-and-prin
cipe/">Sao Tome and Principe</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/saudi-arabia/">Saudi Arabia</li><li>ca href="/
travel/travel-advice/a-z-list-of-countries/senegal/">Senegal</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/serbia/">Serbia
</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/seychelles/">Seychelles</li><li>ca href="/travel/travel-advice/a-z-list-
of-countries/sierra-leone/">Sierra Leone</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/singapore/">Singapore</li><li>ca
 href="/travel/travel-advice/a-z-list-of-countries/slovakia/">Slovak Republic (Slovakia)</li><li>ca href="/travel/travel-advice/a-z-list-o
f-countries/slovenia/">Slovenia</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/solomon-islands/">Solomon Islands</li><li
>ca href="/travel/travel-advice/a-z-list-of-countries/somalia/">Somalia</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/south
-africa/">South Africa</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/south-sudan/">South Sudan</li><li>ca href="/travel
/travel-advice/a-z-list-of-countries/spain/">Spain</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/sri-lanka/">Sri Lanka</li>
<li>ca href="/travel/travel-advice/a-z-list-of-countries/sudan/">Sudan</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/su
riname/">Suriname</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/swaziland/">Swaziland</li><li>ca href="/travel/travel-a
dvice/a-z-list-of-countries/sweden/">Sweden</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/switzerland/">Switzerland</li>
<li>ca href="/travel/travel-advice/a-z-list-of-countries/syria/">Syria</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/tajik
istan/">Tajikistan</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/tanzania/">Tanzania</li><li>ca href="/travel/travel-ad
vice/a-z-list-of-countries/thailand/">Thailand</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/the-netherlands/">Netherlands<
/a><li><li>ca href="/travel/travel-advice/a-z-list-of-countries/timor-leste/">Timor Leste</li><li>ca href="/travel/travel-advice/a-z-list
-of-countries/togo/">Togo</li><li>ca href="/travel/travel-advice/a-z-list-of-countries/trinidad-and-tobago/">Trinidad & Tobago</li>
```

Exemple (simpliste)

```
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]more countries.html | sed -e 's/<li>/@/g' | tr '@' '\n' | grep a-z-list | head -n 10
      <li class="en lang-switcher secondary-nav-link"><a href="/travel/travel-advice/a-z-list-of-countries/">English -->
English -->
      <li class="en lang-switcher secondary-nav-link"><a href="/travel/travel-advice/a-z-list-of-countries/afghanistan/">Afghanistan</a></li>
i> <!-- English -->
<a href="/travel/travel-advice/a-z-list-of-countries/afghanistan/">Afghanistan</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/albania/">Albania</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/algeria/">Algeria</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/andorra/">Andorra</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/angola/">Angola</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/anguilla/">Anguilla</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/antigua-and-barbuda/">Antigua & Barbuda</a></li>
<a href="/travel/travel-advice/a-z-list-of-countries/argentina/">Argentina</a></li>
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]
```

Exemple (simpliste)

```
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]more countries.html | sed -e 's/<li>/@/g' | tr '@' '\n' | grep a-z-list | grep '^<a href="/trav
```

- [Afghanistan](/travel/travel-advice/a-z-list-of-countries/afghanistan/)
- [Albania](/travel/travel-advice/a-z-list-of-countries/albania/)
- [Algeria](/travel/travel-advice/a-z-list-of-countries/algeria/)
- [Andorra](/travel/travel-advice/a-z-list-of-countries/andorra/)
- [Angola](/travel/travel-advice/a-z-list-of-countries/angola/)
- [Anguilla](/travel/travel-advice/a-z-list-of-countries/anguilla/)
- [Antigua & Barbuda](/travel/travel-advice/a-z-list-of-countries/antigua-and-barbuda/)
- [Argentina](/travel/travel-advice/a-z-list-of-countries/argentina/)
- [Armenia](/travel/travel-advice/a-z-list-of-countries/armenia/)
- [Australia](/travel/travel-advice/a-z-list-of-countries/australia/)

```
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]
```

Exemple (simpliste)

```
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]more countries.html | sed -e 's/<li>/@/g' | tr '@' '\n' | grep a-z-list | grep '^/a>)*>\([^>]*\)<.*^1/1' | head -n 20
Afghanistan
Albania
Algeria
Andorra
Angola
Anguilla
Antigua & Barbuda
Argentina
Armenia
Australia
Austria
Azerbaijan
Bahamas
Bahrain
Bangladesh
Barbados
Belarus
Belgium
Belize
Benin
[felipe,/Users/felipe/home/Papers/Talks/Ecole-IVADO-2018]
```

- on peut espérer apprendre une règle du genre :
 $\langle \text{ul} \rangle \dots \langle \text{li} \rangle \langle \text{a href="*countries*" } \rangle \times \langle \text{/a} \rangle \langle \text{/li} \rangle \dots \langle \text{/ul} \rangle \Rightarrow X$

- acronymes :

NP ([A-Z]+)

⇒ This is called Open Information retrieval (OIE),
a task ...

- hyperonymes [**Hearst, 1992**] :

NP such as NP

NP is a NP

NP, including {NP, }* {or | and} NP

NP, especially NP

...

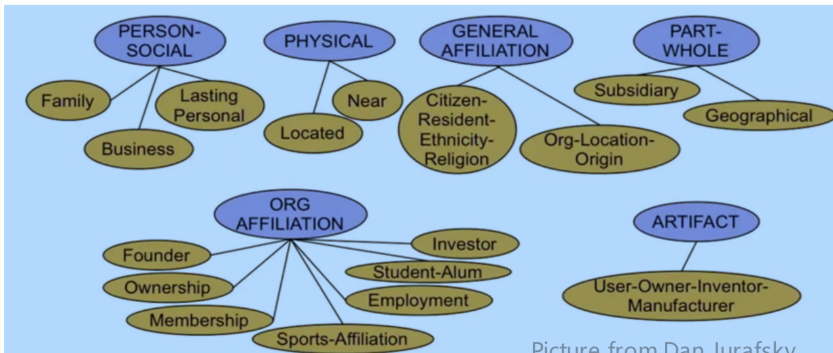
- ⇒
- riz, blé ∈ céréale
 - nitrogène, phosphore ∈ substance

Extraction d'information

Supervisée

Automated Context Extraction (ACE)

- beaucoup de travaux dédiés à reconnaître quelques dizaines de relations :



Picture from Dan Jurafsky

[GuoDong et al., 2005]

Features	P	R	F
Words	69.2	23.7	35.3
+Entity Type	67.1	32.1	43.4
+Mention Level	67.1	33.0	44.2
+Overlap	57.4	40.9	47.8
+Chunking	61.5	46.5	53.0
+Dependency Tree	62.1	47.2	53.6
+Parse Tree	62.3	47.6	54.0
+Semantic Resources	63.1	49.5	55.5

Table 2: Contribution of different features over 43 relation subtypes in the test data

- requiert des centaines d'exemples d'entraînement pour chaque relation
- pour des performances finalement modestes

Extraction d'information

Semi-supervisée

DIPRE (Brin, 1998)

in : \mathcal{D} une collection web de textes et **S** quelques exemples de paires d'entités validant la relation d'intérêt

- recherche d'auteurs et leurs livres, à l'aide du seed S :

Isaac Asimov	The Robots of Dawn
David Brin	Startide Rising
James Gleick	Chaos : Making a New Science
William Shakespeare	The Comedy of Errors

- depuis une collection \mathcal{D} de 24 millions de page Web (échantillonnées à chaque itération)

DIPRE (Brin, 1998)

in : \mathcal{D} une collection web de textes et **S** quelques exemples de paires d'entités validant la relation d'intérêt

1. **Chercher les occurrences** des paires de **S** dans \mathcal{D} ;
pour chacune d'elle retenir un 7-tuple
 $\langle \text{arg}_1, \text{arg}_2, \text{order}, \text{url}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - ✓ read The Adventures of Sherlock Holmes by Arthur Conan Doyle online or in your mail ...
 - ✓ I recommend you to read Startide by David Brin online at ...

in : \mathcal{D} une collection web de textes et S quelques exemples de paires d'entités validant la relation d'intérêt

1. **Chercher les occurrences** des paires de S dans \mathcal{D} ;
pour chacune d'elle retenir un 7-tuple
 $\langle \text{arg}_1, \text{arg}_2, \text{order}, \text{url}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - $\langle \text{ACD}, \text{TASH}, \text{false}, \text{url}_1, \text{read}, \text{by}, \text{online or} \rangle$
 - $\langle \text{DB}, \text{Startide}, \text{false}, \text{url}_2, \text{you to read}, \text{by}, \text{online at} \rangle$

in : \mathcal{D} une collection web de textes et S quelques exemples de paires d'entités validant la relation d'intérêt

1. **Chercher les occurrences** des paires de S dans \mathcal{D} ;
pour chacune d'elle retenir un 7-tuple
 $\langle \text{arg}_1, \text{arg}_2, \text{order}, \text{url}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - $\langle \text{ACD}, \text{TASH}, \text{false}, \text{url}_1, \text{read}, \text{by}, \text{online or} \rangle$
 - $\langle \text{DB}, \text{Startide}, \text{false}, \text{url}_2, \text{you to read}, \text{by}, \text{online at} \rangle$
2. **Induire des patrons** $\langle \text{order}, \text{urlprefix}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - en groupant les 7-tuples qui ont le même `order` et `middle`
 - et en gardant les “parties correspondantes” (texte et url)
 - ✓ $\langle \text{false}, \text{url}', \text{read}, \text{by}, \text{online} \rangle$
où url' est le préfixe commun à url_1 et url_2

- les 3 patrons générés lors de la première itération :
 - `< false, www.sff.net/locus/c., , + by, ()`
 `△ title by author (`
 - `< false,`
 `dns.city-net.com/~Imann/awards/hugos/1984.html`
 `<i>, </i> by, ()`
 `△ <i> title </i> by author (`
 - `< true,`
 `dolphin.upenn.edu/~dcummins/texts/sf-award.`
 `html, , ||, || ()`
 `△ author||title|| (`

in : \mathcal{D} une collection web de textes et S quelques exemples de paires d'entités validant la relation d'intérêt

1. **Chercher les occurrences** des paires de S dans \mathcal{D} ;
pour chacune d'elle retenir un 7-tuple
 $\langle \text{arg}_1, \text{arg}_2, \text{order}, \text{url}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - $\langle \text{ACD}, \text{TASH}, \text{false}, \text{url}_1, \text{read}, \text{by}, \text{online or} \rangle$
 - $\langle \text{DB}, \text{Startide}, \text{false}, \text{url}_2, \text{you to read}, \text{by}, \text{online at} \rangle$
2. **Induire des patrons** $\langle \text{order}, \text{urlprefix}, \text{prefix}, \text{middle}, \text{suffix} \rangle$
 - en groupant les 7-tuples qui ont le même `order` et `middle`
 - et en gardant les “parties correspondantes” (texte et url)
 $\checkmark \langle \text{false}, \text{url}', \text{read}, \text{by}, \text{online} \rangle$
où url' est le préfixe commun à url_1 et url_2
3. **Appliquer les patrons** à la collection $\mathcal{D} \Rightarrow S$.

Aller en 1 tant que nécessaire.

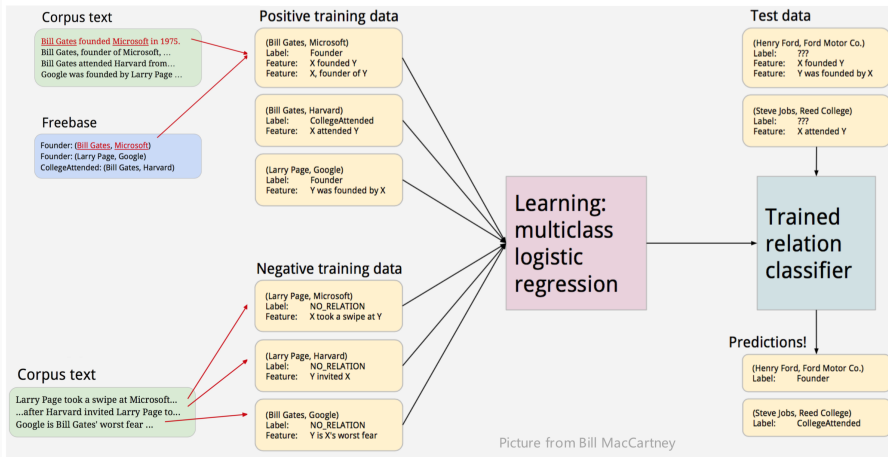
- 3 itérations sur un échantillon de la collection :
 1. 199 occurrences, 3 patrons, 4047 paires
 2. sur 5M. de documents : 3972 occurrences, 105 patrons, 9369 paires
 3. sur un corpus dédié de 156k documents : 15257 paires

H. D. Everett	The Death-Mask and Other Ghosts
H. G. Wells	First Men in the Moon
H. G. Wells	Science Fiction: Volume 2
H. G. Wells	The First Men in the Moon
H. G. Wells	The Invisible Man
H. G. Wells	The Island of Dr. Moreau
H. G. Wells	The Science Fiction Volume 1
H. G. Wells	The Shape of Things to Come: The Ultimate Revolution
H. G. Wells	The Time Machine
H. G. Wells	The War of the Worlds
H. G. Wells	When the Sleeper Wakes
H. M. Hoover	Journey Through the Empty
H. P. Lovecraft & August Derleth	The Lurker at the Threshold
H. P. Lovecraft	At the Mountains of Madness and Other Tales of Terror
H. P. Lovecraft	The Case of Charles Dexter Ward
H. P. Lovecraft	The Doom That Came to Sarnath and Other Stories

Très similaire à DIPRE, mais :

- patrons intégrant les types des entités en présence
 - ✓ ⟨org⟩'s headquarters in ⟨loc⟩
- mesure de confiance
 - des patrons détectés,
 - des paires initiales
 - **note** : requiert un corpus d'entraînement
- appliqué à l'extraction de paires (organisation, siège social)
 - ✓ (Microsoft, Redmond)
 - ✓ (Boeing, Seattle)

Supervision distante (Mintz et al. 2009)



Extraction d'information ouverte

Extraction d'information ouverte

Extracteurs

- **Idée** : OIE = **pas** de liste de relations pré-établie
- Différents types d'extracteurs (disponibles souvent pour l'anglais)
 - POS-based : TEXTRUNNER, WOE^{pos}, REVERB, SONEX
 - Dependency-based : WOE^{parse}, OLLIE, TREEKERNEL, PATTY, CLAUSIE, MINIE, GRAPHENE, OPENIE
 - Semantic role labeling : LUND, SWIRL, EXAMPLAR
- Le premier : TEXTRUNNER [**Banko et al., 2007**]

Des sorties variées [Niklaus et al., 2018]

OLLIE:

- (1) (Republican candidate Mitt Romney; will be elected President in; 2008)[enabler=If he w
- (2) (Republican candidate Mitt Romney; will be elected; President)[enabler=If he wins five
- (3) (Mitt Romney; be candidate of; Republican)
- (4) (Mitt Romney; be candidate for; Republican)
- (5) (he; wins; five key states)

ReVerb:

- (6) (he; wins; five key states)
- (7) (Republican candidate Mitt Romney; will be elected President in; 2008)

PredPatt:

- (8) (he; wins; five key states)
- (9) (Republican candidate Mitt Romney; will be elected President in; 2008)

ClausIE:

- (10) (he; wins; five key states)
- (11) (Republican candidate Mitt Romney; will be elected; President in 2008 If he wins five
- (12) (Republican candidate Mitt Romney; will be elected; President in 2008)

OpenIE 5.0:

- (13) (Republican candidate Mitt Romney; will be elected; President; T:in 2008)
- (14) (he; wins; five key states)

Graphene:

- (15) #1 CORE (Mitt Romney; will be elected; President)
- ("a) CONTEXT:NOUN_BASED Mitt Romney was a republican candidate .
- ("b) CONTEXT:TEMPORAL in 2008 .
- ("c) CONTEXT:CONDITION #3
- ("d) CONTEXT:NOUN_BASED #2
- (16) #2 CORE (Mitt Romney; was; a republican candidate)
- (17) #3 CONTEXT (he; wins; five key states)

Figure 1: Comparison of the output generated by different Open IE systems for the sentence "Mitt Romney will be elected President in 2008 if he wins five key states, Republican candidate Mitt Romney will be elected President in 2008".

- **idée** : utiliser les **infobox** de WIKIPEDIA comme supervision pour entraîner un extracteur

Pierre Lapointe

[✎](#) Pour les articles homonymes, voir [Lapointe](#).

Pierre Lapointe (né le 25 mai 1961 à Alma, au Québec) est un auteur-compositeur-interprète canadien (québécois).

Son œuvre s'inscrit d'une part dans la tradition de la chanson francophone, qu'elle soit québécoise ou française, ce qui se vérifie dans des textes souvent très littéraires ^{[[ren sealer](#)]}. Pierre Lapointe est également influencé par la musique pop qu'il entend utiliser pour renouveler la première chère [[]. Les arts graphiques contemporains, en particulier l'art numérique, sont très présents dans son univers via ses vidéos, et participent à la création d'un univers onirique et paradoxal, entre chansons mélancoliques et obscures, et scénographies colorées voire provocatrices.

Se définissant lui-même comme « chanteur populaire » ^[2,3], il s'est construit un personnage de dandy égoïste de fille artistique **montrealaise**⁴, lui permettant d'imprimer un décalage volontaire entre l'artiste sur scène et sa production largement biographique.

Ses disques enregistrent un succès critique et commercial au Canada ⁴.

Sommaire [[modifier](#)]

- 1 Biographie [[modifier](#)]
 - 1.1 Origines
 - 1.2 Débuts
 - 1.3 Les premiers succès
 - 1.4 La confirmation
- 2 Analyse de l'œuvre [[modifier](#)]
 - 2.1 Thèmes
 - 2.2 Influences
- 3 Discographie [[modifier](#)]
 - 3.1 Albums
 - 3.2 Participations
- 4 Distinctions
- 5 Notes et références [[modifier](#)]
- 6 Lien externe

Biographie [[modifier](#)] [[ajouter le code](#)]

Cette section ne cite pas suffisamment ses sources. Pour l'améliorer, ajoutez des références vérifiables ou les modèles [[référence nécessaire]] ou [[référence souhaitée]] sur les passages nécessitant une source.

Pierre Lapointe



Pierre Lapointe lors d'un spectacle donné au Collège de la Cité-Québec de Montréal (mars 2015).

Informations générales

Naissance	25 mai 1961 (55 ans) Alma, Québec, Canada
Activité principale	Chanteur, auteur-compositeur-interprète
Genre musical	Chanson française
Instruments	guitare, voix
Années actives	Depuis 1993
Labels	Audiogram
Site officiel	www.pierrelapointe.com

- WOE est décrit en 2 saveurs :
 - WOE^{parse} qui utilise un analyseur en dépendance
 - WOE^{pos} qui utilise les POS (comme TEXTRUNNER)

Woe : 1 - distant supervision

Pour chaque article a de WIKIPEDIA, pour chaque paire attribut/valeur (t, v) dans l'infobox de a :

1. matcher dans les phrases de a le titre de a et v

- présence de v :
 - correspondance exacte ou synonyme (via redirect, in-links) :
UK vs United Kingdom
- présence de *titre* :
 - correspondance exacte ou synonyme (redirects)
 - partial match : *Pierre vs Pierre Lapointe*
 - pronom le plus fréquent : *He vs Pierre Lapointe*
sauf si le pronom le plus fréquent est it

2. filtrer les exemples positifs

- une seule phrase dans a doit matcher v
- titre et v dans la même clause, etc.

Woe : 1 - distant supervision

- pour 1M. d'articles WIKIPEDIA : 301 962 exemples positifs

Pierre Lapointe



Pierre Lapointe lors d'un spectacle donné au Café théâtre le Côté-Cour de Jonquière (mars 2011).

Informations générales

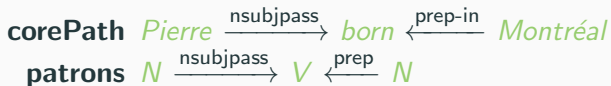
Naissance	23 mai 1981 (37 ans) Alma,  Québec,  Canada
Activité principale	Chanteur, auteur-compositeur-interprète
Genre musical	Chanson française
Instruments	piano, voix
Années actives	Depuis 2000
Labels	Audiogram
Site officiel	www.pierrelapointe.ca

- (activité principale, auteur-compositeur-interprète)
▲ **Pierre Lapointe**, né le 23 mai 1981 à Alma, au Québec, est un **auteur-compositeur-interprète** canadien (québécois).
- (labels, audiogram)
▲ En été 2003, **il** signe avec **Audiogram**.

Woe : 2 - extracteur (cas Woe^{parse})

1. analyser syntaxiquement (en dépendance) les phrases positives

Pierre was not born in Montreal



2. une base de patrons avec leur fréquence
 - 15 333 patrons, 185 ayant une fréq. ≥ 100 , 1929 ≥ 5
 - prob d'un patron est une fonction de la fréquence de ce patron
3. les dépendances non *core* (adjectivales, négations, etc.) sont ajoutées lors de l'extraction $\langle X, \text{was not born in}, Y \rangle$

- 900 phrases : 300 du WSJ, de WIKIPEDIA et du Web
- **référence** : annotation manuelle de tous les triplets
 - $WOE^{parse} > WOE^{pos} > \text{TEXTRUNNER}$
 - nb moy. de tuples par phrase :
 - WOE^{parse} 1.42
 - WOE^{pos} 1.05
 - TEXTRUNNER 0.75
 - temps moyen pour traiter une phrase :
 - WOE^{parse} 0.679s
 - WOE^{pos} et TEXTRUNNER : 0.022s
(WOE^{parse} est 30 fois plus lent)

- 3 modules
 1. identification de relation verbale
 2. identification des arguments (NP) à gauche et à droite
 3. estimation de confiance (maxent)
- Évaluation
 - 500 pages Web retournées par <http://random.yahoo.com/bin/ryl> soumises à plusieurs extracteurs et dont les sorties sont évaluées manuellement

ReVerb : Extraction des relations

- contraintes **syntaxiques** :

V | V P | V W* P

V = verb particle? adv?

W = (noun | adj | adv | pron | det)

P = (prep | particle | inf. marker)

invented (V), *located in* (VP), *has atomic weight of* (VW*P)

+ joindre deux relations adjacentes

ReVerb : Extraction des relations

- contraintes **syntaxiques** :

V | V P | V W* P

V = verb particle? adv?

W = (noun | adj | adv | pron | det)

P = (prep | particle | inf. marker)

invented (V), *located in* (VP), *has atomic weight of* (VW*P)

+ joindre deux relations adjacentes

- contraintes **lexicales** :

- éliminer les relations trop spécifiques

is offering only modest greenhouse gas reduction targets at

- une relation doit avoir plusieurs arguments dans une collection
 1. extraire toutes les relations et leurs arguments
 2. garder celles qui ont au moins 20 paires d'arguments différentes

ReVerb : Extraction des arguments

Pour chaque relation r :

- trouver x le NP le plus proche de r à gauche
 - qui n'est pas un pronom relatif, un adverbe WHO, ou there
- trouver y le NP le plus proche à droite de r

TOKS	Elle reprend ses études à l' Université McGill et obtient un doctorat en psychologie (1965) .
POS	CLS V DET NC P DET NC NPP CC V DET NC P NC PONCT NC PONCT PONCT
CHNK	B-NP B-VN B-NP I-NP B-PP B-NP I-NP I-NP B-COORD B-VN B-NP I-NP B-PP B-NP O B-NP O O

ReVerb : Extraction des arguments

Pour chaque relation r :

- trouver x le NP le plus proche de r à gauche
 - qui n'est pas un pronom relatif, un adverbe WHO, ou there
- trouver y le NP le plus proche à droite de r

TOKS	Elle reprend ses études à l' Université McGill et obtient un doctorat en psychologie (1965) .
POS	CLS V DET NC P DET NC NPP CC V DET NC P NC PONCT NC PONCT PONCT
CHNK	B-NP B-VN B-NP I-NP B-PP B-NP I-NP I-NP B-COORD B-VN B-NP I-NP B-PP B-NP O B-NP O O
EXTR	Elle == reprend == ses études

ReVerb : Extraction des arguments

Pour chaque relation r :

- trouver x le NP le plus proche de r à gauche
 - qui n'est pas un pronom relatif, un adverbe WHO, ou there
- trouver y le NP le plus proche à droite de r

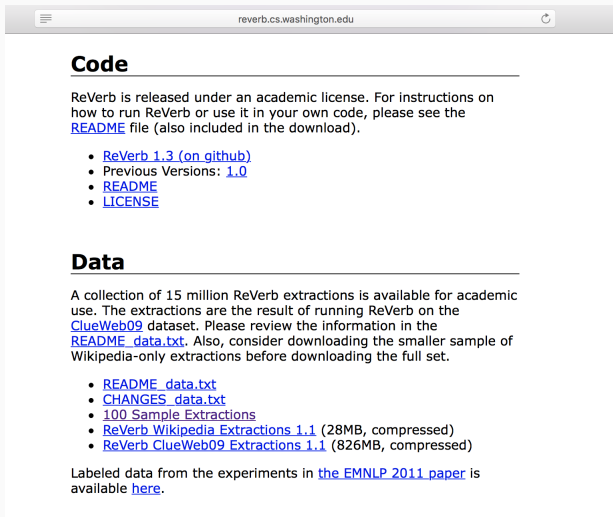
TOKS	Elle reprend ses études à l' Université McGill et obtient un doctorat en psychologie (1965) .
POS	CLS V DET NC P DET NC NPP CC V DET NC P NC PONCT NC PONCT PONCT
CHNK	B-NP B-VN B-NP I-NP B-PP B-NP I-NP I-NP B-COORD B-VN B-NP I-NP B-PP B-NP O B-NP O O
EXTR	Elle == reprend == ses études
EXTR	ses études == obtient == un doctorat

ReVerb : Estimateur de confiance

- **in** : (x, r, y) **out** : $p(\text{ok})$
- entraînement : référence manuelle de 1000 phrases
- 19 features qui ne dépendent pas d'une relation en particulier :

1.16 (x, r, y) couvre tous les mots de s
0.50 la dernière prép dans r est *for*
0.49 la dernière prép dans r est *on*
0.46 la dernière prép dans r est *of*
0.43 $\text{length}(s) \leq 10$ mots
-0.93 conjonction de coord. à gauche de r dans s
⋮

ReVerb est disponible (pour l'anglais)



The screenshot shows a web browser window with the address bar containing "reverb.cs.washington.edu". The page content is divided into two main sections: "Code" and "Data".

Code

ReVerb is released under an academic license. For instructions on how to run ReVerb or use it in your own code, please see the [README](#) file (also included in the download).

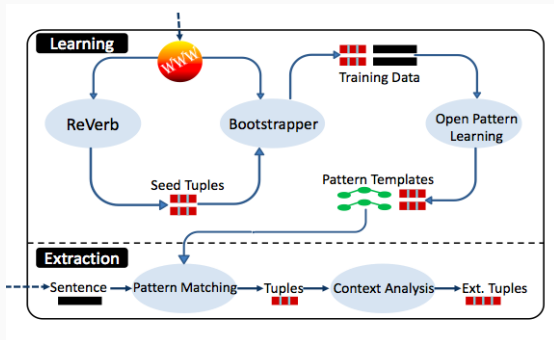
- [ReVerb 1.3 \(on github\)](#)
- Previous Versions: [1.0](#)
- [README](#)
- [LICENSE](#)

Data

A collection of 15 million ReVerb extractions is available for academic use. The extractions are the result of running ReVerb on the [ClueWeb09](#) dataset. Please review the information in the [README_data.txt](#). Also, consider downloading the smaller sample of Wikipedia-only extractions before downloading the full set.

- [README_data.txt](#)
- [CHANGES_data.txt](#)
- [100 Sample Extractions](#)
- [ReVerb Wikipedia Extractions 1.1](#) (28MB, compressed)
- [ReVerb ClueWeb09 Extractions 1.1](#) (826MB, compressed)

Labeled data from the experiments in [the EMNLP 2011 paper](#) is available [here](#).



- requiert des tuples extraits par REVERB
- relations nominales : *Microsoft co-founder Bill Gates said...*
 - ⟨Bill Gates, is cofounder of, Microsoft⟩

- ~110k tuples fiables extraits par REVERB depuis CLUEWEB :
 - $\text{freq} \geq 2$
 - args = noms propres
 - confiance élevée (selon REVERB)
- 18M des phrases de CLUEWEB contiennent ces tuples
- 4M une fois filtrées (mots de têtes des args reliés par un chemin de dépendances de long. au plus 4)
- **hyp** : ces phrases expriment des relations pertinentes
vrai à 90% selon une éval. manuelle sur 100 phrases

Extraction Template	Open Pattern
1. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓{prep.*}↓ {arg2}
2. (arg1; {rel}; arg2)	{arg1} ↑nsubj↑ {rel:postag=VBD} ↓dobj↓ {arg2}
3. (arg1; be {rel} by; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓agent↓ {arg2}
4. (arg1; be {rel} of; arg2)	{rel:postag=NN; type=Person } ↑nn↑ {arg1} ↓nn↓ {arg2}
5. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {slot:postag=VBN; lex ∈ announce name choose... } ↓dobj↓ {rel:postag=NN} ↓{prep.*}↓ {arg2}

pris de [Mausam et al., 2012]

Analyse en dépendance des 4M de phrases retenues et extraction de patrons syntaxiques (droite) exprimant des relations (gauche) :

- des patrons purement syntaxiques (ex : 1 à 3)
 - généralisation agressive (*on*, *on*, *for*, ... = prep)

Extraction Template	Open Pattern
1. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓{prep.*}↓ {arg2}
2. (arg1; {rel}; arg2)	{arg1} ↑nsubj↑ {rel:postag=VBD} ↓dobj↓ {arg2}
3. (arg1; be {rel} by; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓agent↓ {arg2}
4. (arg1; be {rel} of; arg2)	{rel:postag=NN; type=Person } ↑nn↑ {arg1} ↓nn↓ {arg2}
5. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {slot:postag=VBN; lex ∈ announce name choose... } ↓dobj↓ {rel:postag=NN} ↓{prep.*}↓ {arg2}

pris de [Mausam et al., 2012]

Analyse en dépendance des 4M de phrases retenues et extraction de patrons syntaxiques (droite) exprimant des relations (gauche) :

- des patrons purement syntaxiques (ex : 1 à 3)
 - généralisation agressive (*on, on, for, ...* = prep)
- des patrons lexicalisés (ex : 5)

Extraction Template	Open Pattern
1. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓{prep.*}↓ {arg2}
2. (arg1; {rel}; arg2)	{arg1} ↑nsubj↑ {rel:postag=VBD} ↓dobj↓ {arg2}
3. (arg1; be {rel} by; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓agent↓ {arg2}
4. (arg1; be {rel} of; arg2)	{rel:postag=NN; type=Person } ↑nn↑ {arg1} ↓nn↓ {arg2}
5. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {slot:postag=VBN; lex ∈ announce name choose... } ↓dobj↓ {rel:postag=NN} ↓{prep.*}↓ {arg2}

pris de [Mausam et al., 2012]

Analyse en dépendance des 4M de phrases retenues et extraction de patrons syntaxiques (droite) exprimant des relations (gauche) :

- des patrons purement syntaxiques (ex : 1 à 3)
 - généralisation agressive (*on, on, for, ...* = prep)
- des patrons lexicalisés (ex : 5)
- des patrons typés (ex : 4)
 - généralisation des patrons lexicalisés grâce à WORDNET

Ollie : Extraction

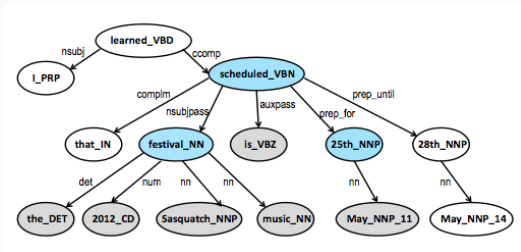


Fig 4 dans l'article

I learned that the 2012 Sasq. music festival is scheduled for May 25th until May 28th

- template : (arg ; be {rel} {prep} ; arg2)
- patron : {arg1} $\xrightarrow{nsubjpass}$ {rel :postag=VBN} \xleftarrow{prep} {arg2}

1. match arg1=*festival* arg2=*25th* et rel=*scheduled*
<festival, be scheduled for, 25th>

Ollie : Extraction

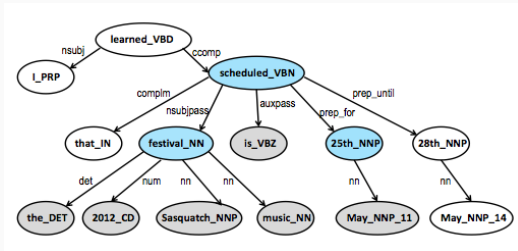


Fig 4 dans l'article

I learned that the 2012 Sasq. music festival is scheduled for May 25th until May 28th

- template : (arg ; be {rel} {prep} ; arg2)
- patron : {arg1} $\xrightarrow{nsbjpass}$ {rel :postag=VBN} \xleftarrow{prep} {arg2}

1. match $arg1=festival$ $arg2=25th$ et $rel=scheduled$
<festival, be scheduled for, 25th>
2. étendre les nœuds *args* et *rel* à certains mots gouvernés
<the Sasquatch music festival, be scheduled for, May 25th>

OLLIE détecte 2 cas d'extractions non factuelles :

- *Early astronomers believed that the earth is the center of the universe*

R : ⟨the earth, be the center of, the universe⟩

O : ⟨the earth, be the center of, the universe⟩

AttributedTo : believe ; Early astronomers

- *If he wins five states, Romney will be elected President*

R : ⟨Romney, will be elected, President⟩

O : (⟨Romney, will be elected, President⟩

ClausalModifier if ; he wins five key states)

- simplifier la structure syntaxique pour faciliter l'extraction

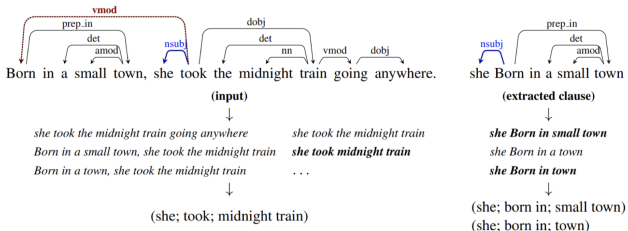


Figure 4: An illustration of Stanford Open IE's approach. From left to right, a sentence yields a number of independent clauses. From top to bottom, each clause produces a set of entailed shorter utterances, and segments the ones which match an atomic pattern into a relational triple (Angeli et al., 2015).

Pris de [Niklaus et al., 2018]

Triplet extractors

Chilly Gonzales (born Jason Charles Beck; 20 March 1972) is a Canadian musician who resided in Paris, France for several years, and now lives in Cologne, Germany.

Extract the triples

Extract with Ollie debug

Extract with Distylium debug

REVERB

20 March 1972 is a Canadian musician 0,56

a Canadian musician resided in Paris 0,79

France now lives in Cologne 0,87

OLLIE

20 March 1972 is a Canadian musician for several years 0,86

20 March 1972 now lives in Cologne 0,65

a Canadian musician resided in Paris 0,76

20 March 1972 is a Canadian musician who resided in Paris 0,74

DISTYLIUM

a Canadian musician resided in Paris 56832,97

a Canadian musician resided for several years 56832,80

a Canadian musician now lives in Cologne 22315,53

Chilly Gonzales is a Canadian musician 10000,53

Chilly Gonzales is a Canadian musician 38,63

a Canadian musician resided lives 24,87

Triplet extractors

Chilly Gonzales (born Jason Charles Beck; 20 March 1972) is a Canadian musician who resided in Paris, France for several years, and now lives in Cologne, Germany.

Extract the triples

Extract with Ollie debug

Extract with Distylum debug

CLAUSIE

Chilly Gonzales born Jason Charles Beck 20 March 1972 is a Canadian musician -254,85

a Canadian musician resided in Paris France for several years -254,85

a Canadian musician lives in Cologne Germany now -254,85

MINIE

Chilly Gonzales is Canadian musician 0,00

Canadian musician resided in Paris for QUANT_O_1 years | QUANT:[O_1=several] 0,00

Canadian musician resided in Paris 0,00

Canadian musician lives in Cologne 0,00

Paris is in France 0,00

Cologne is in Germany 0,00

STANFORD

Gonzales is Canadian 1,00

Triplet extractors

Chilly Gonzales (born Jason Charles Beck; 20 March 1972) is a Canadian musician who resided in Paris, France for several years, and now lives in Cologne, Germany.

Extract the triples

Extract with Ollie debug

Extract with Distylium debug

OPENIE

Chilly Gonzales (born is a Canadian musician who resided in Paris, France for several years 0,97

a Canadian musician resided in Paris, France for several years 0,91

PROPS

subj:Chilly Gonzales (born Jason Charles Beck ; 20 March 1972) is a Canadian musician and now lives in Cologne , Germany resided prep_in:Paris , France -2

subj:Chilly Gonzales (born Jason Charles Beck ; 20 March 1972) is a Canadian musician and now lives in Cologne , Germany resided prep_for:several years -2

Extraction d'information ouverte

Limites

- environ 50% des tuples

Jardin zoologique du Québec , reopened in 2002 after two years of restorations but closed in 2006 after a political decision . It featured 750 specimens of 300 different species of animals . **The zoo specialized in winged fauna and garden themes** , but also presented several species of mammals .

⟨The zoo, specialized in, winged fauna and garden themes⟩

- environ 50% des tuples

Jardin zoologique du Québec , reopened in 2002 after two years of restorations but closed in 2006 after a political decision . It featured 750 specimens of 300 different species of animals . **The zoo specialized in winged fauna and garden themes** , but also presented several species of mammals .

⟨The zoo, specialized in, winged fauna and garden themes⟩

- environ 50% des tuples

Porte St-Louis and Porte St-Jean are the main gates through the walls from the modern section of downtown ; the Kent Gate was a gift to the province from Queen Victoria and **the foundation stone was laid by the Queen 's daughter** , Princess Louise , Marchioness of Lorne , on June 11 , 1879 .

⟨**the foundation stone**, was laid by, the Queen 's daughter⟩

- environ 50% des tuples

Porte St-Louis and Porte St-Jean are the main gates through the walls from the modern section of downtown ; **the Kent Gate** was a gift to the province from Queen Victoria and **the foundation stone was laid by the Queen 's daughter** , Princess Louise , Marchioness of Lorne , on June 11 , 1879 .

⟨**the foundation stone**, was laid by, the Queen 's daughter⟩

Quebec City 's skyline is dominated by the massive Château Frontenac Hotel , perched on top of Cap-Diamant . **It was designed by architect Bruce Price** , as one of a series of “*château* ” style hotels built for the Canadian Pacific Railway company .

⟨**It**, was designed by, architect Bruce Price⟩

Quebec City 's skyline is dominated by the massive **Château Frontenac Hotel** , perched on top of Cap-Diamant . **It was designed by architect Bruce Price** , as one of a series of “**château** ” style hotels built for the Canadian Pacific Railway company .

⟨**It**, was designed by, architect Bruce Price⟩

D'autres problèmes

- tuples non. informatifs :
 - ⟨The town, distinguished, itself⟩
 - ⟨One-quarter of the people, were members of, religious⟩
 - ⟨The team, has, size league titles⟩
- syntaxe compliquée

The English-speaking community peaked in relative terms during the 1860s , when 40 % of Quebec City 's residents were Anglophone .

⟨The English-speaking community, peaked in, relative terms⟩

Et pourtant ça marche !

- si appliqué à de larges collections (M. de phrases)
- en appliquant des filtres agressifs :
tuples vu au moins n fois, score de confiance élevé, etc.

⟨Quebec, was founded by, Samuel de Champlain⟩

⟨Quebec City, is located on, the north bank of the Saint Lawrence River⟩

⟨Quebec City, is an important hub in, the province 's autoroute system⟩

⟨Quebec City, is located in, the Saint Lawrence River valley⟩

⟨Quebec City, was struck by, the 1925 Charlevoix-Kamouraska earthquake⟩

Intermède

Intermède

deFacto

1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
48h plus tard : 20M tuples une fois filtrés

Extraction des tuples

...

CANO l' activité biologique == se dérouler dans == le lac

CANO les sédiments == être appeler == relargage

CANO sédiments == jouer == un rôle majeur

CANO l' eau == recouvrer == les sédiments

CANO les sédiments == contenir de == l' oxygène

CANO les échanges de P == être généralement == unidirectionnels de l' eau

CANO les conditions anoxiques == modifier == la nature des échanges chimiques

CANO la décomposition == anaérobie de == la matière organique

CANO La diminution du potentiel d' oxydoréduction == découler de == cette anoxie

CANO cette anoxie == favoriser == la remise

CANO biologiques == contribuer à == la libération du P

CANO Cette combinaison de facteurs == faire en == sorte

CANO les sédiments == ne pas être uniquement == un lieu d' accumulation du P

CANO le P == pouvoir être == maintes reprises

CANO Le lac Nairne == être affecter par == des épisodes de floraison de cyanobactéries annuellement

CANO des floraisons == être observer sur == toute la superficie du lac

CANO Certains épisodes == survenir en == été

CANO Ces floraisons massives == être == un symptôme de l' eutrophisation du lac

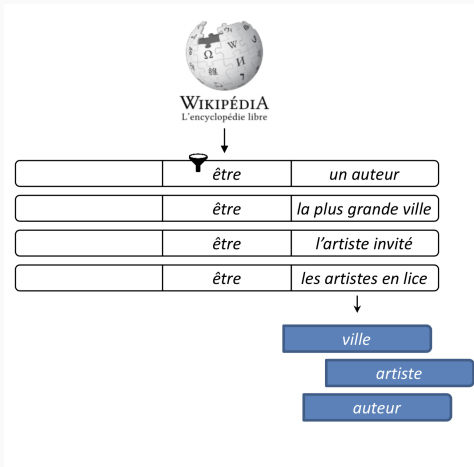
CANO Le lac Nairne == être à == coeur de la municipalité de Saint-Aimé-des-Lacs

CANO un plan d' action == être mettre sur == pied

CANO plusieurs actions == être entreprendre comme == la caractérisation

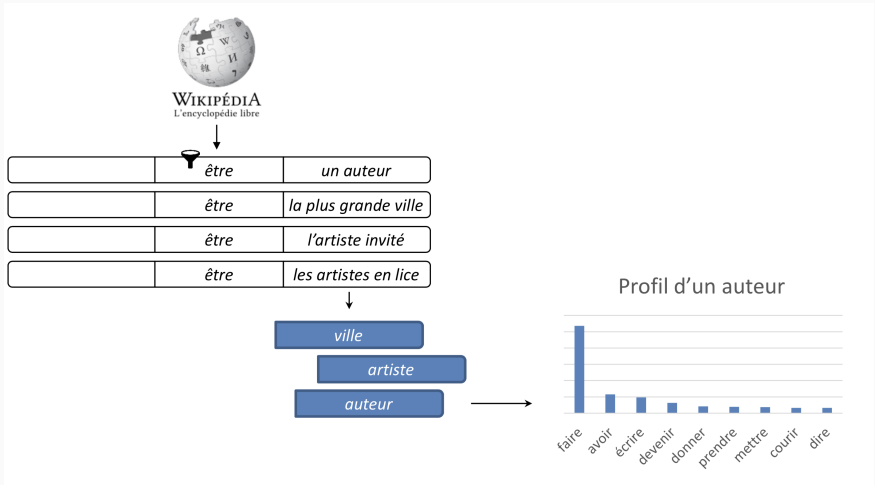
1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
2. arg_{2s} les plus fréquents de la relation *être* promus *catégories*

Sélection de catégories

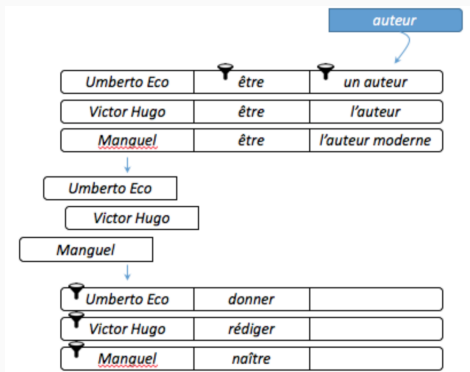


1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
2. arg_{2s} les plus fréquents de la relation *être* promus **catégories**
3. calcul du **profil relationnel** de chaque catégorie : les relations les plus discriminantes impliquant des instances de la catégorie **Philosophe** → *affirmer, appeler, considérer, décrire, défendre, développer, fonder, publier, reprendre, écrire, etc.*

Profil relationnel



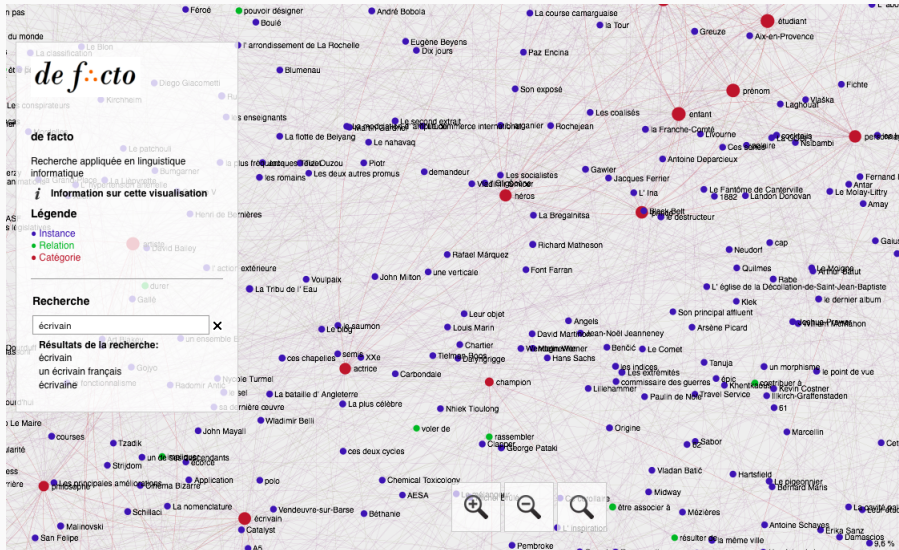
Profil relationnel

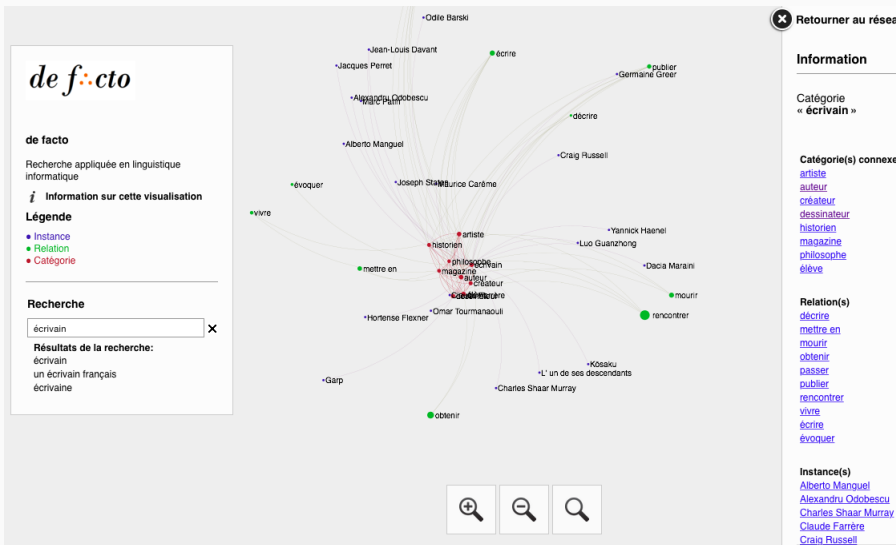


1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
2. arg_{2s} les plus fréquents de la relation *être* promus **catégories**
3. calcul du **profil relationnel** de chaque catégorie : les relations les plus discriminantes impliquant des instances de la catégorie
4. sélection des **instances** les plus représentatives de chaque catégorie à l'aide de leur profil
Acteur → David Arquette, Jeremy Bulloch, Fernand Gravey, etc.

1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
2. arg₂s les plus fréquents de la relation *être* promus **catégories**
3. calcul du **profil relationnel** de chaque catégorie : les relations les plus discriminantes impliquant des instances de la catégorie
4. sélection des **instances** les plus représentatives de chaque catégorie à l'aide de leur profil
5. *clustering* (K-means) des catégories à l'aide de leur profil
Écrivain → *artiste, auteur, créateur, dessinateur, historien, magazine, philosophe, élève*

1. extraction de tuples depuis 31M de phrases de WIKIPEDIA-FR (2014) à l'aide d'une adaptation de REVERB au français
2. arg_{2s} les plus fréquents de la relation *être* promus *catégories*
3. calcul du **profil relationnel** de chaque catégorie : les relations les plus discriminantes impliquant des instances de la catégorie
4. sélection des *instances* les plus représentatives de chaque catégorie à l'aide de leur profil
5. *clustering* (K-means) des catégories à l'aide de leur profil
6. génération d'un graphe navigable avec un *plugin* (*sigma.js*) modifié pour l'occasion





- 3 types de jeux pour aider à nettoyer la base de DEFACTO



Érudit célèbre ses 20 ans !

Recherche

Par auteur, titre, mots-clés...



- [Revues](#)
- [Livres et actes](#)
- [Thèses et mémoires](#)
- [Rapports de recherche](#)

[Recherche avancée](#)

Derniers numéros

Géographie



Les enjeux sociaux de l'eau : comparaisons

Philosophie, Théologie



Volume 74, numéro 1, février 2018

De la revue [Laval théologique et philosophique](#)

Anthropologie et ethnologie



Récits de savoirs partagés par l'art et la création en milieux autochtones

Géographie, Études urbaines



La présence - absence des études urbaines en France
Volume 13, 2018

De la revue [Environnement Urbain](#)

Histoire



Numéro 180, mai-août 2018

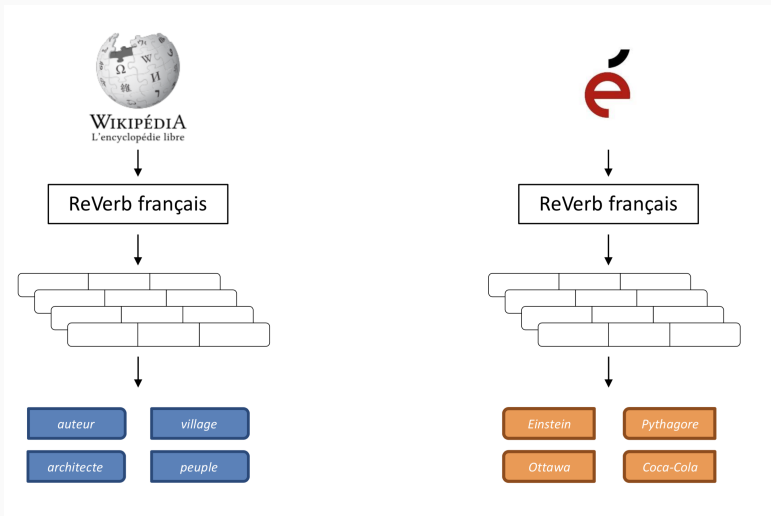
Education



Rapport à l'écrit, contextes de for (Vol. 2)
Volume 20, num 2017

De la revue [Now](#)

- Peut-on ré-utiliser le typage émergeant de DEFACTO (WIKIPEDIA) pour étiqueter des entités extraites d'Érudit ?
 - 3M de phrases (français)
 - REVERB : 5M de triplets
 - TEST : 100 entités vues au moins 50 fois dans Érudit étiquetées à la main





<i>Michel Foucault</i>	<i>dégager</i>	<i>4 caractères</i>
<i>le village</i>	<i>compter</i>	<i>3000 habitants</i>
<i>Michel Foucault</i>	<i>commenter</i>	<i>la philo. pénale</i>
<i>le terme</i>	<i>renvoyer</i>	<i>à cette personne</i>



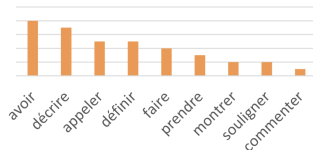
1960

Belgique

Michel Foucault



Profil de Michel Foucault



70% de bon étiquetage

Instance	Type véritable	Type prédit
1960	année, période	période
Aragon	artiste, auteur, écrivain	auteur
article	étude, livre, recueil	étude
Belgique	pays, lieu, lieux, toponyme	ville
Canadiens français	peuple, groupe	personnages
Allemagne	pays, lieu, lieux, toponyme	ville
Santé Canada	organisation, association, organisme, groupe	médecin
Michel Foucault	philosophe, écrivain, auteur	auteur
Premier ministre	ministre, président, maître	ministre
Socrate	philosophe, écrivain, auteur	ville
texte	étude, livre, recueil	livre

Intermède

Allium



Prototype **en cours** de développement réalisé dans le cadre d'un projet FCI piloté par Vincent Larivière (EBSI) dans le but d'améliorer la plateforme Érudit.

6 | EN COUVERTURE | ENTREVUE



SYLVAIN L'ESPÉRANCE

« LA GRÈCE, C'EST CE QUI NOUS ATTEND »

Dévoilé aux Rencontres internationales du documentaire de Montréal, en novembre 2016, **Combat au bout de la nuit** est un film-fleuve porté par le désir de lutter contre les mesures d'austérité des gouvernements. Sylvain L'Espérance s'est tenu à le tourner en Grèce, mais il le considère comme un écho de l'état du monde. Sa présentation dans la section Panorama de la 67^e Berlinale, en février 2017, le ravit particulièrement parce qu'on lui a annoncé avoir voulu faire, avec cette sélection, « un geste politique ». Le documentariste québécois ne pouvait mieux être servi, lui qui cherche, derrière sa caméra imprégnée de poésie, à suggérer des formes de résistance.

JÉRÔME DELGADO



Numéro 307, Mars, 2017, p. 6-9

Combat au bout de la nuit

Tous droits réservés © La revue Séquences Inc., 2017



CONCEPTS

AUTEURS

CONNEXES

LIEUX

Grèce 43,7

Athènes 25,4

AUTRES CONCEPTS

Espérance 47,8

Combat 40,9

Sylvain L' Espérance 35,8

film 24,2

révolte 19,6

bout de la nuit 16,4

▲ CONCEPTS AUTEURS CONNEXES

Aussi par Jérôme Delgado

Alexandre David

Par Jérôme Delgado

The Clock

Par Jérôme Delgado

Suzie

Par Jérôme Delgado

Chaotic Anna

Par Jérôme Delgado




▲ CONCEPTS AUTEURS **CONNEXES**

Documents apparentés

Entretien avec Sylvain L'Espérance Par Robert Daudelin et Gérard Grugeau	Le siècle des migrants Par Robert Daudelin	Les écarts essentiels : le cinéma de Sylvain L'Espérance Par Robert Daudelin	TOP 10 – 2016
--	--	--	----------------------

< ————— >

6 | EN COUVERTURE | ENTREVUE



**SYLVAIN
L'ESPÉRANCE
« LA GRÈCE,
C'EST CE QUI**

🔖
👤
➦
▲

Structuration

- Extracteur similaire à OLLIE où les arguments sont restreints aux seules entités nommées de YAGO
- tuples généralisés en patrons Syntactic-Ontological-Lexical
 - séquence de mots, POS et *
 - *⟨person⟩'s [adj] voice * ⟨song⟩* matche :
Amy Winehouse's soft voice in 'Rehab'
Elvis Presley's solid voice in this song 'All shook up'
- structuration selon la généralité des patrons, et leur synonymie
⟨person⟩ winner of ⟨award⟩ ⇒ ⟨person⟩ nominated for ⟨award⟩

Relation	Paraphrases	Precision	Sample Paraphrases
DBPedia/artist	83	0.96±0.03	[adj] studio album of, [det] song by ...
DBPedia/associatedBand	386	0.74±0.11	joined band along, plays in ...
DBPedia/doctoralAdvisor	36	0.558±0.15	[det] student of, under * supervision ...
DBPedia/recordLabel	113	0.86±0.09	[adj] artist signed to, [adj] record label ...
DBPedia/riverMouth	31	0.83±0.12	drains into, [adj] tributary of ...
DBPedia/team	1,108	0.91±0.07	be * traded to, [prp] debut for ...
YAGO/actedIn	330	0.88±0.08	starred in * film, [adj] role for ...
YAGO/created	466	0.79±0.10	founded, 's book ...
YAGO/isLeaderOf	40	0.53±0.14	elected by, governor of ...
YAGO/holdsPoliticalPosition	72	0.73±0.10	[prp] tenure as, oath as ...

Table 6: Sample Results for Relation Paraphrasing

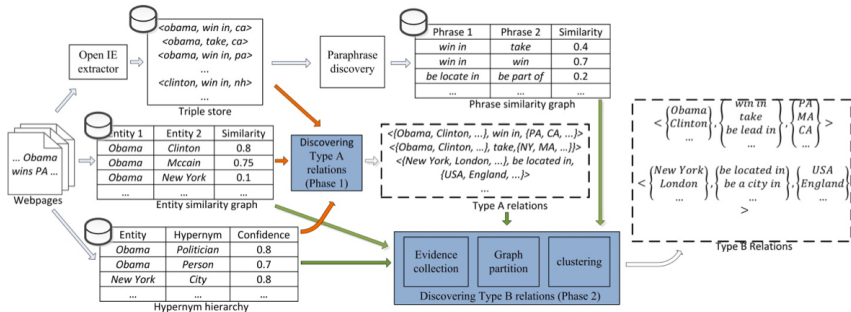


Figure 1. Overview of the WEBRE algorithm (Illustrated with examples sampled from experiment results). The tables and rectangles with a database sign show knowledge sources, shaded rectangles show the 2 phases, and the dotted shapes show the system output, a set of Type A relations and a set of Type B relations. The orange arrows denote resources used in phase 1 and the green arrows show the resources used in phase 2.

- 14.7M de triplets distincts extraits par REVERB depuis CLUEWEB (filtres appliqués)
 - 1.3M relations verbales, 3.3M entités (args)
- 0.2M relations de type A, 84 000 de type B

Argument 1	Relation phrase	Argument 2
<i>marijuana, cafeine, nicotine...</i>	<i>result in, be risk factor for, be major cause of...</i>	<i>insomnia, emphysema, breast cancer, ...</i>
<i>C# 2.0, php5, java, c++, ...</i>	<i>allow the use of, also use, introduce the concept of...</i>	<i>destructors, interfaces, template, ...</i>
<i>clinton, obama, mccain, ...</i>	<i>win, win in, take, be lead in, ...</i>	<i>ca, dc, fl, nh, pa, va, ga, il, nc, ...</i>

Table 3. Sample Type B relations extracted.

- Extraction de schémas d'événements à partir de textes

Actor	Rel	Actor
A1:<person>	failed	A2:test
A1:<person>	was suspended for	A3:<time period>
A1:<person>	used	A4:<substance, drug>
A1:<person>	was suspended for	A5:<game, activity>
A1:<person>	was in	A6:<location>
A1:<person>	was suspended by	A7:<org, person>
Actor Instances:		
A1: {Murray, Morgan, Governor Bush, Martin, Nelson}		
A2: {test}		
A3: {season, year, week, month, night}		
A4: {cocaine, drug, gasoline, vodka, sedative}		
A5: {violation, game, abuse, misfeasance, riding}		
A6: {desert, Simsbury, Albany, Damascus, Akron}		
A7: {Fitch, NBA, Bud Selig, NFL, Gov Jeb Bush}		

Table 1: An event schema produced by our system, represented as a set of (*Actor*, *Rel*, *Actor*) triples, and a set of instances for each actor *A1*, *A2*, etc. For clarity we show unstemmed verbs.

Applications

Never Ending Language Learning (NELL)

rtw.ml.cmu.edu/rtw/

Read the Web

Research Project at Carnegie Mellon University

Home

Project Overview

Resources & Data

Publications

People

NELL: Never-Ending Language Learning

Can computers learn to read? We think so. "Read the Web" is a research project that attempts to create a computer system that learns over time to read the web. Since January 2010, our computer system called NELL (Never-Ending Language Learner) has been running continuously, attempting to perform two tasks each day:

- First, it attempts to "read," or extract facts from text found in hundreds of millions of web pages (e.g., `playsInstrument(George_Harrison, guitar)`).
- Second, it attempts to improve its reading competence, so that tomorrow it can extract more facts from the web, more accurately.

So far, NELL has accumulated over 50 million candidate beliefs by reading the web, and it is considering these at different levels of confidence. NELL has high confidence in 2,810,379 of these beliefs — these are displayed on this website. It is not perfect, but NELL is learning. You can track NELL's progress below or [@cmunell on Twitter](#), browse and download its [knowledge base](#), read more about our [technical approach](#), or join the [discussion group](#).



Recently-Learned Facts [twitter](#)

Refresh

instance	iteration	date learned	confidence	
manuel_hedilla_larrey is a South American person	1111	06-jul-2018	93.7	
gloria_mactaggart is an author in the scientific field of machine learning	1111	06-jul-2018	96.3	
_199_00 is an award, championship, or tournament trophy	1111	06-jul-2018	100.0	

Never Ending Language Learning (NELL)

NELL Knowledge Base Browser

CMU Read the Web Project

log in | preferences | help/instruction

- room
- officebuildingroom
- website
 - blog
 - politicsblog
 - url
- river
- attraction
 - transportation
 - museum
 - port
 - monument
 - skiarea
 - aquarium
 - stadiumorevenue
 - park
 - zoo
- highway
- trail
- planet
- street
- beach
- mountainrange
- lake
- zipcode
- caf_
- cave
- geolocatablething
 - building (+)
 - mountain
 - river
 - stateorprovince
 - country
 - oilgasfield
 - radiostation
 - street

transportation

(category)

a mode of travel, such as rail transport or automobile. It is a category means of moving from one location

View list | [map](#) | [metadata](#)
6,114 instances, 2 pages: 1 2 next

instance	iteration	date learned
aalborg_airport	1098	22-jan-2018
aberdeen_airport	1110	24-jun-2018
aberdeen_dyce_airport	1103	18-mar-2018
abilene_airport	1110	24-jun-2018
abu_dhabi_airport	1114	25-aug-2018
acadiana_regional_airport	1103	18-mar-2018
adelaide_airport	1115	03-sep-2018
adelaide_international_airport	1103	18-mar-2018
agadir_airport	1105	31-mar-2018
agen_airport	1096	19-jan-2018
agra_airport	1108	11-jun-2018
aguadilla_airport	1112	24-jul-2018
ahmedabad_airport	1098	22-jan-2018
airlake_airport	1116	12-sep-2018
ajaccio_airport	1103	18-mar-2018
akureyri_airport	1095	18-jan-2018
albany_international_airport	1113	15-aug-2018
albuquerque_airport	1110	24-jun-2018
albuquerque_international_airport	1110	24-jun-2018
albury_airport	1100	02-feb-2018
alexandria_airport	1112	24-jul-2018
alghero_airport	1103	18-mar-2018
alicante_airport	1103	18-mar-2018

Never Ending Language Learning (NELL)

NELL Knowledge Base Browser

CMU Read the Web Project

log in | preferences | help/instruction

categories

relations











- everypromotedthing
- location
 - building
 - airport
 - bridge
 - hotel
 - placeofworship
 - retailstore
 - museum
 - monument
 - restaurant
 - stadiumoreventvenue
 - shoppingmall
 - skyscraper
 - hospital
 - trainstation
 - geopoliticallocation
 - county
 - continent
 - stateorprovince
 - country
 - city
 - island
 - mountain
 - farm
 - landscapefeatures
 - room
 - officebuildingroom
 - website
 - blog
 - politicsblog
 - url
 - river
 - attraction
 - transportation

aberdeen_airport (airport)

literal strings: [aberdeen_airport](#), [aberdeen_airport](#), [Aberdeen Airport](#), [Aberdeen airport](#)

Help NELL Learn!

NELL wants to know if these beliefs are correct.
If they are or ever were, click thumbs-up. Otherwise, click thumbs-down.

- [aberdeen_airport](#) is a [building](#)  
- [aberdeen_airport](#) is an [airport](#)  
- [aberdeen_airport](#) is a [transportation system](#)  
- [aberdeen_airport](#) is a [tourist attraction](#)  
- [aberdeen_airport](#) is an airport [in the city aberdeen](#) (city)  

categories

- [building](#)(100.0%)
 - OE @827 (100.0%) on 30-mar-2014 [] using aberdeen_airport
 - CPL @1096 (97.4%) on 18-jan-2018 ["passenger numbers at _ "stay car park at _ "book flights from _ "you are flying "airport transfers to _ " is here Directions Maps Public transport." "Airport car parking at _ "car parking facilities near parking at _ "Secure parking at _ "Long Stay car park at _ "discounted car hire at _ "Locate flights from _ " _ providing kitchen" "Approx distances to _ "airport locations at _ "D-minute drive from _ "on-airport parking at _ "cheap car hire "airport car parking facilities near _ "airport car parking facilities at _ "United Kingdom Book hotels serving _ "Book o at _ "Cheap car hire at _ "parking options available at _ "parking option at _ "Kingdom Book hotels serving _ "you 'r from _ "airport car parking at _ " _ has a range duty free shops" "return flights from _ "parking providers at _ "parking at _ "Stay parking at _ "Book hotels serving _ " include FlyBe" "airport car parks at _ "airport car parks near _ "main at _ "airport lounges at _ "Terminal building at _ "cheap hotel near _ "rental car companies at _ "off-site parking at "available rental car companies at _ "available car parks at _ "Servisair lounge at _ "] using aberdeen_airport
 - SEAL @180 (100.0%) on 16-dec-2010 [1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 32 33 34 35 36 37 38 39 40 41 42 43 44 45] using aberdeen_airport
 - CMC @266 (99.9%) on 23-may-2011 [PREFIX=ai LASTSUFFIX=rport LAST_WORD=airport LASTPREFIX=air LASTPREFIX=LASTPREFIX=ai] using aberdeen_airport

Never Ending Language Learning (NELL)

So far, NELL has accumulated over 50 million candidate beliefs by reading the web, and it is considering these at different levels of confidence. NELL has high confidence in 2,810,379 of these beliefs — these are displayed on this website. It is not perfect, but NELL is learning. You can track NELL's progress below or @cmunell on Twitter, browse and download its knowledge base, read more about our technical approach, or join the discussion group.

- basé sur les types de Yago
- aide au développement d'extracteurs

The screenshot displays a web interface for exploring text corpora, organized into four main sections:

- X-Type**: A search interface for X-types. The search bar contains "Search for x-types..." and a "Search" button. The "Active:" section shows "computer.software" selected. A list of other X-types includes "book.book_subject", "cvg.cvg_platform", "computer.operating_system", "computer.web_browser", and "network.node".
- Patterns**: A search interface for patterns. The search bar contains "Search for patterns..." and a "Search" button. The "Active:" section shows several patterns: "[X] support.v [Y]", "open [Y] file in [X]", "[X] play [Y] file", and "[X] read [Y] file". Other patterns include "[Y] file in [X]" and "[X] use [Y]".
- Y-Type**: A search interface for Y-types. The search bar contains "Search for y-types..." and a "Search" button. The "Active:" section shows "computer.file_format" selected. A list of other Y-types includes "computer.internet_protocol", "computer.software_genre", "book.book_binding", "internet.protocol", and "network.node".
- Output**: A table displaying extracted information. The table has three columns: "Subject", "Object", and "Example".

Subject	Object	Example
Excel	SpreadsheetML	Excel supports SpreadsheetML for both import and export, providing a complete pathway for information.
Dia	SVG	My diagram editor of choice is Dia , which supports export-to-SVG, so one approach is attaching the .
Google Earth	KML	I opened the KML file of this jog in Google Earth, and so now I've associated ...

At the bottom left of the Output section, there is a small icon resembling a link or refresh symbol.

- extraction of genomic knowledge from [PubMed](#) articles
- available in the cloud for browsing, searching and reasoning

The Literome Project Welcome 100.68.140.29 Microsoft Research

Search for directed genic interactions:

ABL1 → **CTNNB1** (1 - 4 of 4)

Direct Search

ABL1 → **CTNNB1** (4)

CTNNB1 → **ABL1** (1)

Possible intermediates for ABL1 → **CTNNB1**

- ABL1 → ACD → CTNNB1
- ABL1 → BCL2 → CTNNB1
- ABL1 → BCR → CTNNB1
- ABL1 → BMP2 → CTNNB1
- ABL1 → CCND1 → CTNNB1
- ABL1 → CD4 → CTNNB1
- ABL1 → CD79A → CTNNB1
- ABL1 → CD79B → CTNNB1
- ABL1 → CDK5 → CTNNB1
- ABL1 → CEBPA → CTNNB1
- ABL1 → CISH → CTNNB1
- ABL1 → CRK → CTNNB1
- ABL1 → CSF1R → CTNNB1
- ABL1 → CSF3 → CTNNB1
- ABL1 → CTNNB1 → CTNNB1
- ABL1 → CTNND1 → CTNNB1
- ABL1 → CXCL12 → CTNNB1

PMID: 17318191
Bcr-Abl stabilizes beta-catenin in chronic myeloid leukemia through its tyrosine phosphorylation.

PMID: 17618275
Cables links Robo-bound Abl kinase to N-cadherin-bound beta-catenin to mediate Slit-induced modulation of adhesion and transcription.

Bcr-Abl stabilizes **beta-catenin** in ... (details)

Bcr-Abl physically ... is **required** to phosphorylate **beta-catenin** at ... (details)

... the **Bcr-Abl** triggered **Y** ... of **beta-catenin** as ... (details)

... in **Abl** mediated phosphorylation of **beta-catenin** on ... (details)

About Contact Us Terms of Use Trademarks Privacy Statement © 2015 Microsoft Corporation. All rights reserved.

- QA \equiv séquence d'opérateurs qui transforment une question en une requête à un *tuple-store*
 - 4 \neq tuple-stores : FREEBASE et 3 autres extraits automatiquement
- une fonction donnant un score à chaque opération est apprise sur un corpus de questions et leur réponse.
- répondre \equiv chercher dans l'espace des opérations

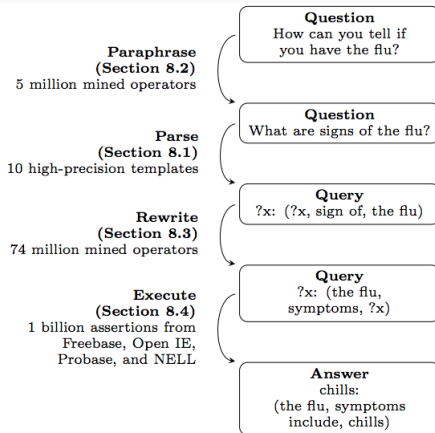


Figure 1: OQA automatically mines millions of operators (left) from unlabeled data, then learns to compose them to answer questions (right) using evidence from multiple knowledge bases.

Open-domain QA [Fader et al., 2014]

Input	What are some examples of building maintenance jobs?
Parse	?x: (?x, example of, building maintenance jobs)
Rewrite	?x: (?x, is-a, building maintenance job)
Execute	{changing light bulb, is-a, small building maintenance job}

Input	What animal represents California?
Paraphrase	What are California's symbols?
Parse	?x: (california, symbols, ?x)
Execute	{California Water Service, Trading symbol, CWT}

Conclusion

Conclusion

- Extraire et exploiter des connaissances extraites de textes est un **enjeu** académique et industriel

Conclusion

- Extraire et exploiter des connaissances extraites de textes est un **enjeu** académique et industriel
- Nombreuses ressources disponibles :
 - [Yago](#) [Freebase](#), [ConcepNet](#), etc.
 - s'y lier est important (*entity linking*)

Conclusion

- Extraire et exploiter des connaissances extraites de textes est un **enjeu** académique et industriel
- Nombreuses ressources disponibles :
 - [Yago](#) [Freebase](#), [ConcepNet](#), etc.
 - s'y lier est important (*entity linking*)
- Évaluer la technologie d'extraction est difficile
 - arbitraire, **rappel** difficile à mesurer

Conclusion

- Extraire et exploiter des connaissances extraites de textes est un **enjeu** académique et industriel
- Nombreuses ressources disponibles :
 - [Yago](#) [Freebase](#), [ConcepNet](#), etc.
 - s'y lier est important (*entity linking*)
- Évaluer la technologie d'extraction est difficile
 - arbitraire, **rappel** difficile à mesurer
- pas ou peu de normalisation des relations extraites

Conclusion

- Extraire et exploiter des connaissances extraites de textes est un **enjeu** académique et industriel
- Nombreuses ressources disponibles :
 - [Yago](#) [Freebase](#), [ConcepNet](#), etc.
 - s'y lier est important (*entity linking*)
- Évaluer la technologie d'extraction est difficile
 - arbitraire, **rappel** difficile à mesurer
- pas ou peu de normalisation des relations extraites
- Beaucoup d'information textuelle difficile d'accès
 - importance des informations temporelles

Chilly Gonzales

"Chilly Gonzales (born Jason Charles Beck; 20 March 1972) is a Grammy-winning Canadian musician who resided in Paris, France for several years, and now lives in Cologne, Germany."

Chilly Gonzales

"Chilly Gonzales (born Jason Charles Beck ; 20 March 1972) is a Grammy-winning Canadian musician who resided in Paris, France for several years, and now lives in Cologne, Germany."

⟨Chilly Gonzales, [was] born, Jason Charles Beck⟩

⟨Chilly Gonzales, [was] born [on], 20 March 1972⟩

⟨Chilly Gonzales, is a, musician⟩

⟨Chilly Gonzales, [won], Grammy [award]⟩

⟨Chilly Gonzales, is, Canadian⟩

⟨Chilly Gonzales, resided in, Paris⟩ **for several years**

⟨Paris, [is in], France⟩

⟨Chilly Gonzales, lives in, Cologne⟩

⟨Cologne, [is in], Germany⟩



Agichtein, E. and Gravano, L. (2000).

Snowball : Extracting relations from large plain-text collections.

In *Fifth ACM Conference on Digital Libraries*, pages 85–94.



Akbik, A., Michael, T., and Boden, C. (2014).

Exploratory relation extraction in large text corpora.

In *25th International Conference on Computational Linguistics*, pages 2087–2096.



Angeli, G., Johnson Premkumar, M. J., and Manning, C. D. (2015).

Leveraging linguistic structure for open domain information extraction.

In *53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on*

Natural Language Processing (Volume 1 : Long Papers), pages 344–354.



Balasubramanian, N., Soderland, S., Mausam, and Etzioni, O. (2013).

Generating coherent event schemas at scale.

In *EMNLP'13*, pages 1721–1731.



Banko, M., Cafarella, M. J., Soderland, S., Broadhead, M., and Etzioni, O. (2007).

Open information extraction from the web.

In *IN IJCAI*, pages 2670–2676.



Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. (2018).

Think you have solved question answering? try arc, the AI2 reasoning challenge.

CoRR.



Fader, A., Soderland, S., and Etzioni, O. (2011).

Identifying relations for open information extraction.

In *Empirical Methods in Natural Language Processing*, EMNLP '11, pages 1535–1545.



Fader, A., Zettlemoyer, L., and Etzioni, O. (2014).

Open question answering over curated and extracted knowledge bases.

In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '14, pages 1156–1165.



Forand, K. and Langlais, P. (2017).

Curating an open information extraction knowledge base using games with a purpose.

In *Game4NLP, workshop at EACL 2017*, Valencia, Spain.



Gotti, F. and Langlais, P. (2016).

From french wikipedia to erudit : A test case for cross-domain open information extraction.

Computational Intelligence.



GuoDong, Z., Jian, S., Jie, Z., and Min, Z. (2005).

Exploring various knowledge in relation extraction.

In 43rd Annual Meeting on Association for Computational Linguistics, pages 427–434.



Hearst, M. A. (1992).

Automatic acquisition of hyponyms from large text corpora.

In Proceedings of the 14th Conference on Computational Linguistics - Volume 2, COLING '92, pages 539–545.



Lai, G., Xie, Q., Liu, H., Yang, Y., and Hovy, E. (2017).

Race : Large-scale reading comprehension dataset from examinations.

In *Conference on Empirical Methods in Natural Language Processing*, pages 785–794. Association for Computational Linguistics.



Mausam, Schmitz, M., Bart, R., Soderland, S., and Etzioni, O. (2012).

Open language learning for information extraction.

In *Joint EMNLP and CoNLL*, pages 523–534.



Min, B., Shi, S., Grishman, R., and Lin, C.-Y. (2012).

Ensemble semantics for large-scale unsupervised relation extraction.

In *Proceedings of the 2012 Joint EMNLP and CoNLL*, pages 1027–1037.



Nakashole, N., Weikum, G., and Suchanek, F. (2012).

Patty : A taxonomy of relational patterns with semantic types.

In *Joint EMNLP and CoNLL*, pages 1135–1145.



Niklaus, C., Cetto, M., Freitas, A., and Handschuh, S. (2018).
A survey on open information extraction.

In *27th International Conference on Computational Linguistics*, pages 3866–3878.



Poon, H., Quirk, C., DeZiel, C., and Heckerman, D. (2014).
Literome : Pubmed-scale genomic knowledge base in the cloud.

Bioinformatics, 30(19) :2840–2842.



Wu, F. and Weld, D. S. (2010).

Open information extraction using wikipedia.

In *48th Annual Meeting of the Association for Computational Linguistics*, pages 118–127.