

# Problem 2

For this problem set, we will use

<https://app.sketchengine.eu/#dashboard?corpname=preloaded%2F covid19>

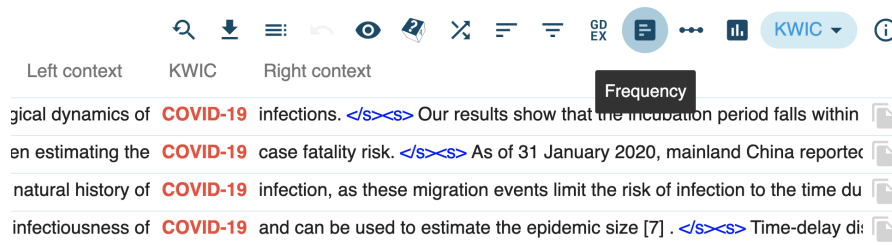
## Problem 2.1

Click the above link, and follow this: Dashboard -> Concordance -> Advanced -> CQL.

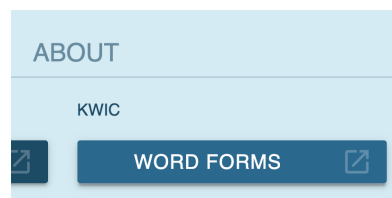
Now write a query to find sentences containing all forms of covid and execute it. Some forms include covid-19, covid19, COVID19, covid-36, covid-54.

Once you get the sentences, click `Frequency -> KWIC > WORD FORMS` to generate the frequency of words. These steps are shown below:

Step 1:



Step 2:



Step 3: The word list looks something like this:

	Word	↓ Frequency	Per million tokens
1	<input type="checkbox"/> COVID-19	20,773	73.99
2	<input type="checkbox"/> Covid-19	429	1.53
3	<input type="checkbox"/> COVID19	169	0.60
4	<input type="checkbox"/> COVID-2019	157	0.56
5	<input type="checkbox"/> CoVID-19	32	0.11

**What is the CQL query that you used for getting all forms of covid (i.e. the query that is used to generate the above figure)?**

Answer: [word = "[cC][oO][vV][iI][dD]-?\d+"]

**Include the snapshot of the top 20 words (5 words are shown above)?**

Answer:

(31 items, 21,642 total frequency)

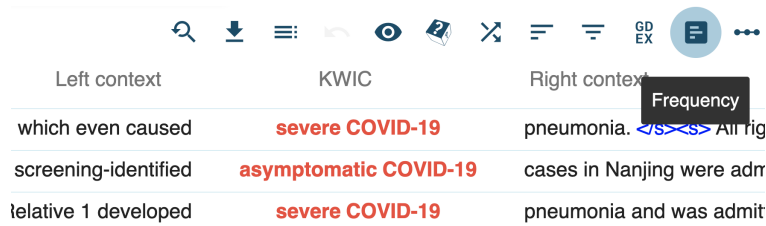
		Word	Frequency ↓	Relative ?		
1	<input type="checkbox"/>	COVID-19	20,773	73.99	<div></div>	...
2	<input type="checkbox"/>	Covid-19	429	1.53	<div></div>	...
3	<input type="checkbox"/>	COVID19	169	0.60	<div></div>	...
4	<input type="checkbox"/>	COVID-2019	157	0.56	<div></div>	...
5	<input type="checkbox"/>	CoVID-19	32	0.11	<div></div>	...
6	<input type="checkbox"/>	covid-19	30	0.11	<div></div>	...
7	<input type="checkbox"/>	CoViD-19	10	0.04	<div></div>	...
8	<input type="checkbox"/>	COVID-10	7	0.02	<div></div>	...
9	<input type="checkbox"/>	COVID-9	7	0.02	<div></div>	...
10	<input type="checkbox"/>	Covid-2019	4	0.01	<div></div>	...
11	<input type="checkbox"/>	covid19	3	0.01	<div></div>	...
12	<input type="checkbox"/>	Covid19	2	< 0.01	<div></div>	...
13	<input type="checkbox"/>	covid-10	1	< 0.01	<div></div>	...
14	<input type="checkbox"/>	COVID-138	1	< 0.01	<div></div>	...
15	<input type="checkbox"/>	Covid-10	1	< 0.01	<div></div>	...
16	<input type="checkbox"/>	Covid-56	1	< 0.01	<div></div>	...
17	<input type="checkbox"/>	COVID-173	1	< 0.01	<div></div>	...
18	<input type="checkbox"/>	COVID-27	1	< 0.01	<div></div>	...
19	<input type="checkbox"/>	COVID-110	1	< 0.01	<div></div>	...
20	<input type="checkbox"/>	COVID-2	1	< 0.01	<div></div>	...

## Problem 2.2

Let's write CQL queries to find interesting words that occur in specific syntactic relations with covid (all forms). We did similar things in class. You will have to use tag and lemma in CQL queries. This [tagset](#) could be useful

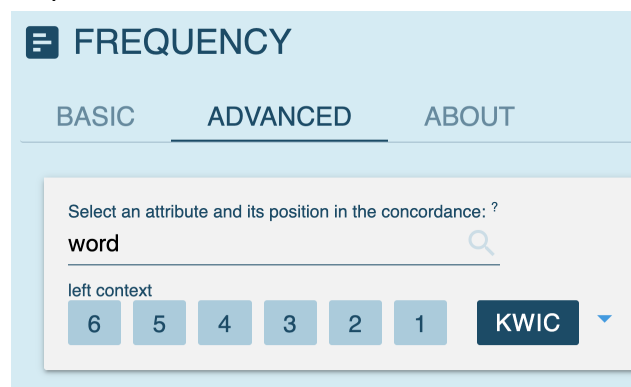
I will demonstrate how to get the modifiers of covid:

Step 1: First write a CQL query that produces concordance (examples) like this:



Left context	KWIC	Right context
which even caused	severe COVID-19	pneumonia. </s></s> All rig
screening-identified	asymptomatic COVID-19	cases in Nanjing were adr
relative 1 developed	severe COVID-19	pneumonia and was admit

Step 2:



**FREQUENCY**

BASIC **ADVANCED** ABOUT

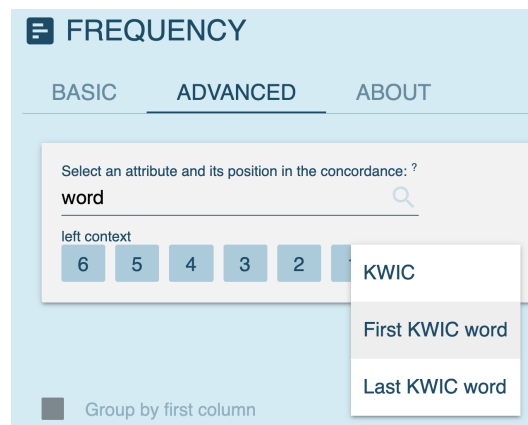
Select an attribute and its position in the concordance: ?

word

left context

6 5 4 3 2 1 KWIC

Step 3:



**FREQUENCY**

BASIC **ADVANCED** ABOUT

Select an attribute and its position in the concordance: ?

word

left context

6 5 4 3 2 1 KWIC

First KWIC word

Last KWIC word

Group by first column

Step 4:

		Word	↓ Frequency	Per million tokens
1	<input type="checkbox"/>	severe	298	1.06
2	<input type="checkbox"/>	confirmed	115	0.41
3	<input type="checkbox"/>	current	103	0.37

**What is the CQL query for modifiers of covid (all forms)?**

Answer: [tag = "J.\*"] [word = "[cC][oO][vV][iI][dD]-?\d+"]

**Include the snapshot of modifiers (top three are shown above)**

(307 items, 1,712 total frequency)

		Word	Frequency ↓	Relative ?		
1	<input type="checkbox"/>	severe	298	1.06	<div><div></div></div>	...
2	<input type="checkbox"/>	confirmed	115	0.41	<div><div></div></div>	...
3	<input type="checkbox"/>	current	103	0.37	<div><div></div></div>	...
4	<input type="checkbox"/>	suspected	81	0.29	<div><div></div></div>	...
5	<input type="checkbox"/>	laboratory-confirmed	69	0.25	<div><div></div></div>	...
6	<input type="checkbox"/>	ongoing	64	0.23	<div><div></div></div>	...
7	<input type="checkbox"/>	new	43	0.15	<div><div></div></div>	...
8	<input type="checkbox"/>	first	42	0.15	<div><div></div></div>	...
9	<input type="checkbox"/>	mild	40	0.14	<div><div></div></div>	...
10	<input type="checkbox"/>	reported	31	0.11	<div><div></div></div>	...
11	<input type="checkbox"/>	critical	27	0.10	<div><div></div></div>	...
12	<input type="checkbox"/>	potential	26	0.09	<div><div></div></div>	...
13	<input type="checkbox"/>	global	24	0.09	<div><div></div></div>	...
14	<input type="checkbox"/>	ill	22	0.08	<div><div></div></div>	...
15	<input type="checkbox"/>	asymptomatic	18	0.06	<div><div></div></div>	...
16	<input type="checkbox"/>	early	17	0.06	<div><div></div></div>	...
17	<input type="checkbox"/>	moderate	16	0.06	<div><div></div></div>	...
18	<input type="checkbox"/>	active	15	0.05	<div><div></div></div>	...
19	<input type="checkbox"/>	novel	14	0.05	<div><div></div></div>	...
20	<input type="checkbox"/>	recent	14	0.05	<div><div></div></div>	...

**What is the CQL query of words that are modified by covid (all forms)?**

Answer: [word = "[cC][oO][vV][iI][dD]-?\d+"] [tag = "N.\*"]

**Include the snapshot of those words**

(549 items, 8,856 total frequency)

	Word	Frequency ↓	Relative ?		
1	<input type="checkbox"/> patients	1,720	6.13	<div><div></div></div>	...
2	<input type="checkbox"/> cases	954	3.40	<div><div></div></div>	...
3	<input type="checkbox"/> outbreak	721	2.57	<div><div></div></div>	...
4	<input type="checkbox"/> infection	696	2.48	<div><div></div></div>	...
5	<input type="checkbox"/> epidemic	540	1.92	<div><div></div></div>	...
6	<input type="checkbox"/> pneumonia	496	1.77	<div><div></div></div>	...
7	<input type="checkbox"/> pandemic	409	1.46	<div><div></div></div>	...
8	<input type="checkbox"/> resource	396	1.41	<div><div></div></div>	...
9	<input type="checkbox"/> virus	153	0.54	<div><div></div></div>	...
10	<input type="checkbox"/> case	147	0.52	<div><div></div></div>	...
11	<input type="checkbox"/> infections	144	0.51	<div><div></div></div>	...
12	<input type="checkbox"/> transmission	141	0.50	<div><div></div></div>	...
13	<input type="checkbox"/> disease	125	0.45	<div><div></div></div>	...
14	<input type="checkbox"/> patient	104	0.37	<div><div></div></div>	...
15	<input type="checkbox"/> spread	63	0.22	<div><div></div></div>	...
16	<input type="checkbox"/> diagnosis	58	0.21	<div><div></div></div>	...
17	<input type="checkbox"/> outbreaks	55	0.20	<div><div></div></div>	...
18	<input type="checkbox"/> testing	55	0.20	<div><div></div></div>	...
19	<input type="checkbox"/> treatment	51	0.18	<div><div></div></div>	...
20	<input type="checkbox"/> mortality	50	0.18	<div><div></div></div>	...

**What is the CQL query for words that occur in right coordination with covid (all forms)**  
(e.g., in COVID-19 , SARS-2002 , and HCoV-NL63, the words iSARS-2002 and HCoV-NL63 are the right conjuncts/coordinates).

Answer:

```
[word = "[cC][oO][vV][iI][dD]-?\d+" & tag = "N.*"] [word = "\", " | tag = "N.*"] {0,} [tag = "CC"] [tag = "N.*"]
```

**Include the snapshot of those words**

(210 items, 426 total frequency)

	Word	Frequency ↓	Relative ?	
1	<input type="checkbox"/> SARS	31	0.11	...
2	<input type="checkbox"/> MERS-COV	14	0.05	...
3	<input type="checkbox"/> MERS	13	0.05	...
4	<input type="checkbox"/> Treatment	11	0.04	...
5	<input type="checkbox"/> SARS-CoV-2	11	0.04	...
6	<input type="checkbox"/> control	10	0.04	...
7	<input type="checkbox"/> prevention	9	0.03	...
8	<input type="checkbox"/> H1N1	9	0.03	...
9	<input type="checkbox"/> treatment	8	0.03	...
10	<input type="checkbox"/> HAPE	6	0.02	...
11	<input type="checkbox"/> death	6	0.02	...
12	<input type="checkbox"/> deaths	6	0.02	...
13	<input type="checkbox"/> influenza	5	0.02	...
14	<input type="checkbox"/> SARS-CoV	5	0.02	...
15	<input type="checkbox"/> mortality	5	0.02	...
16	<input type="checkbox"/> RoS	4	0.01	...
17	<input type="checkbox"/> cancer	4	0.01	...
18	<input type="checkbox"/> SARS-2002	4	0.01	...
19	<input type="checkbox"/> diagnosis	4	0.01	...
20	<input type="checkbox"/> B	4	0.01	...

**What is the CQL query for verbs that can take covid (all forms) as subject?**

Answer: [word = "[cC][oO][vV][iI][dD]-?d+" & tag = "N.\*"][]{}[0,3][tag = "VV.\*"]

**Include the snapshot of verbs that take covid as subject**

(1,436 items, 8,273 total frequency)

	Word	Frequency ↓	Relative ?	
1	<input type="checkbox"/> remains	444	1.58	...
2	<input type="checkbox"/> reported	228	0.81	...
3	<input type="checkbox"/> confirmed	191	0.68	...
4	<input type="checkbox"/> caused	179	0.64	...
5	<input type="checkbox"/> spread	138	0.49	...
6	<input type="checkbox"/> using	133	0.47	...
7	<input type="checkbox"/> based	113	0.40	...
8	<input type="checkbox"/> including	105	0.37	...
9	<input type="checkbox"/> identified	98	0.35	...
10	<input type="checkbox"/> found	90	0.32	...
11	<input type="checkbox"/> admitted	74	0.26	...
12	<input type="checkbox"/> showed	72	0.26	...
13	<input type="checkbox"/> compared	69	0.25	...
14	<input type="checkbox"/> spreading	69	0.25	...
15	<input type="checkbox"/> included	66	0.24	...
16	<input type="checkbox"/> include	64	0.23	...
17	<input type="checkbox"/> associated	60	0.21	...
18	<input type="checkbox"/> occurred	56	0.20	...
19	<input type="checkbox"/> according	56	0.20	...
20	<input type="checkbox"/> become	51	0.18	...

**What is the CQL query for verbs that can take covid (all forms) as object?**

Answer: [tag = "VV.\*"][]{0,5}[word = "[cC][oO][vV][iI][dD]-?\d+" & tag = "N.\*"]

**Include the snapshot of verbs that take COVID as object.**

(1,637 items, 11,070 total frequency)

		Word	Frequency ↓	Relative ?		
1	<input type="checkbox"/>	confirmed	663	2.36	<div><div></div></div>	...
2	<input type="checkbox"/>	reported	252	0.90	<div><div></div></div>	...
3	<input type="checkbox"/>	diagnosed	242	0.86	<div><div></div></div>	...
4	<input type="checkbox"/>	infected	211	0.75	<div><div></div></div>	...
5	<input type="checkbox"/>	used	130	0.46	<div><div></div></div>	...
6	<input type="checkbox"/>	suspected	128	0.46	<div><div></div></div>	...
7	<input type="checkbox"/>	associated	127	0.45	<div><div></div></div>	...
8	<input type="checkbox"/>	related	107	0.38	<div><div></div></div>	...
9	<input type="checkbox"/>	caused	103	0.37	<div><div></div></div>	...
10	<input type="checkbox"/>	control	99	0.35	<div><div></div></div>	...
11	<input type="checkbox"/>	named	95	0.34	<div><div></div></div>	...
12	<input type="checkbox"/>	hospitalized	94	0.33	<div><div></div></div>	...
13	<input type="checkbox"/>	found	90	0.32	<div><div></div></div>	...
14	<input type="checkbox"/>	including	87	0.31	<div><div></div></div>	...
15	<input type="checkbox"/>	declared	85	0.30	<div><div></div></div>	...
16	<input type="checkbox"/>	prevent	83	0.30	<div><div></div></div>	...
17	<input type="checkbox"/>	treat	83	0.30	<div><div></div></div>	...
18	<input type="checkbox"/>	treating	78	0.28	<div><div></div></div>	...
19	<input type="checkbox"/>	regarding	77	0.27	<div><div></div></div>	...
20	<input type="checkbox"/>	contain	73	0.26	<div><div></div></div>	...

## Problem 2.3

What are the most important words that form collocations with COVID (where covid is the right word)?

You can generate collocations as follows: First get concordance of all forms of covid.

Step 1:

Left context      KWIC      Right

epidemiological dynamics of **COVID-19** infect  
red when estimating the **COVID-19** case  
adv the natural history of **COVID-19** infect

Step 2:

Left context      KWIC      Right context      Collocations

ical dynamics of **COVID-19** infections. </s><s> Our results show that the

Step 3:

COLLOCATIONS

BASIC      ADVANCED      ABOUT

Attribute ?      Range ?

word      -5 -4 -3 -2 -1 KWIC 1 2 3 4 5

Step 4:

	Word	Cooccurrences ?
1	<input type="checkbox"/> confirmed	458
2	<input type="checkbox"/> suspected	133

Show the collocations sorted according to what you think is the best metric (T-Score, MI, LogDice). Indicate the metric you used.

	Word	Cooccurrences ?	Candidates ?	T-score	MI	LogDice ↓
1	<input type="checkbox"/> confirmed	458	65,495	21.17	6.50	7.43 ...
2	<input type="checkbox"/> suspected	133	21,439	11.39	6.33	6.66 ...
3	<input type="checkbox"/> laboratory-confirmed	75	3,601	8.63	8.08	6.61 ...
4	<input type="checkbox"/> severe	298	112,078	16.76	5.11	6.19 ...
5	<input type="checkbox"/> ongoing	64	12,451	7.88	6.06	5.94 ...
6	<input type="checkbox"/> treat	57	14,546	7.40	5.67	5.69 ...
7	<input type="checkbox"/> treating	46	9,478	6.67	5.98	5.60 ...
8	<input type="checkbox"/> current	103	50,596	9.76	4.72	5.55 ...
9	<input type="checkbox"/> declared	33	4,219	5.69	6.66	5.39 ...
10	<input type="checkbox"/> towards	44	17,999	6.42	4.99	5.18 ...

	Word	Cooccurrences ?	Candidates ?	T-score	MI	LogDice ↓
11	<input type="checkbox"/> hospitalized	40	15,225	6.14	5.09	5.15 ...
12	<input type="checkbox"/> named	33	9,921	5.61	5.43	5.10 ...
13	<input type="checkbox"/> about	114	98,409	9.97	3.91	4.96 ...
14	<input type="checkbox"/> non-severe	21	825	4.57	8.37	4.94 ...
15	<input type="checkbox"/> with	1,999	2,412,053	40.55	3.43	4.75 ...
16	<input type="checkbox"/> mild	40	30,282	5.96	4.10	4.66 ...
17	<input type="checkbox"/> against	151	180,158	11.16	3.44	4.62 ...
18	<input type="checkbox"/> contracted	17	1,822	4.09	6.92	4.57 ...
19	<input type="checkbox"/> having	34	27,817	5.46	3.99	4.49 ...
20	<input type="checkbox"/> of	5,930	8,766,274	68.23	3.13	4.47 ...



I used LogDice.

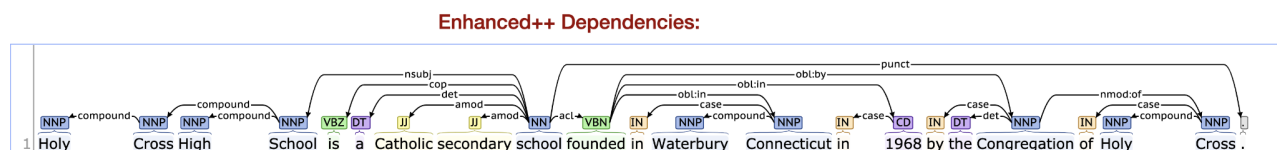
## Problem 3:

Write [SemGreX](#) regular expressions that can detect organizations and their founders. Make use of <https://corenlp.run> to parse sentences to syntactic graphs and for running SemGreX expressions.

Here is an example:

**Holy Cross High School** is a Catholic secondary school founded in Waterbury Connecticut in 1968 by the **Congregation of Holy Cross** .

The corresponding Enhanced++ Dependencies syntactic graph is as follows:



The below SemGreX pattern extracts the headword of the organization and the headword of the founder.

```
{}=organization <nsubj ({} >cop {} >acl ({}lemma:found} >/obl:by/ {}=founder))
```

**CoreNLP Tools:**

TokensRegex

Semgrex

Tregex

Enter a **Semgrex** expression to run against the "enhanced dependencies" above:

{}=organization <nsubj ({} >cop {} >acl ({}lemma:found} >/obl:by/ {}=founder))

Match

match

organization

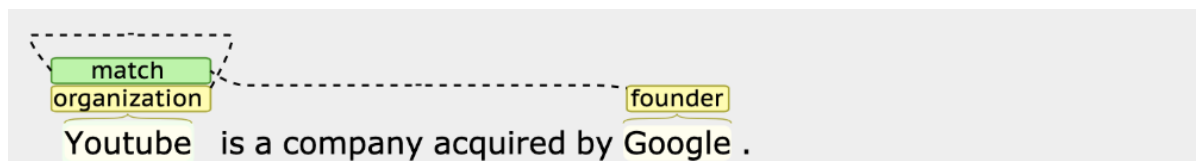
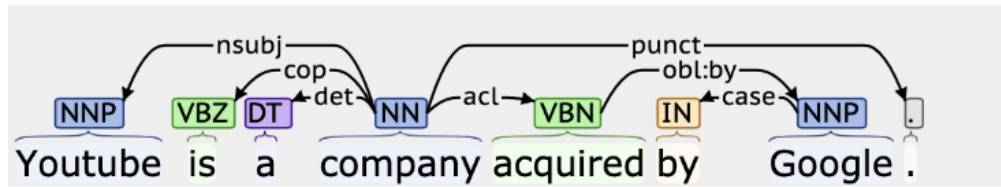
1 Holy Cross High School is a Catholic secondary school founded in Waterbury Connecticut in 1968 by the Congregation of Holy Cross .

founder

This pattern can be read as the “organization” that is a subject of something, and this something is founded by the founder.

Here it extracts School (i.e., the headword of Holy Cross High School) as the organization and Congregation (i.e., the headword of the Congregation of Holy Cross) as the founder.

Your goal is to write SemGreX expressions that can generalize to multiple sentences but at the same time don't match incorrect sentences. For example, if you don't use {lemma:found} in the above sentence, your pattern will also match a sentence like "Youtube is a company acquired by Google" (see below.)

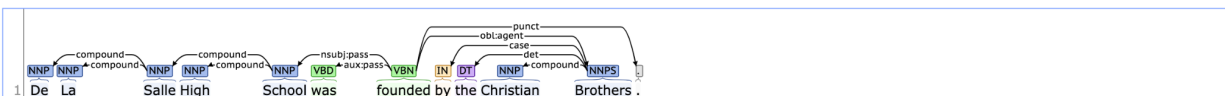


### Problem 3.1

Write the SemGreX patterns for the following sentences that extract the organization name (headword is enough) and its founder (headword is enough). Sentences that can make use of the same expression should be in the same snapshot (containing Enhanced++ Dependencies, Semgrex expression, and the matchings):

**De La Salle High School** was founded by **the Christian Brothers** .

Enhanced++ Dependencies:



CoreNLP Tools:

TokensRegex Semgrex Tregex

Enter a Semgrex expression to run against the "enhanced dependencies" above:

{pos:/NNP\*/}=organization </nsubj:pass/ ({lemma:found} >/obl:agent/ {pos:/NNP\*/}=founder

Match

1 De La Salle High School was founded by the Christian Brothers .

match organization founder

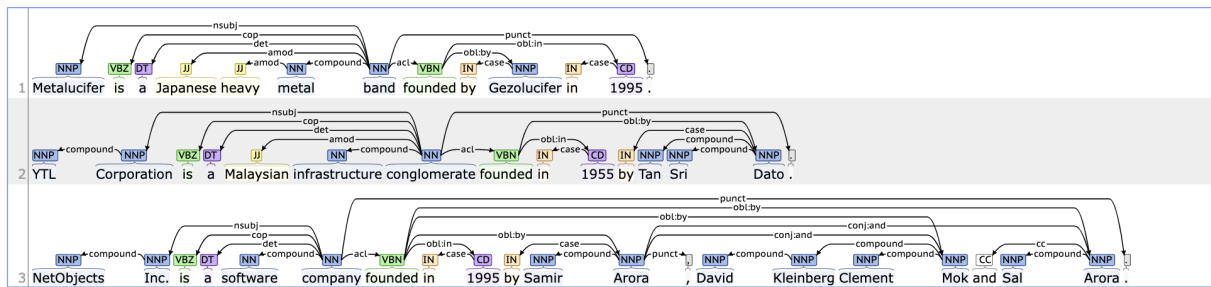
**Metalucifer** is a Japanese heavy metal band founded by **Gezolucifer** in 1995 .

**YTL Corporation** is a Malaysian infrastructure conglomerate founded in 1955 by **Tan Sri Dato**.

**NetObjects Inc.** is a software company founded in 1995 by **Samir Arora, David Kleinberg Clement Mok and Sal Arora** .

(If there are multiple founders, you have to extract headword corresponding to each founder)

#### Enhanced++ Dependencies:



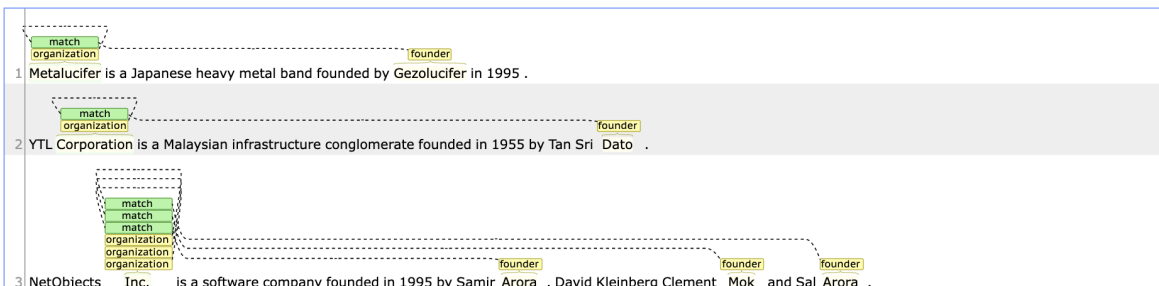
#### CoreNLP Tools:

TokensRegex Semgrep Tregex

Enter a Semgrep expression to run against the "enhanced dependencies" above:

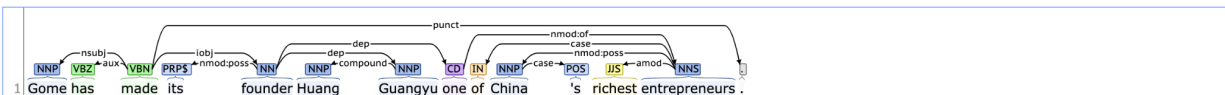
{pos:/NNP\*/}=organization <nsbj {pos:/NNS\*/} >acl {lemma:found} >/obl:by/ {pos:/NNP\*/}=founder)

Match



Gome has made its founder Huang Guangyu one of China's richest entrepreneurs.

#### Enhanced++ Dependencies:



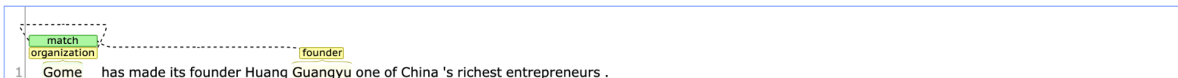
#### CoreNLP Tools:

TokensRegex Semgrep Tregex

Enter a Semgrep expression to run against the "enhanced dependencies" above:

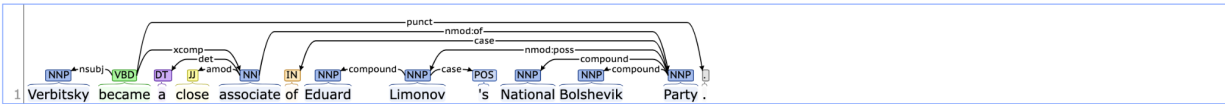
{pos:/NNP\*/}=organization <nsbj {lemma:make} >obj {lemma:founder} >dep {pos:/NNP\*/}=founder)

Match



Verbitsky became a close associate of Eduard Limonov's National Bolshevik Party.

### Enhanced++ Dependencies:



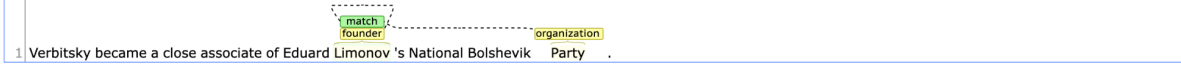
### CoreNLP Tools:

TokensRegex Semgrep Tregex

Enter a Semgrep expression to run against the "enhanced dependencies" above:

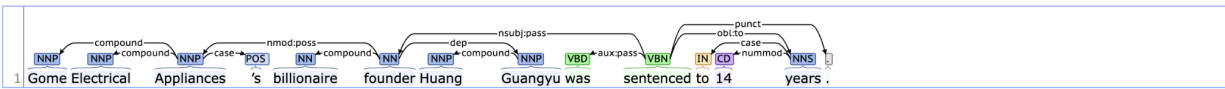
{pos:/NNP.\*}=founder </nmod:poss/ {pos:/NNP.\*}=organization

Match



Gome Electrical Appliances's billionaire founder Huang Guangyu was sentenced to 14 years.

### Enhanced++ Dependencies:



### CoreNLP Tools:

TokensRegex Semgrep Tregex

Enter a Semgrep expression to run against the "enhanced dependencies" above:

{pos:/NNP.\*}=organization </nmod:poss/ ((lemma:founder) >dep {pos:/NNP.\*}=founder

Match

