

DAGOBDAH

Table and Graph Contexts For Efficient Semantic Annotation Of Tabular Data



Viet-Phi Huynh



Jixiong Liu



Yoan Chabot
[@yoan_chabot](#)



Frédéric Deuzé



Thomas Labbé
[@tau_labbe](#)



Pierre Monnin
[@piermonn](#)



Raphaël Troncy
[@rtroncy](#)



IC
29/06/2022



On the Importance of Interpreting Tabular Data

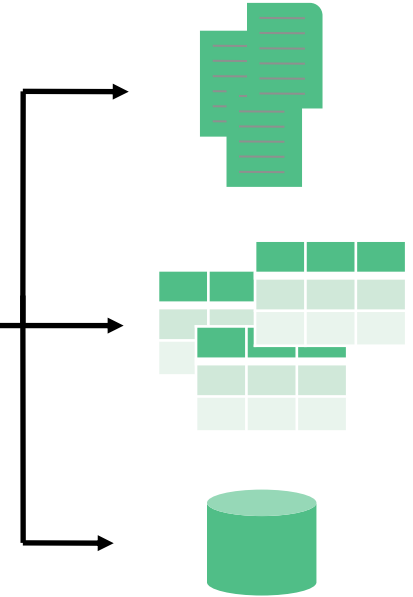
Multiservice operator



- Telecommunications (services + infrastructure)
- Video On Demand, Music
- Bank
- Cybersecurity
- ...

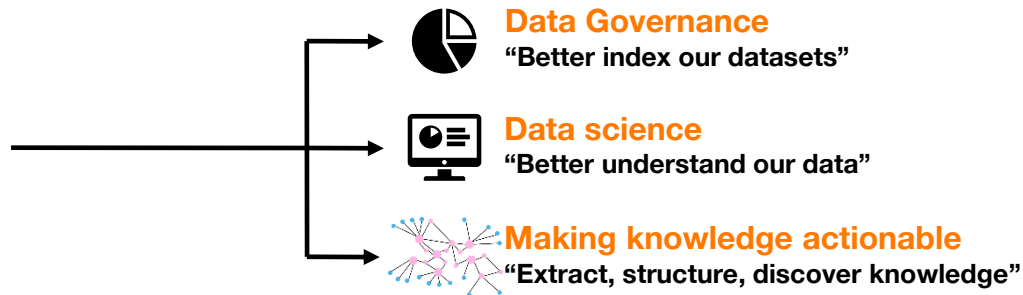
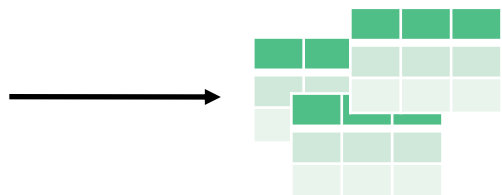
- 140,000 employees worldwide
- Heterogeneous client portfolio (individuals, companies)

<https://www.orange.com/en>
<https://hellofuture.orange.com/en/>



Huge amount of generated data

On the Importance of Interpreting Tabular Data



- Volume curse
- 7 languages (FR, EN, PL, ES, ...)
- Heterogeneously structured tables
- Little context

Cell types	Deployment environment	Max. number of users	Output power (mW)	Max. distance from base station	
5G NR FR2	Femtocell	Homes, businesses	Home: 4–8 Businesses: 16–32	indoors: 10–100 outdoors: 200–1000	tens of meters
	Pico cell	Public areas like shopping malls, airports, train stations, skyscrapers	64 to 128	indoors: 100–250 outdoors: 1000–5000	tens of meters
	Micro cell	Urban areas to fill coverage gaps	128 to 256	outdoors: 5000–10000	few hundreds of meters
	Metro cell	Urban areas to provide additional capacity	more than 250	outdoors: 10000–20000	hundreds of meters
Wi-Fi (for comparison)	Homes, businesses		indoors: 20–100		

Frequency (MHz)	Cell radius (km)	Cell area (km ²)	Relative cell count
450	48.9	7521	1
950	26.9	2269	3.3
1800	14.0	618	12.2
2100	12.0	449	16.2

	2012 T1	2012 T2	2012 T3
Réseaux en fibre optique jusqu'à l'abonné (FTTH) - déploiement et mutualisation			
Remarque liminaire			
Le total de locaux raccordables par zone ne correspond pas à la somme des lignes raccordables de la zone. Le nombre de lignes raccordables par opérateur d'infrastructure (OI) prend en compte l'ensemble des lignes raccordables alors que le total des locaux raccordables tient compte des déploiements multiples qui sont déduits de la somme.			
<i>Dans le cas d'un opérateur verticalement intégré, le nombre de lignes raccordables déployées par son activité d'opérateur d'infrastructure (OI) sur le marché de gros est différent du nombre de lignes éligibles aux offres de son activité d'opérateur commercial (OC) sur le marché de détail.</i>			
France	T1	T2	T3
Total des locaux raccordables et détail des lignes raccordables par opérateur.	1 776 000	1 873 000	2 157 000
Orange (OI - initiative privée)	875 000	918 000	1094 000
Alice France (OI - initiative privée)	369 000	388 000	449 000
Free Infrastructure (OI - initiative privée)	207 000	218 000	221 000
Autres OI	325 000	349 000	393 000
Lignes éligibles par nombre d'opérateurs présents au point de mutualisation :			
Au moins un opérateur présent au PM via la mutualisation passive	1580 000	1749 000	1960 000
	704 000	785 000	933 000
	311 000	349 000	539 000
	30 000	39 000	93 000

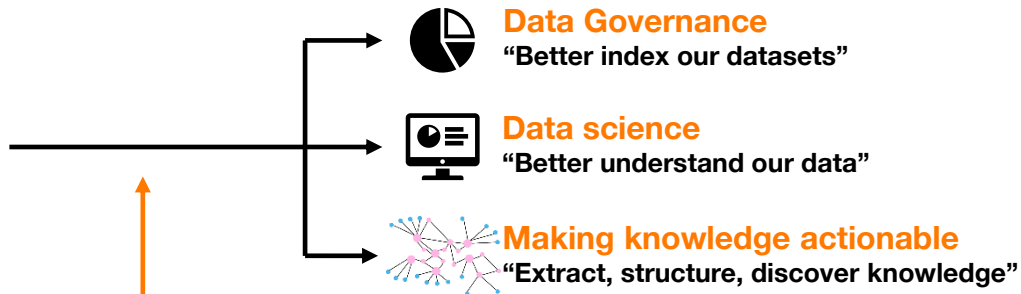
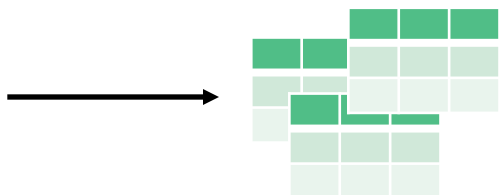
Indicateurs clés de l'activité des opérateurs de communications électroniques		2015
Emplois (champ : ancien cadre réglementaire)	Unités	
Emplois directs	Unités	110 470
<i>Source ARCEP - Enquêtes annuelles 1998 à 2010, enquêtes trimestrielles 2020</i>		
Investissements (champ : ancien cadre réglementaire)	Millions d'€	
Investissements au cours de l'exercice*	Millions d'€	10 630
dont investissements hors achats de fréquences mobiles	Millions d'€	7 931

* Cet indicateur intègre les montants payés pour l'achat de fréquences mobiles.

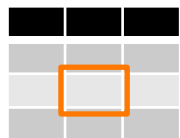
On the Importance of Interpreting Tabular Data



Documents



CEA
Cell-Entity Annotation

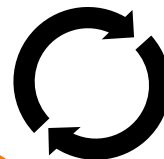


CTA
Column-Type Annotation

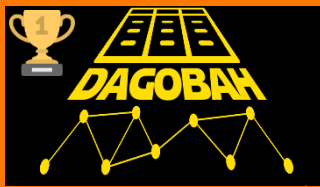


3 Tasks*

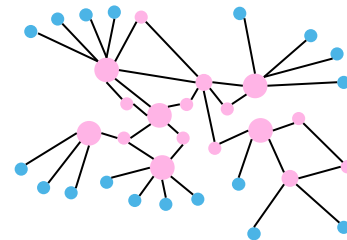
Virtuous loop



Semantic Table Interpretation



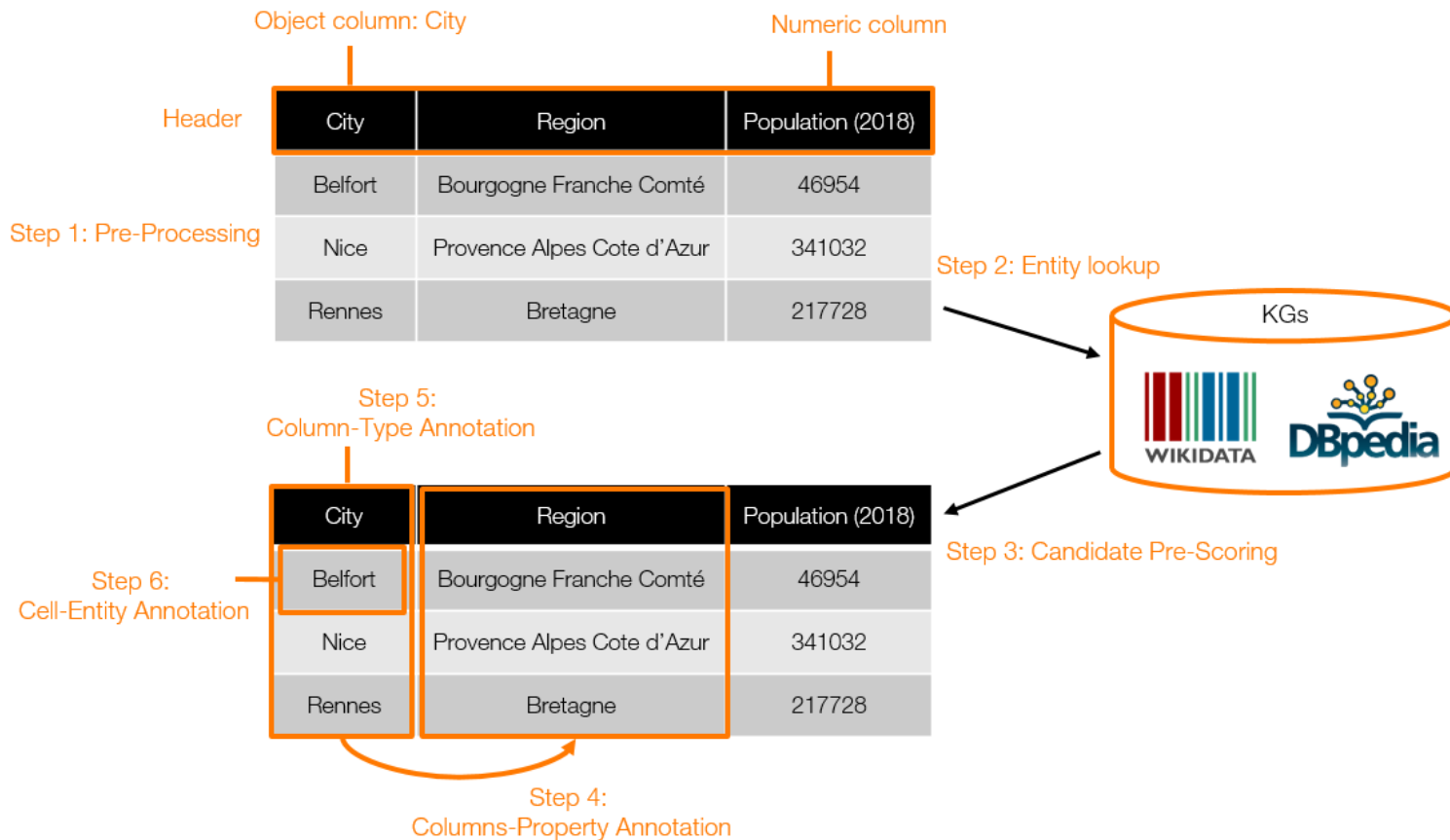
1st system in Accuracy
SemTab 2021



* Other possible tasks

Liu et al. From Tabular Data to Knowledge Graphs: A Survey of Semantic Table Interpretation Tasks and Methods. Submitted to Journal of Web Semantics, 2022.

DAGOBDAH: Annotation Workflow



Making Sense of Tables: DAGOBDAH

- Raw table:

Year	English title	Original title	Director(s)	Ref
2000	Lijmen/Het Been	Robbe De Hert	[22]	
2001	No Man's Land	NiÅ\u008Dija zemlja	Danis TanoviÅ\u008D	[23]
2002	The Son	Le Fils	Jean-Pierre and Luc Dardenne	[19]
2003	On the Run	Cavale	Lucas Belvaux	[24]
2003	An Amazing Couple	Un couple Å©patant	Lucas Belvaux	[24]

- Re-structured table:

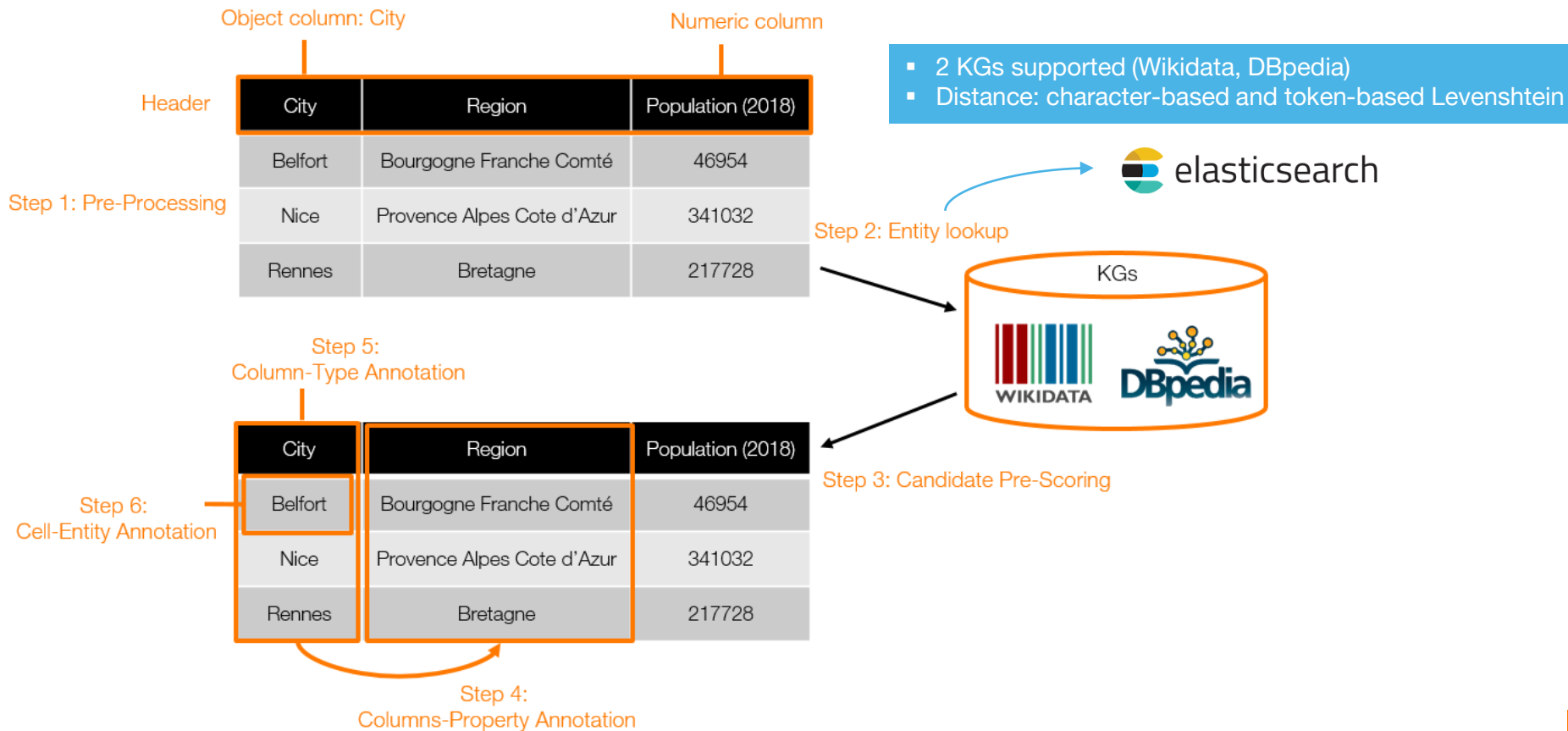
Year	English title	Original title	Director(s)	Ref
2000	Lijmen/Het Been		Robbe De Hert	[22]
2001	No Man's Land	NiÅ\u008Dija zemlja	Danis TanoviÅ\u008D	[23]
2002	The Son	Le Fils	Jean-Pierre and Luc Dardenne	[19]
2003	On the Run	Cavale	Lucas Belvaux	[24]
2003	An Amazing Couple	Un couple Å©patant	Lucas Belvaux	[24]

Pre-processing

Ad-hoc rules, spaCy NER

- Title: {}
- Orientation: HORIZONTAL, *Conf.=1.0*
- Header: ['Year', 'English title', 'Original title', 'Director(s)', 'Ref']
Conf= 0.6
- Key column: index=1, *Conf.=1.0*
- Data Type: [DateTime (*Conf* 1.0), String (1.0), String (1.0), String (1.0), String_Number(1.0)]
- Primitive typing: [DateTime (1.0), Unk, Unk, Person (1.0), Unk]
- Possible row re-alignment.

DAGOBDAH: Annotation Workflow



Making Sense of Tables: DAGOBAH

City	Region	Population (2018)
Belfort	Bourgogne Franche Comté	46954
Nice	Provence Alpes Cote d'Azur	341032
Rennes	Bretagne	217728

e_m

Context - $\mathcal{N}_{table}(e_m)$

$e_c \in \mathcal{E}_c$

- Q171545 : (Belfort) city in France
- Q4991979 : (Belfort) Neighborhood in Maastricht

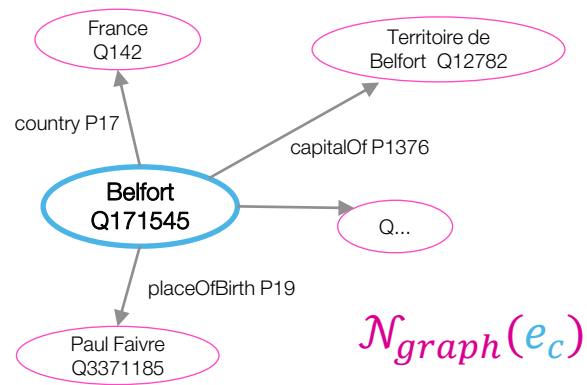
Candidate pre-scoring

\mathcal{E}_c set of candidates for cell e_m
 → confidence score for each $e_c \in \mathcal{E}_c$:

$$PSc(e_c, e_m) = Sc_{context}(\mathcal{N}_{table}(e_m), \mathcal{N}_{graph}(e_c)) * Sc_{sim}(e_c, e_m)^x$$

Context score of a candidate entity

→ quantifying the correspondance of neighboring cells with neighboring entities



$\mathcal{N}_{graph}(e_c)$

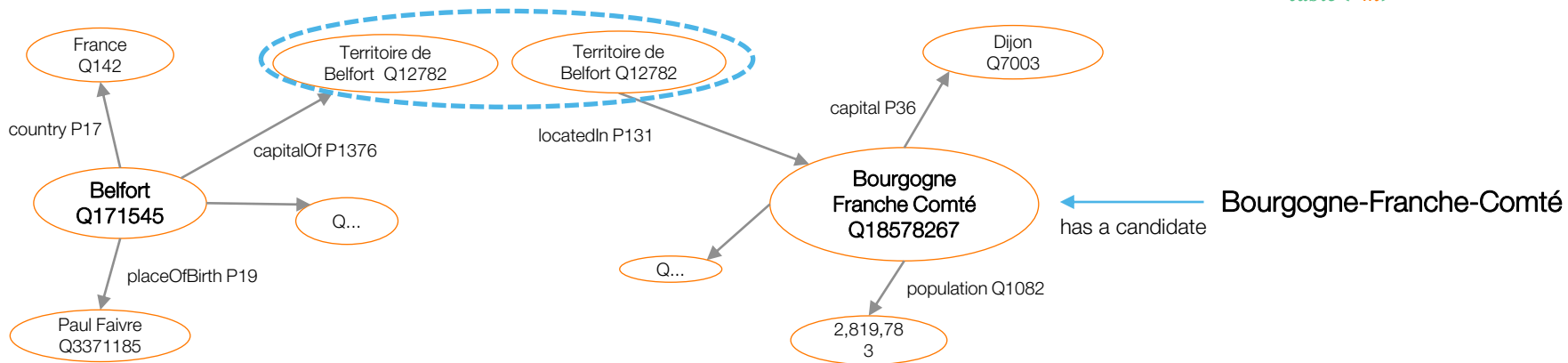
Effective Candidate Scoring

Relevance of cell **Bourgogne-Franche-Comté** w.r.t. **Q171545 (Belfort)** is determined through graph contexts discovery

City	Region	Population (2018)
Belfort	Bourgogne Franche Comté	46954
Nice	Provence Alpes Cote d'Azur	341032
Rennes	Bretagne	217728

e_m

Context - $\mathcal{N}_{table}(e_m)$



2-hop predicate path: (Q171545) Belfort is the capital of (Q12782) Territoire de Belfort which is located in (Q18578267) Bourgogne-Franche-Comté

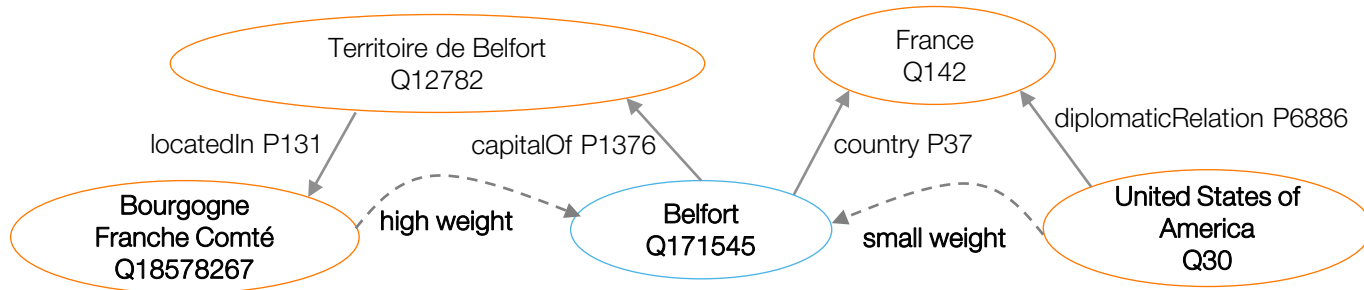
Candidate Soft Scoring

- Each neighbor does not have an equal contribution e_m

- ✓ **Object neighbors** are more important than literals (datetime, unit, plain number...)
- ✓ **Cells in the left side** of the table should have more weight (i.e. key column)
- ✓ **Columns highly connected** to target column should have more weight
- ✓ KG predicate paths provide **different information content**

City	Region	Population (2018)
Belfort	Bourgogne Franche Comté	46954
Nice	Provence Alpes Cote d'Azur	341032
Rennes	Bretagne	217728

Context - $\mathcal{N}_{table}(e_m)$



$generality(Q12782) \ll generality(Q142)$ → Q18578267 more informative

Making Sense of Tables: DAGOBAH

Candidate pre-scoring

\mathcal{E}_c set of candidates for cell e_m
 \rightarrow confidence score for each $e_c \in \mathcal{E}_c$:

$$PSc(e_c, e_m) = Sc_{context}(\mathcal{N}_{table}(e_m), \mathcal{N}_{graph}(e_c)) * Sc_{sim}(e_c, e_m)^x$$

City	Region	Population (2018)
Belfort	Bourgogne Franche Comté	46954
Nice	Provence Alpes Cote d'Azur	341032
Rennes	Bretagne	217728

e_m

$e_c \in \mathcal{E}_c$

Context - $\mathcal{N}_{table}(e_m)$

Q171545 : (Belfort) city in France
 Q4991979 : (Belfort) Neighborhood in Maastricht

Columns-Property Annotation

- The **most frequently occurring property** between best entities pre-scored in the two columns
- Supports multiple relation types, e.g. {s,p,o}, {s, p1, p2, o}, {s, p1, (-)p2, o}

Column-Type Annotation

- Majority voting based on pre-scored entities of the same column
- Type hierarchy considered with levels

Cell-Entity Annotation

- Pre-scoring = local information
- Final scoring = local information (pre-scoring) + global information (CTA and CPA)

$$Sc(e_c, e_m) = \frac{PSc(e_c, e_m) + \alpha \times \overline{score_{CPA}} + \beta \times score_{CTA}}{1 + \alpha + \beta}$$

$\alpha = CPA_coverage / 2$

$\beta = CTA_coverage / 2$

Performance: the SemTab 2021 challenge



1st system in Accuracy
SemTab 2021

Settings

Setting	Entity's KG context	Scoring
1	1-hop predicate paths	Hard scoring
2	2-hop predicate paths	Hard scoring
3	2-hop predicate paths	Soft scoring
4	2-hop, but using only uni-directional paths, i.e. sub $\vec{p}_1 \vec{p}_2 obj$	Soft scoring

100 candidates / cell: 49s / table
250 candidates / cell: 92s / table

LeaderBoard

Task	Setting	CTA		CEA		CPA	
		F1	Precision	F1	Precision	F1	Precision
Round 1 – WD Table	2*	0.832	0.832	0.923	0.923	-	-
	Top 1 SemTab 2021	0.728	0.73	0.907	0.907	-	-
Round 1 - DBPTable	2*	0.422	0.424	0.954	0.946	-	-
	Top 1 SemTab 2021	0.46	0.468	0.692	0.692	-	-
Round 2 - BioTable	4*	0.916	0.916	0.970	0.970	0.899	0.899
	Top 1 SemTab 2021	0.956	1	0.964	0.964	0.899	0.899
Round 2 - HardTable	3*	0.976	0.976	0.975	0.976	0.996	0.996
	Top 1 SemTab 2021	0.977	0.977	0.985	0.985	0.997	0.998
Round 3 - BioDivTable	4*	0.381	0.382	0.496	0.497	-	-
	Top 1 SemTab 2021	0.593	0.595	0.602	0.611	0.947	1
Round 3 – HardTable	3*	0.99	0.99	0.974	0.974	0.991	0.995
	Top 1 SemTab 2021	0.984	0.984	0.968	0.968	0.993	0.994

DAGOBAN Science Health Use Case

Context

A **healthcare practitioner** exploring **two datasets** containing information about **drugs**.

health_drugs_class.csv

NOM	PRODUIT	Code EphMRA	Classe EphMRA	Classe ATC
YELLOX 0,9 MG\u002fML COLLYRE FL 1\u002f5 ML	YELLOX	S01R	ANTIINFLAMMATOIRES NON STEROIDIENS OPHTALMIQUES	BROMFENAC SODIUM SESQUIHYDRATE
ESOMEPRAZOLE BIOGARAN 20 MG CPR GASTRORESISTANT 28	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
AVONEX 30 MCG\u002f0,5 ML SOL INJ STYLO 4\u002f0,5 ML	AVONEX	L03B2	INTERFERONS, BETA	INTERFERON BETA%\u0080\u0093A
ESOMEPRAZOLE BIOGARAN 40 MG CPR GASTRORESISTANT 14	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
VALSARTAN\u002fHYDROCHLOROTHIAZIDE CRISTERS	CGP-48933	C09D1	ASSOCIATIONS ANTAGONISTES DE	VALSARTAN ET DIURETIQUES

health_chemical_components.csv

Analyte	CAS Registry Number	Units	LCMRL
oxycodone	76-42-6	nanogram per litre	4,7
sulfamethoxazole	723-46-6	nanogram per litre	6,5
trimethoprim	738-70-5	nanogram per litre	3,5
valsartan	137862-53-4	nanogram per litre	7,2
verapamil	52-53-9	nanogram per litre	7,8

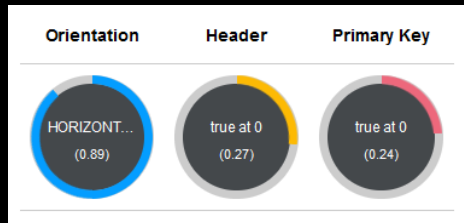
DAGOBAH Science Health Use Case

Context

A **healthcare practitioner** exploring **two datasets** containing information about **drugs**.

DAGOBAH Goals

- Clean datasets** automatically (e.g. encoding issues, misalignement, etc.)



health_drugs_class.csv

NOM	PRODUIT	Code EphMRA	Classe EphMRA	Classe ATC
YELLOX 0,9 MG <u>u002f</u> ML COLLYRE FL 1 <u>u002f</u> 5 ML	YELLOX	S01R	ANTIINFLAMMATOIRES NON STEROIDIENS OPHTALMIQUES	BROMFENAC SODIUM SESQUIHYDRATE
ESOMEPRAZOLE BIOGARAN 20 MG CPR GASTRORESISTANT 28	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
AVONEX 30 MCG <u>u002f</u> 0,5 ML SOL INJ STYLO 4 <u>u002f</u> 0,5 ML	AVONEX	L03B2	INTERFERONS, BETA	INTERFERON BETA% <u>u0080</u> <u>u0093</u> A
ESOMEPRAZOLE BIOGARAN 40 MG CPR GASTRORESISTANT 14	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
VALSARTAN <u>u002f</u> HYDROCHLOROTHIAZIDE CRISTERS	CGP-48933	C09D1	ASSOCIATIONS ANTAGONISTES DE	VALSARTAN ET DIURETIQUES



NOM	PRODUIT	Code EphMRA	Classe EphMRA	Classe ATC
YELLOX 0,9 MG <u>u002f</u> ML COLLYRE FL 1 <u>u002f</u> 5 ML	YELLOX	S01R	ANTIINFLAMMATOIRES NON STEROIDIENS OPHTALMIQUES	BROMFENAC SODIUM SESQUIHYDRATE
YELLOX 0,9 MG <u>u002f</u> ML COLLYRE FL 1 <u>u002f</u> 5 ML				
ESOMEPRAZOLE BIOGARAN 20 MG CPR GASTRORESISTANT 28	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
AVONEX 30 MCG <u>u002f</u> 0,5 ML SOL INJ STYLO 4 <u>u002f</u> 0,5 ML	AVONEX	L03B2	INTERFERONS, BETA	INTERFERON BETA% <u>u0080</u> <u>u0093</u> A
AVONEX 30 MCG <u>u002f</u> 0,5 ML SOL INJ STYLO 4 <u>u002f</u> 0,5 ML				
ESOMEPRAZOLE BIOGARAN 40 MG CPR GASTRORESISTANT 14	ESOMEPRAZOLE BIOGARAN	A02B2	INHIBITEURS DE LA POMPE A PROTONS	ESOMEPRAZOLE
VALSARTAN <u>u002f</u> HYDROCHLOROTHIAZIDE CRISTERS	CGP-48933	C09D1	ASSOCIATIONS ANTAGONISTES DE L'ANGIOTENSINE II ET ANTIHYPERTENSEURS (C2) ET <u>u002f</u> DIURETIQUES	VALSARTAN ET DIURETIQUES
VALSARTAN <u>u002f</u> HYDROCHLOROTHIAZIDE CRISTERS			ASSOCIATIONS ANTAGONISTES DE L'ANGIOTENSINE II ET ANTIHYPERTENSEURS (C2) ET <u>u002f</u> DIURETIQUES	
VALSARTAN <u>u002f</u> HYDROCHLOROTHIAZIDE CRISTERS	CGP-48933	C09D1	ASSOCIATIONS ANTAGONISTES DE L'ANGIOTENSINE II ET ANTIHYPERTENSEURS (C2) ET <u>u002f</u> DIURETIQUES	VALSARTAN ET DIURETIQUES
VALSARTAN <u>u002f</u> HYDROCHLOROTHIAZIDE CRISTERS			ASSOCIATIONS ANTAGONISTES DE L'ANGIOTENSINE II ET ANTIHYPERTENSEURS (C2) ET <u>u002f</u> DIURETIQUES	

DAGOBAH Science Health Use Case

Context

A **healthcare practitioner** exploring **two datasets** containing information about **drugs**

DAGOBAH Goals

1. **Clean datasets** automatically (e.g. encoding issues, misalignement, etc.)
2. Make table **semantics explicit** and **align datasets**

health_drugs_class.csv

NOM	PRODUIT	Code EphMRA	Classe EphMRA	Classe ATC
<p>CTA 0.04 medication Q12140 Coverage 33%</p> <p>ESOMEPRAZOLE BIOGARAN 40 MG CPR GASTRORESISTANT 14</p>	<p>CTA 0.10 pharmaceutical product Q28885102 Coverage 100%</p> <p>ESOMEPRAZOLE BIOGARAN esomeprazole Q553223 0.14</p>	A02B2	<p>CTA 0.06 class Q16889133 Coverage 50%</p> <p>INHIBITEURS DE LA POMPE A PROTONS proton-pump inhibitors Q421704 0.16</p>	<p>CTA 0.15 chemical compound Q11173 Coverage 100%</p> <p>ESOMEPRAZOLE esomeprazole Q553223 0.18</p>
<p>VALSARTAN/HYDROCHLOROTHIAZIDE CRISTERS 0.12 Hydrochlorothiazide / valsartan Q48566694</p>	<p>CGP-48933 0.12 valsartan Q155472</p>	C09D1	<p>ASSOCIATIONS ANTAGONISTES DE L'ANGIOTENSINE II ET ANTIHYPERTENSEURS (C2) ET/OU DIURETIQUES</p>	<p>VALSARTAN ET DIURETIQUES 0.12 valsartan Q155472</p>

health_chemical_components.csv

Analyte	CAS Registry Number	Units	LCMRL
<p>CTA 0.02 medication Q12140 Coverage 78%</p> <p>trimethoprim 0.05 trimethoprim/sulfadoxine Q7842230</p>	<p>CTA 0.04 medication Q12140 Coverage 78%</p> <p>738-70-5 0.07 trimethoprim Q422665</p>	<p>CTA 0.01 SI-accepted non-SI unit Q106839753 Coverage 100%</p> <p>nanogram per litre 0.01 nanogram per litre Q104907186</p>	3,5
<p>valsartan 0.02 valsartan Q155472</p>	<p>137862-53-4 0.00 HD 137862 Q84374092</p>	<p>nanogram per litre 0.01 nanogram per litre Q104907186</p>	7,2

DAGOBAH Science Health Use Case

Context

A **healthcare practitioner** exploring **two datasets** containing information about **drugs**.

DAGOBAH Goals

1. **Clean datasets** automatically (e.g. encoding issues, misalignment, etc.)
2. Make table **semantics explicit** and **align datasets**
3. **Enrich tables** with new information from the Knowledge Graph (fill missing values, new columns)

health_chemical_components.csv

Analyte	CAS Registry Number	Units	LCMRL
oxycodone	76-42-6	nanogram per litre	4,7
sulfamethoxazole	723-46-6	nanogram per litre	6,5
trimethoprim	738-70-5	nanogram per litre	3,5
valsartan	137862-53-4	nanogram per litre	7,2
verapamil	52-53-9	nanogram per litre	7,8



Analyte	CAS Registry Number	Units	LCMRL	mass	chemical formula
sulfamethoxazole Sulfamethoxazole / Trimethoprim Q48566821 0.05	723-46-6 sulfamethoxazole Q415843 0.08	nanogram per litre nanogram per litre Q104907186 0.01	6,5	Knowledge CEA 253.052 253.052253.052	Knowledge CEA <chem>C10H11N3O2S</chem> <chem>C10H11N3O2SC10H11N3O2S</chem>
trimethoprim trimethoprim/sulfadoxine Q7842230 0.05	738-70-5 trimethoprim Q422665 0.07	nanogram per litre nanogram per litre Q104907186 0.01	3,5	Knowledge CEA 290.138 290.138290.138	Knowledge CEA <chem>C11H16N4O3</chem> <chem>C11H16N4O3C11H16N4O3</chem>

Future work

Annotation enhancement

- ✓ Generate global tokens dictionaries (abbreviations, acronyms...) from web data
- ✓ Leverage embeddings-based approach (clustering, language models applied to table...)

Data augmentation: from tables to KGs

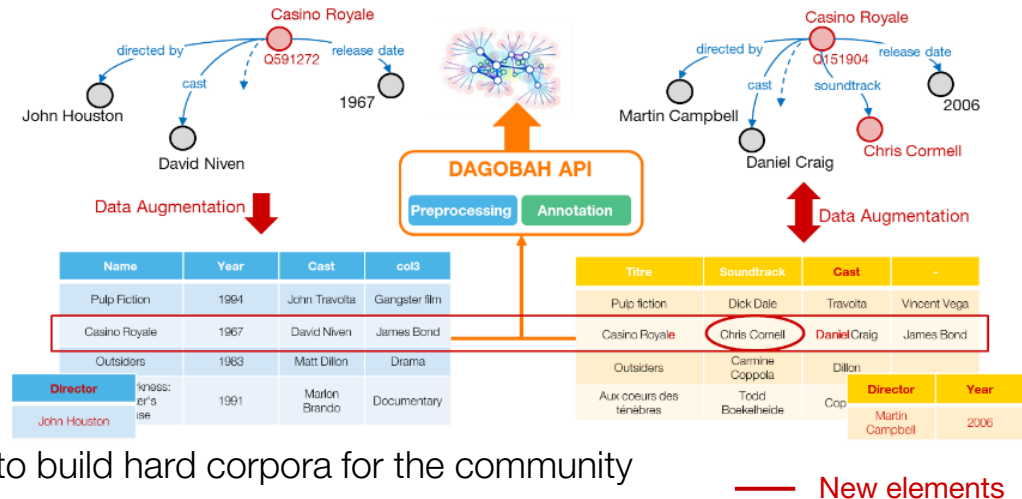
- ✓ Enrich Knowledge graphs leveraging CTA and CPA to add new entities

Towards corpora with new challenges

- ✓ SemTab 2021 challenging corpora (new domains, new target KG)
- ✓ But scores still very high → call to arms to build hard corpora for the community

Towards a KG targeting Orange business areas

- ✓ Current effort to build such a KG supported with DAGOBDAH
- ✓ Challenges: heterogeneous domains and vocabularies



DAGOBAB

DAGOBAB: Table and Graph Contexts For Efficient Semantic Annotation Of Tabular Data

Thanks!



Sarthou-Camy et al. (2022). DAGOBAB UI: A New Hope For Semantic Table Interpretation. In Demo Track ESWC2022.

Huynh et al. (2021). DAGOBAB: Table and Graph Contexts For Efficient Semantic Annotation Of Tabular Data. SemTab@ISWC (pp. 19-31)

Chabot et al. (2021). A Framework for Automatically Interpreting Tabular Data at Orange. In Industry Track ISWC2021.

Huynh et al. (2020). DAGOBAB: Enhanced Scoring Algorithms for Scalable Annotations of Tabular Data. In SemTab@ ISWC (pp. 27-39).

Chabot et al. (2019). DAGOBAB: An End-to-End Context-Free Tabular Data Semantic Annotation System. In SemTab@ ISWC (pp. 41-48).

DAGOBAB: Make Tabular Data Speak Great Again : <https://hellofuture.orange.com/en/dagobah-make-tabular-data-speak-great-again/>