

Generic Framework for the Multidimensional Processing and Analysis of Social Media Content

"A Proxemic Approach"

Cotutelle Computer Science Ph.D. Defense Presented By
Maxime Masson

Supervisors: Dr. Christian Sallaberry, Dr. Rodrigo Agerri




Advisors: Dr. Marie-Noelle Bessagnet, Prof. Philippe Roose, Dr. Annig Le Parc Lacayrelle

Reviewers: Prof. Josiane Mothe, Prof. Elena Cabrio





Examiners: Research Dir. Ana-Maria Olteanu-Raimond, Research Dir. Maguelonne Teisseire

Outline

1 Introduction

-  Context
-  APs project and motivations
-  Framework and contributions overview

2 Contributions

-  **Collect:** Generic and iterative methodology for constructing thematic datasets from social media
-  **Transform:** Optimal strategies for the multilingual analysis of social media content in tourism
-  **Analyze:** Redefining proxemics to model social media entities and their interactions to generate domain-adaptable indicators from social media
-  **Valorize:** Interactive visualization of multidimensional analyses from social media

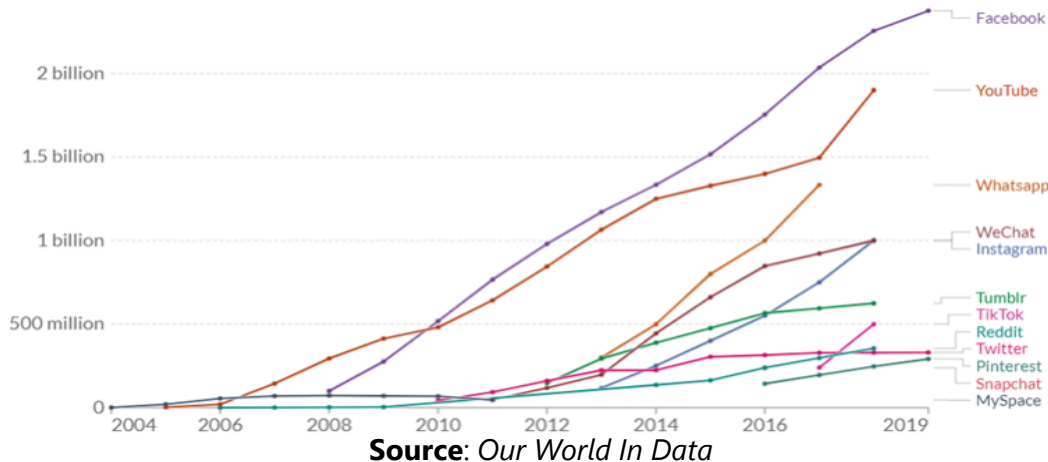
3 Conclusion and future perspectives

Context

User-Generated Content and Social Media

↗ Significant growth in **data sources** available in many domains

🌐 **Web 2.0** and **User-Generated Content (UGC)**



48.3% of the **world's population** using **social media** in 2020



Massive content (per day)

- 500 millions X tweets
- 216 millions *Facebook* messages

Source: Statista

Context

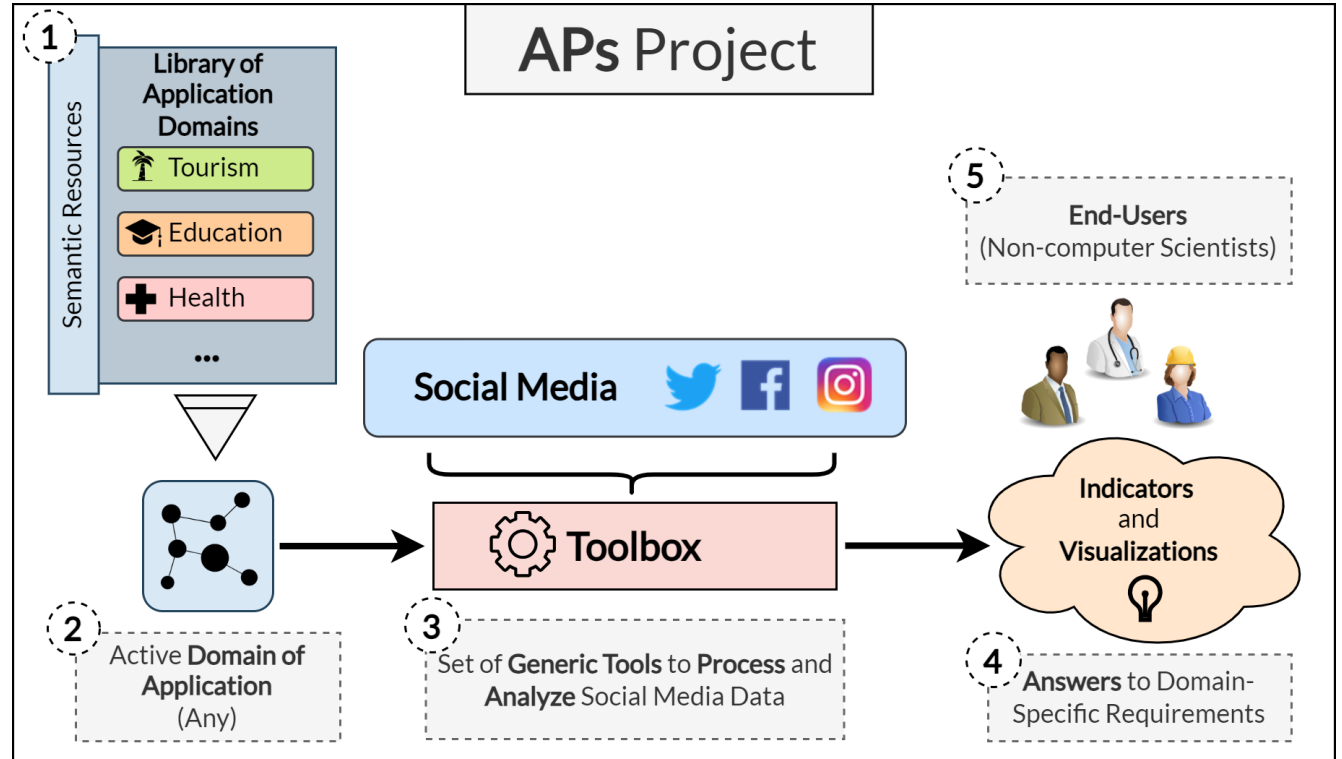
APs Project



Specificity:

Use of the theory of

Proxemics



Motivations I

Decision Support in the Tourism Domain



Tourism Professionals

- Assistance in the **decision-making** process and infrastructure planning
- Understanding the **requirements, practices** and **expectations** of visitors

A

What leisure activities do tourists typically engage in together?

B

Which cities do tourists tend to go to after visiting Bayonne?

C

What are the typical demographics of tourists who visit Biarritz?

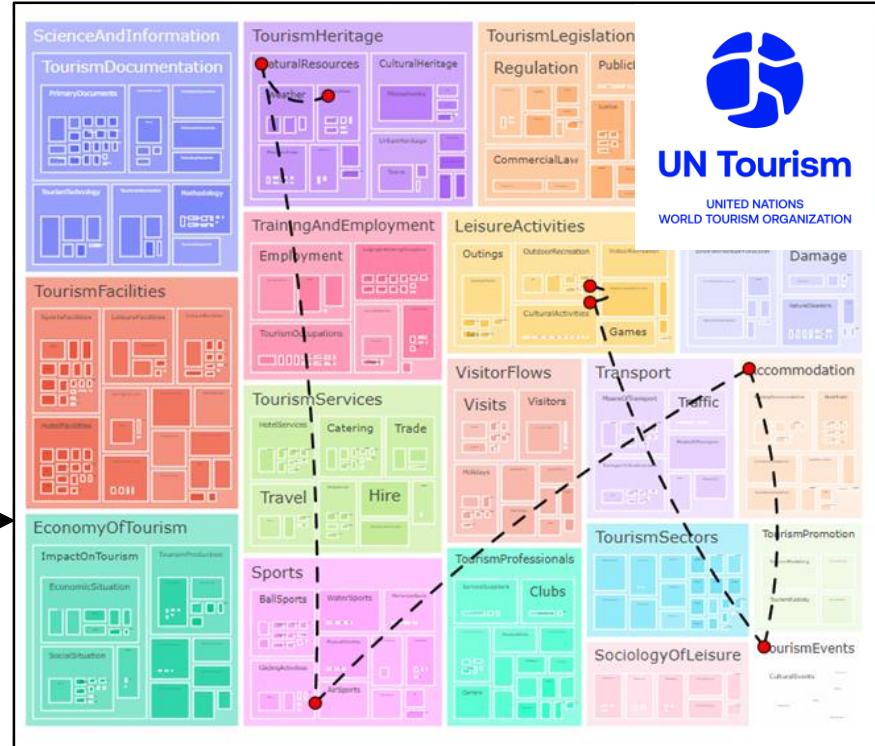
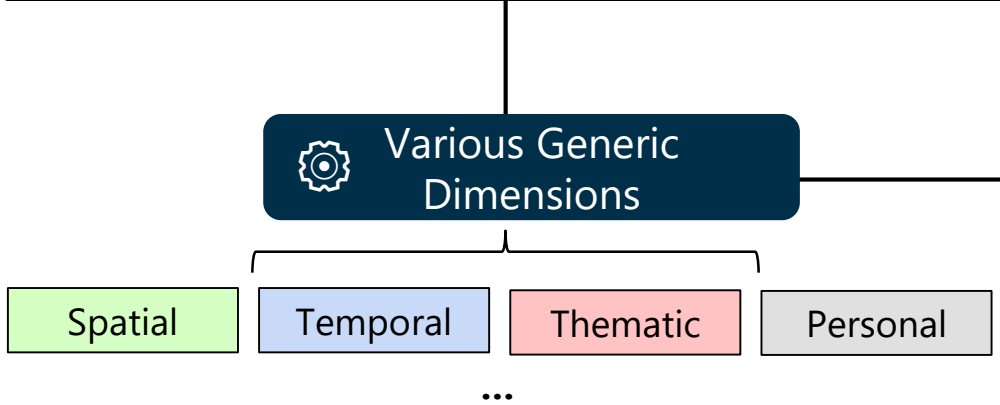
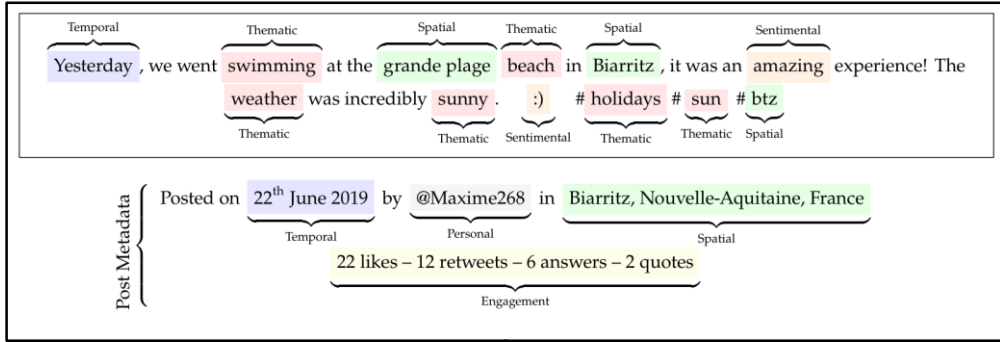
D

What are the common chained tourist activities?



Motivations II

Main Hypothesis



Contributions

APs Framework (Life Cycle)



Double Genericity

- Domain of Application
- Social Media Source



Semantic resource to describe the domain



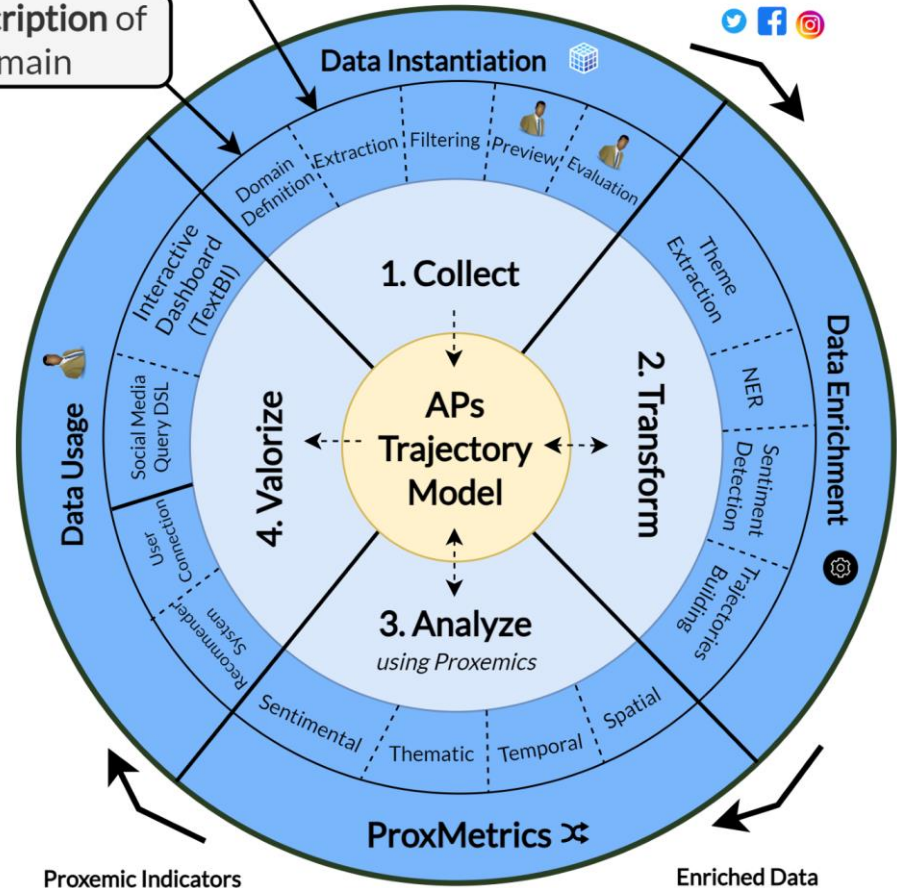
Decision support for non-technical users

- Indicators and Visualizations

   **Social Media Sources**

 **Description of a Domain**

Raw Data
(social media posts and users)



Phase 1: Collect

Raw Dataset

Social Media Posts and Users

Collect

Transform

Analyze

Valorize



International Conference (CORE: B)

International Conference on Web Information
Systems Engineering (WISE 2022)

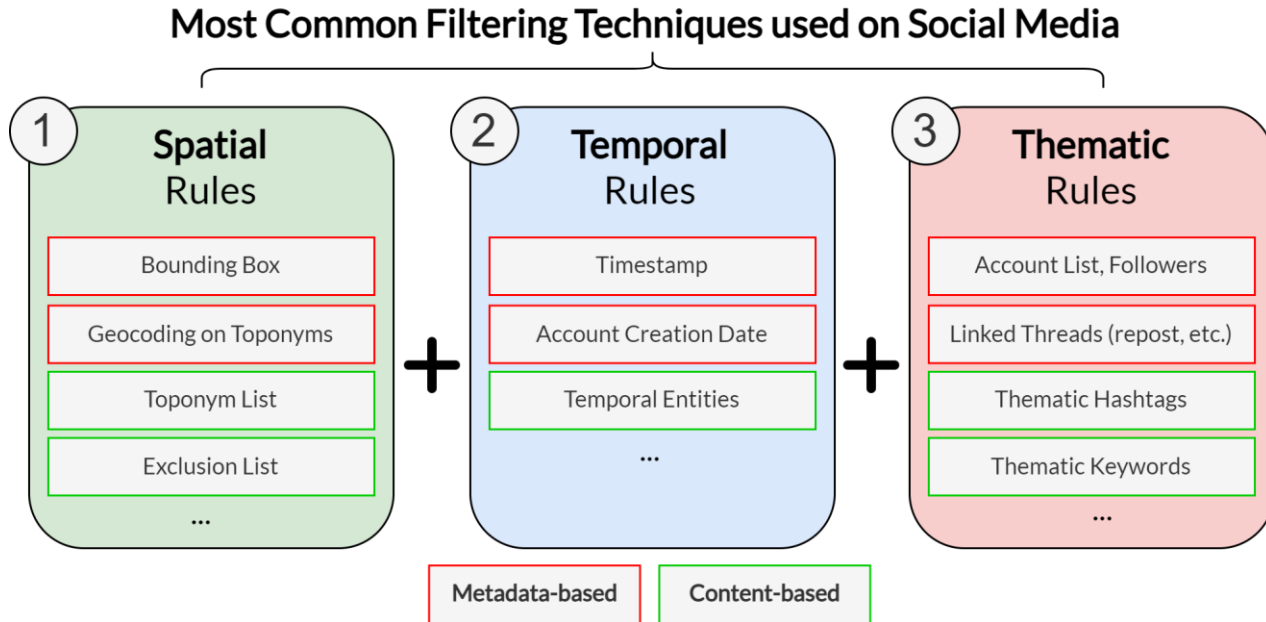
Research Challenge and Hypothesis

Constructing Accurate and Representative Social Media Datasets

- ⚠ **Challenge:** Constructing accurate and representative social media datasets
 - Social media are massive and noisy
 - Applicable across various social media and domain of application
- ❓ **Hypothesis:** High-level generic methodology to build social media datasets
 - Iterative and incremental process
 - Human feedback
 - Semantic domain description
 - Various existing filtering techniques

Related Work

Existing Dataset Building Approaches from Social Media



No **high-level, generic** collection approach. Mostly **ad-hoc implementations**

Generic and Iterative Methodology for Constructing Thematic Datasets from Social Media

Filtering Process and Iterations

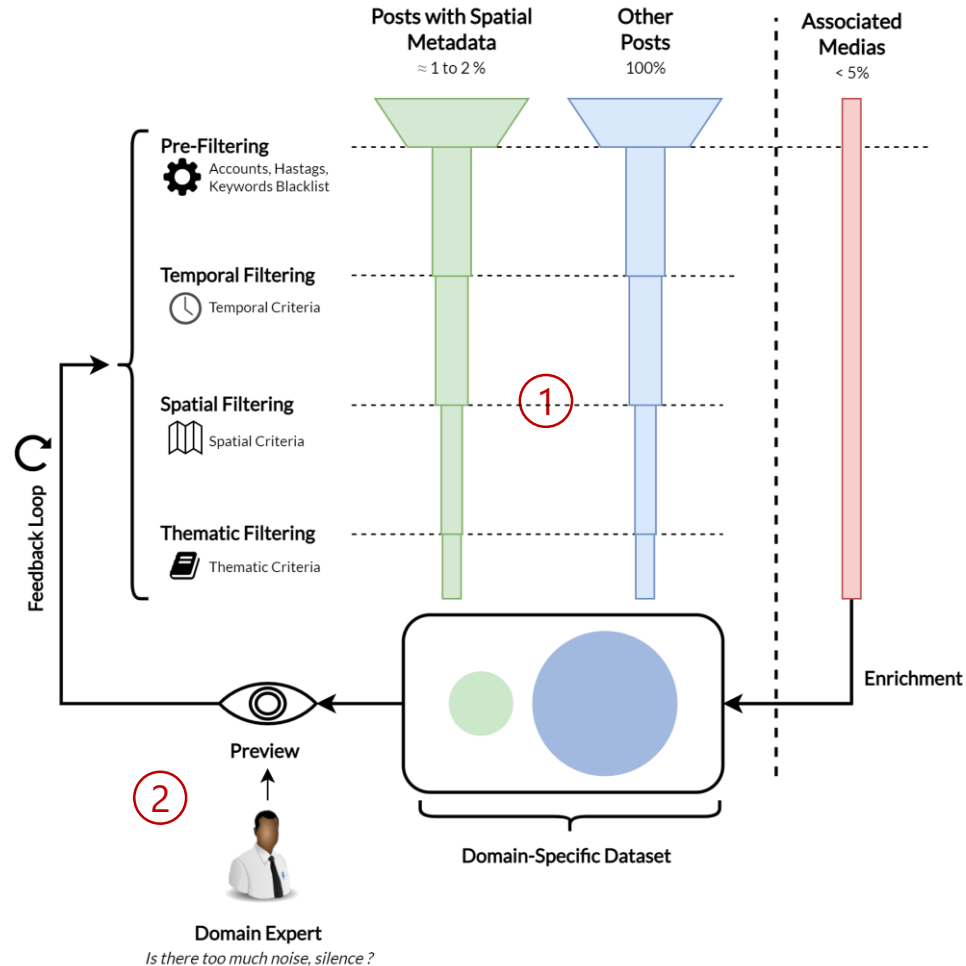
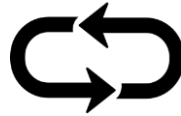


End-User Involvement

- **Evaluate** the resulting dataset
- **Submit feedback**

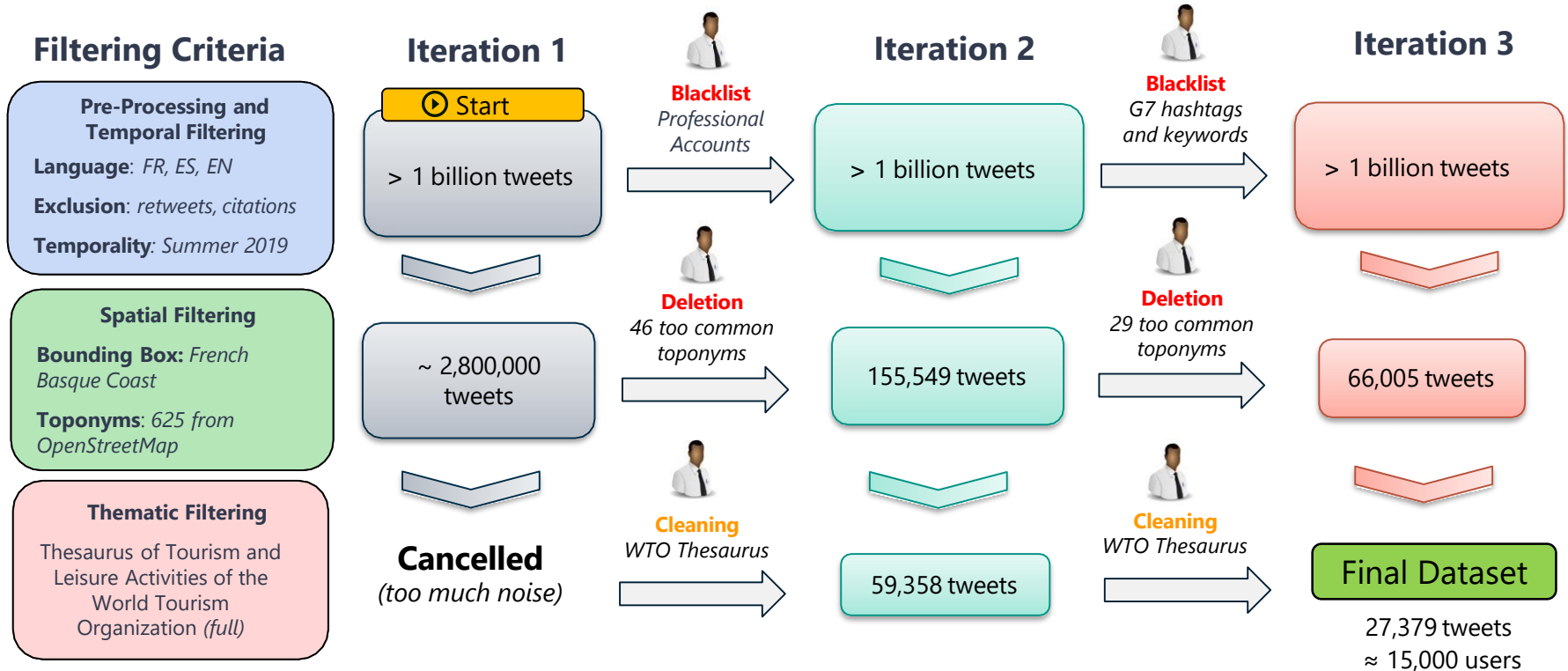
Feedback Loops

- Indefinite number of **iterations**
- Until the dataset is **deemed satisfactory**



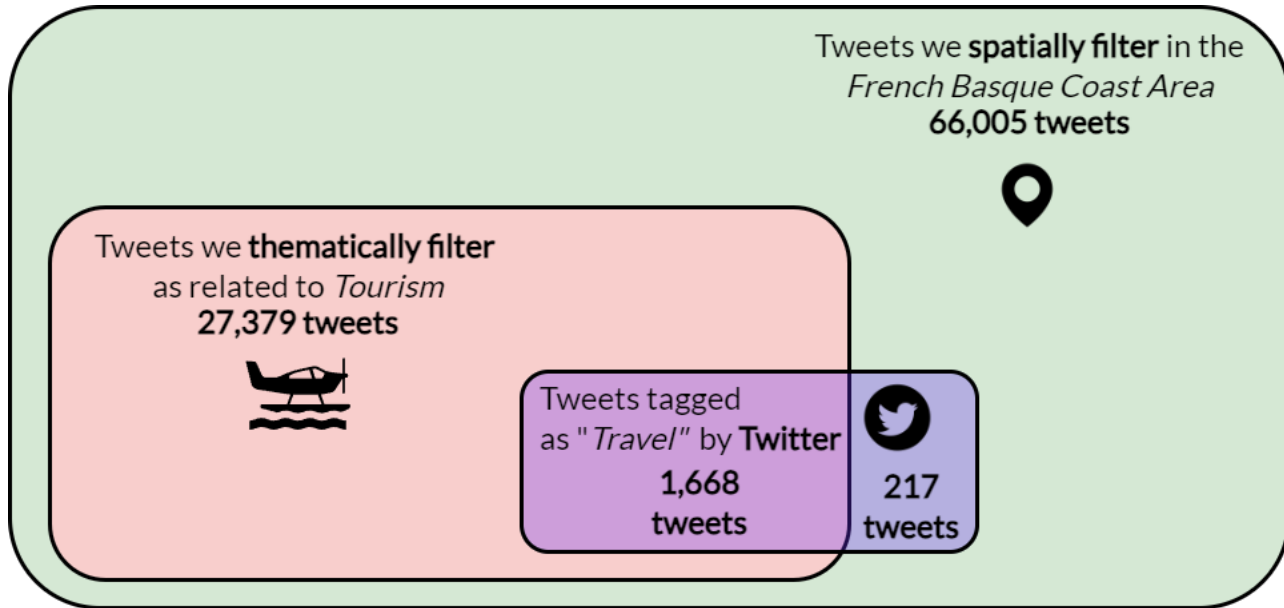
Experimentation

Definition and Filtering Process



Evaluation

Quantitative Analysis: Comparison with X/Twitter's Travel Context Annotations



? Is what we collect additionally **relevant**? Or is it just **noise**?

Evaluation

Qualitative Analysis: Dataset Accuracy

		Iteration 1		Iteration 2		Iteration 3	
		Geotagged	Others	Geotagged	Others	Geotagged	Others
Accuracy	(@ 20)	0.75	< 0.1	0.60	0.30	0.83	0.72
	(@ 50)	0.64		0.60	0.30	0.77	0.74
	(@ 100)	0.52		0.59	0.35	0.74 (κ 0.74)	0.65 (κ 0.48)

- ✓ Potentially **65% to 83% of the tweets collected** could be pertinent
- ✓ Highlight the role of **iterations** in improving **accuracy**

Phase 2: Transform

Raw Dataset
Social Media Posts and Users

Enriched Dataset
Sentiments, Locations, Themes



National Conference

Conference on Computer Science for
Organizations and Information and Decision
Systems (INFORSID 2024)

Research Challenge

Knowledge Extraction for Social Media Content

- ⚠ **Challenge 1:** Identifying best knowledge extraction strategies and models for a given application domain and task
- ⚠ **Challenge 2:** Determining the number of domain-specific annotated examples needed to get satisfying results
 - Manually annotating → **Lengthy, costly, time-consuming**
 - Objective → Get optimal results with **minimal use** of annotated examples
- ❓ **Hypothesis:** Use of a comparative study on optimal strategies for multilingual analysis of social media content based on a novel annotated dataset
- ⓘ **Limitation:** Tourism domain, social media texts

Introduction

Common Knowledge Extraction Tasks

Sentiment Analysis

Positive { Yesterday, we went swimming at the Grande Plage beach in Biarritz, it was an amazing experience! The weather was **incredibly** sunny. :) #holidays #sun #btz

Named Entity

Recognition (NER)

for Locations

Yesterday, we went swimming at the **Grande Plage** beach in **Biarritz**, it was an amazing experience! The weather was incredibly sunny. :) #holidays #sun # **btz**

Location Location Location

Fine-Grained Thematic Concept Extraction

Yesterday, we went **swimming** at the Grande Plage **beach** in Biarritz, it was an amazing experience! The **weather** was incredibly **sunny** . :)

holidays # **sun** #btz

Sports::WaterSports::Swimming NaturalResources::Beaches NaturalResources::Weather ClimaticFactors::Sun VisitorFlows::Holidays ClimaticFactors::Sun

Related Work

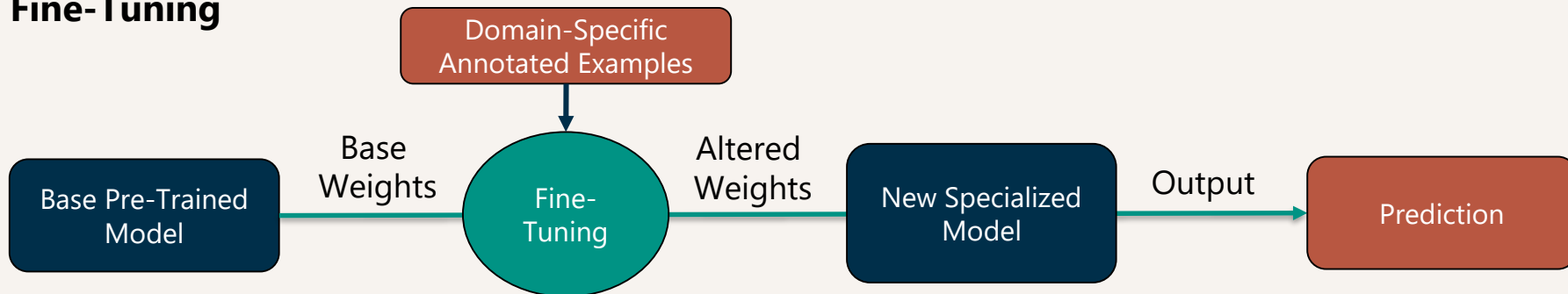
Rule-based Techniques for NER for Locations

Technique based on ...	Advantages (+)	Disadvantages (-)
Lexicon	Easy to implement, easily understandable	Requires a lexicon, limited by lexicon size, ignores context and grammatical structures
Patterns	Precise for well-defined patterns	Missing variations not covered by patterns , requires properly formatted sentences
Syntax and grammar	Exploits linguistic structures for deeper analysis	Complex to maintain , especially in multilingual contexts
Semantics	Can understand nuanced meanings and relationships between terms	Require comprehensive semantic knowledge bases, more computationally intensive

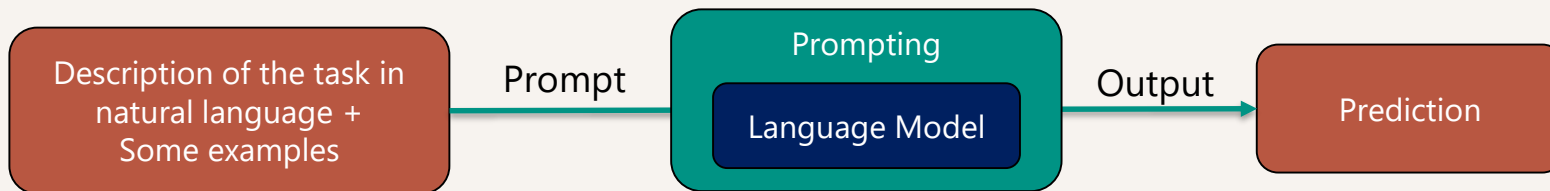
Related Work

Fine-Tuning and Zero-Shot Prompting

Fine-Tuning



Zero-Shot or Few-Shot Prompting



Experimental Setup

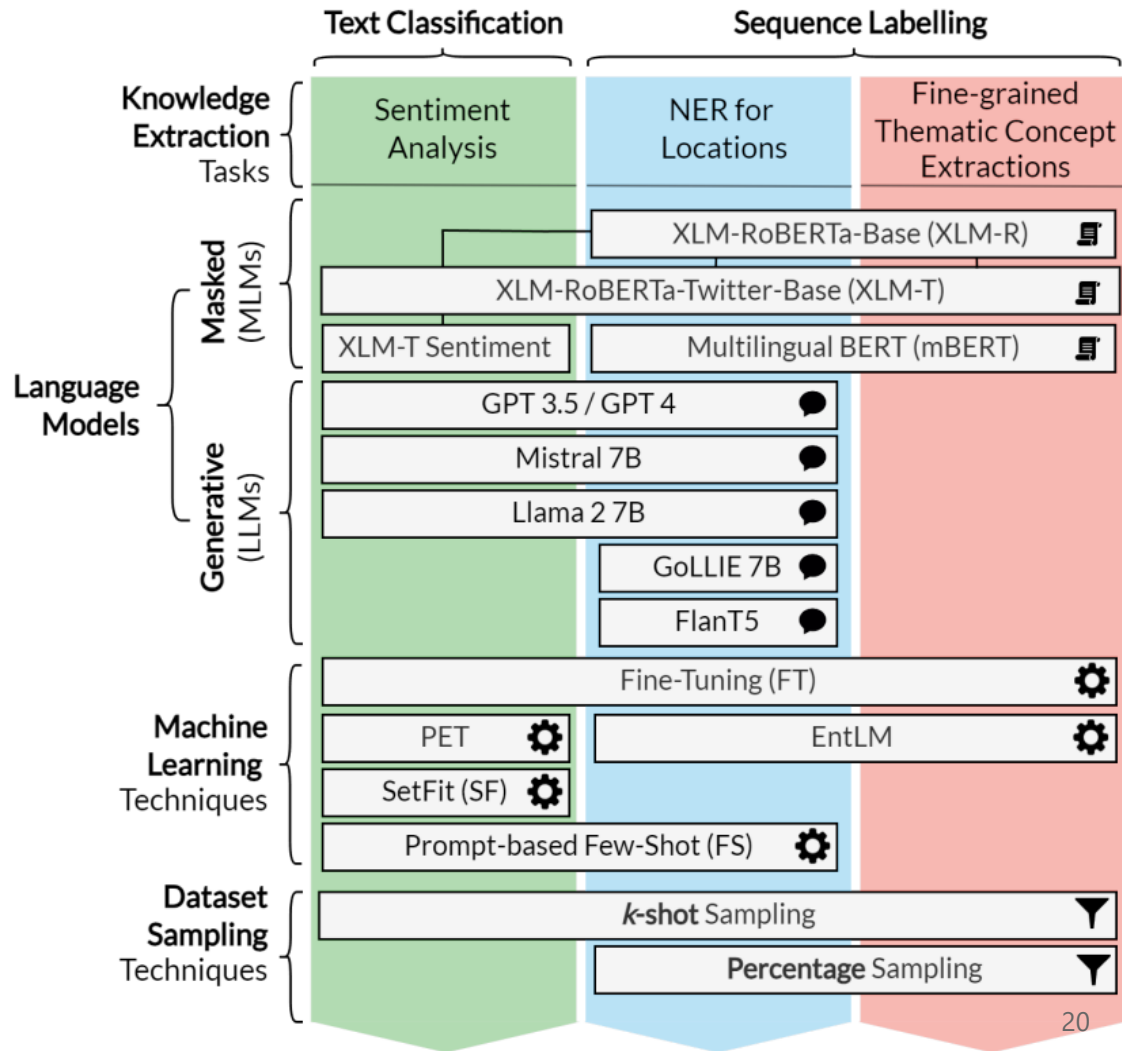
Comparative Analysis

- **Multilingual Dataset**

- 2961 tweets, 624 users



1. Language Models
2. Machine Learning Techniques
3. Dataset Sampling Methods



Results

NER for Locations

- Rules-based F1 : 0,707

- Less than 100 examples

A Few-shot prompting with LLMs

- More than 100 examples

B Fine-Tuning with Large (LLMs) or Masked Language Models (MLMs)

C GoLLIE

Techniques	Examples per class (location) — F1-score									
	0	5	10	20	30	40	50	100	All	
Prompt-based FS	Regular Prompt-based Few-Shot of LLMs									
GPT 3.5	0.694	0.698	0.762	0.762	0.798	0.809	0.828	0.806		
Mistral 7B	0.680	0.704	0.689	0.730	0.749	0.741	0.742	0.739		
LLaMA 2 7B	0.627	0.587	0.615	0.594	0.621	0.580	0.568	0.169		
FT of MLMs	Fine-Tune of Encoder-Only Models (MLMs)									
XLM-T		0.067	0.113	0.001	0.029	0.000	0.067	0.054	0.802	
XLM-R		0.107	0.067	0.130	0.062	0.328	0.133	0.001	0.791	
mBERT		0.115	0.108	0.083	0.007	0.000	0.000	0.000	0.818	
FT of LLMs	Fine-Tune of Encoder-Decoder and Decoder-Only Models (LLMs)									
LLaMA 2 7B		0.000	0.000	0.000	0.000	0.000	0.000	0.228	0.701	
FlanT5		0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.806	
EntLM	Template-Free Few-Shot in Sequence Labeling Tasks for MLMs									
mBERT		0.317	0.385	0.437	0.529	0.562	0.591	0.584	0.788	
GoLLIE	Guideline following model for Information Extraction									
GoLLIE 7B	0.670	0.622	0.632	0.662	0.661	0.694	0.689	0.732	0.832	

Discussion

Named Entity Recognition for Locations

- ✓ Task with **few classes** but **many label words**
 - Low representativeness of label words
- ✓ If it is possible **to annotate many examples, fine-tuning with MLMs** works very well
- ✓ Otherwise, **30 examples are enough** to **obtain satisfactory results** in few-shot with LLMs
- ✓ **Combining** our training dataset with other datasets dedicated to NER
 - **No significant improvements**

Phase 3: Analyze



**International Journal
(SJR: Q1)**

Social Network Analysis
and Mining (2024)



**International Conference
(CORE: B)**

International Symposium on
Intelligent Data Analysis (IDA
2023)



Young Researcher Forum

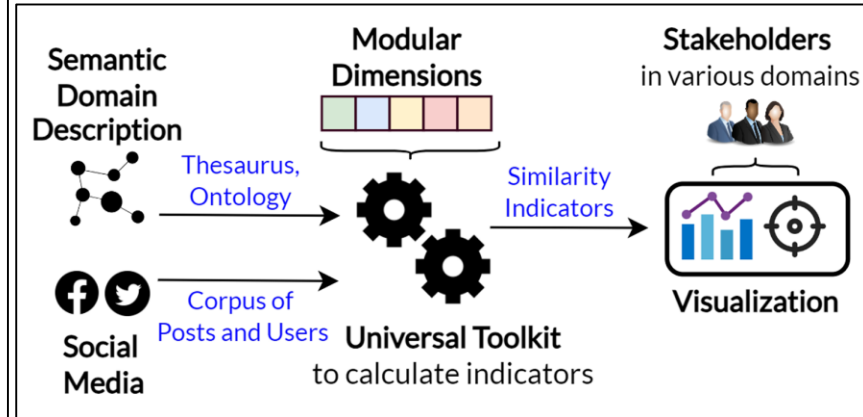
Young Researchers' Forum
at INFORSID 2022

Research Challenge

Social Media Indicators

❗ **Challenge:** Modeling social media entities and interactions in a **domain-agnostic manner** to produce **adaptable indicators** for decision support.

? **Hypothesis:** Redefining the **Proxemics theory** for use in social media and calculating indicators through **similarity measures** based on proxemic dimensions.



Related Work

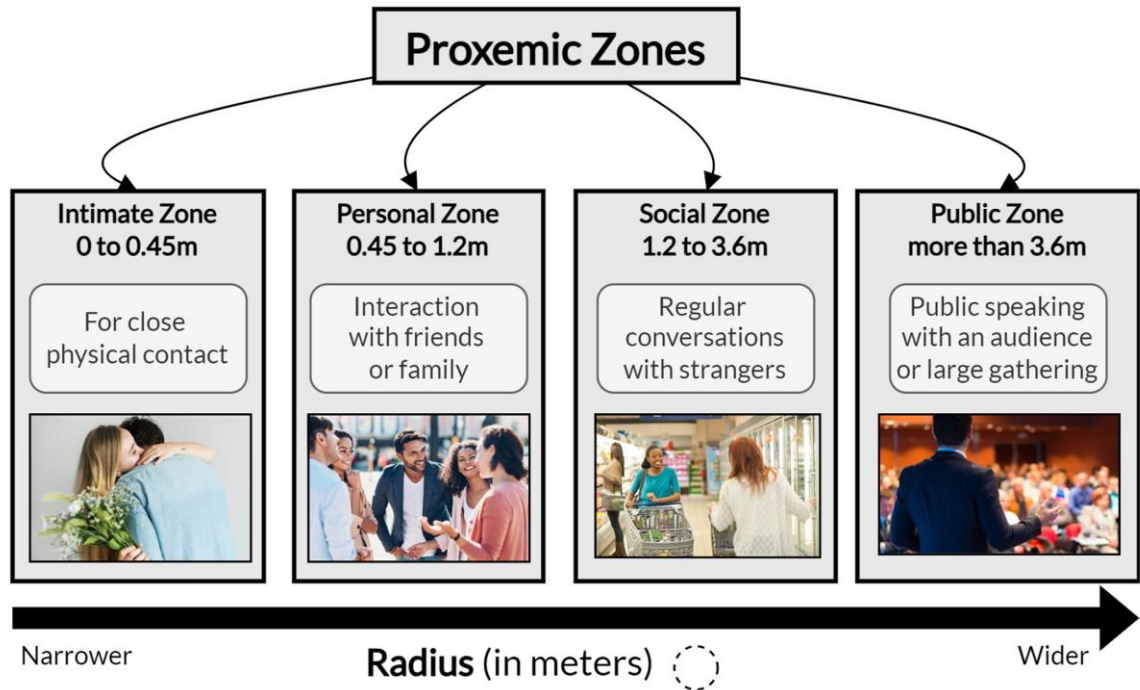
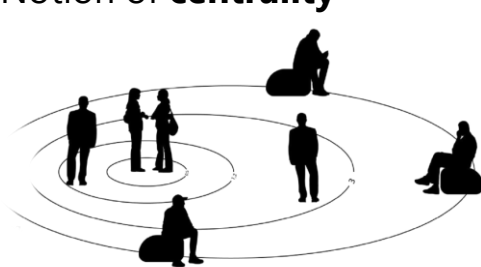
The Proxemics Theory (Hall, 1966)

i **Cultural, social, and physical** factors can affect the definition of proxemic zones

i **Two levels of analysis**

- o Individual
- o Group

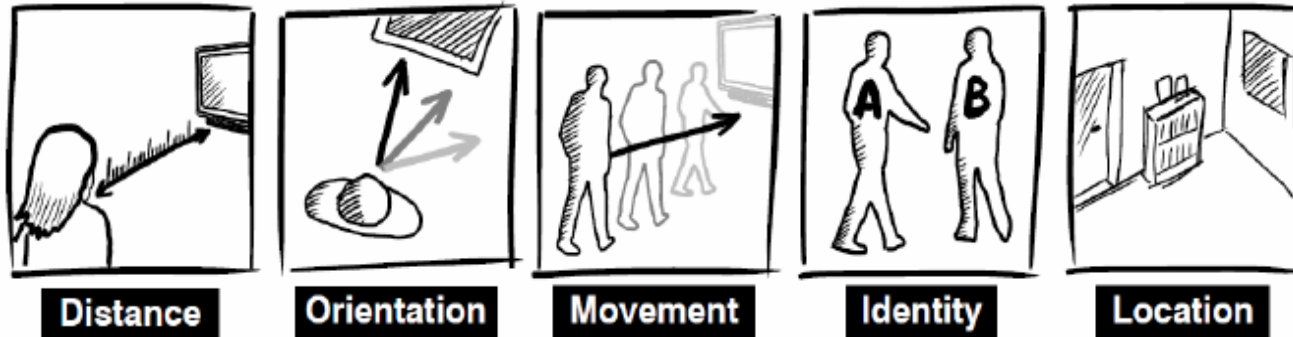
i **Notion of centrality**



Related Work

The 5 Proxemic Dimensions: DILMO (Greenberg, 2011)

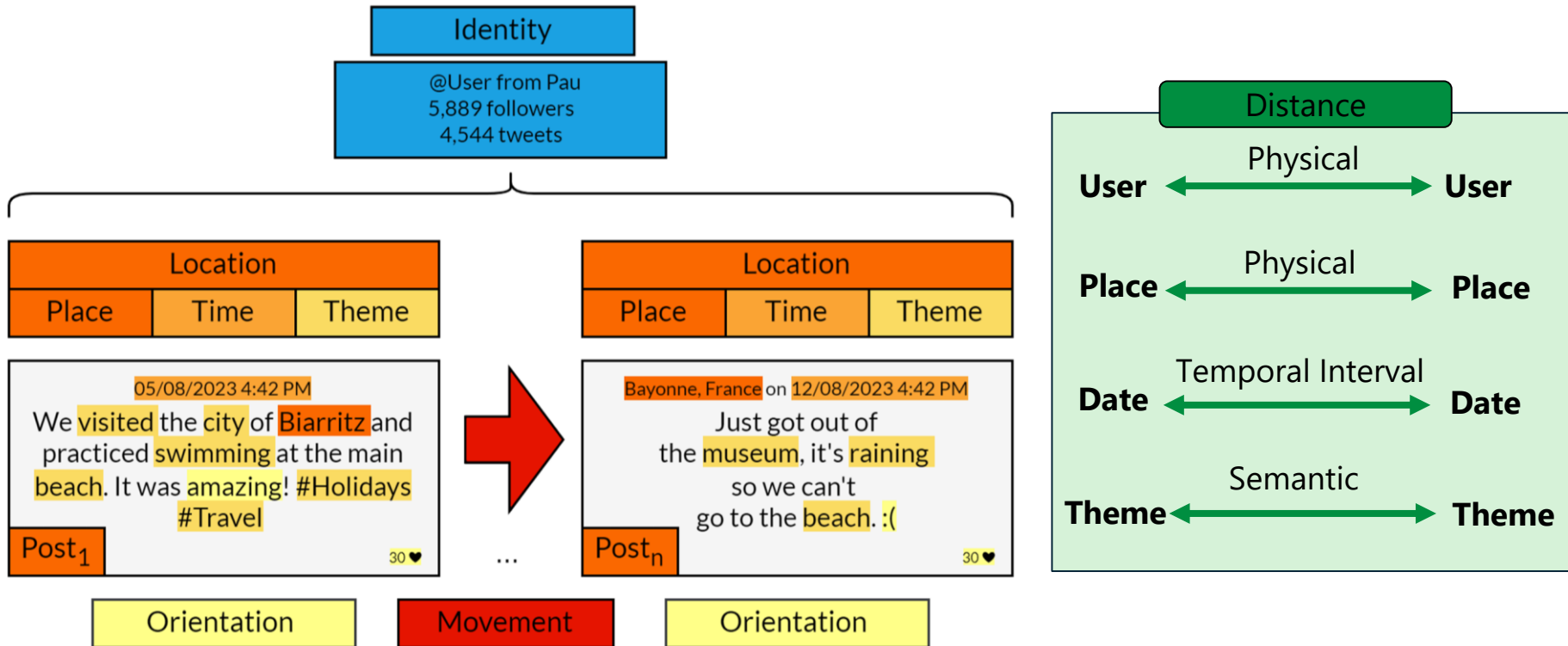
- i **DILMO Dimensions**
- i **Extension** of the theory of proxemics
- i Five dimensions used to **describe proxemic environments**



Source: Greenberg and Marquardt, 2011

Formal Redefinition of *Proxemics* in the Context of Social Media

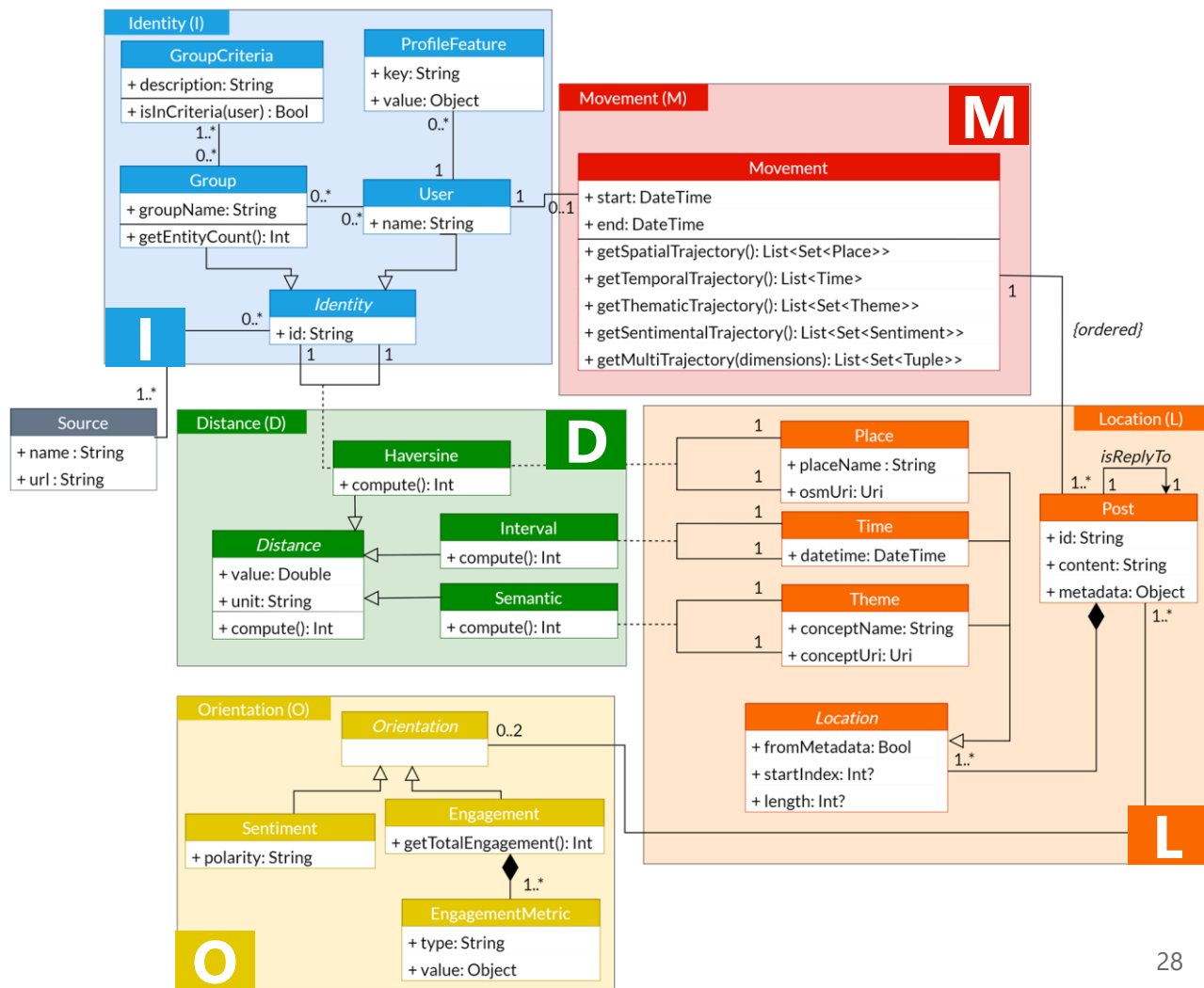
Adapting Proxemic Dimensions to Model Social Media Entities and Interactions



The APs Proxemic Model

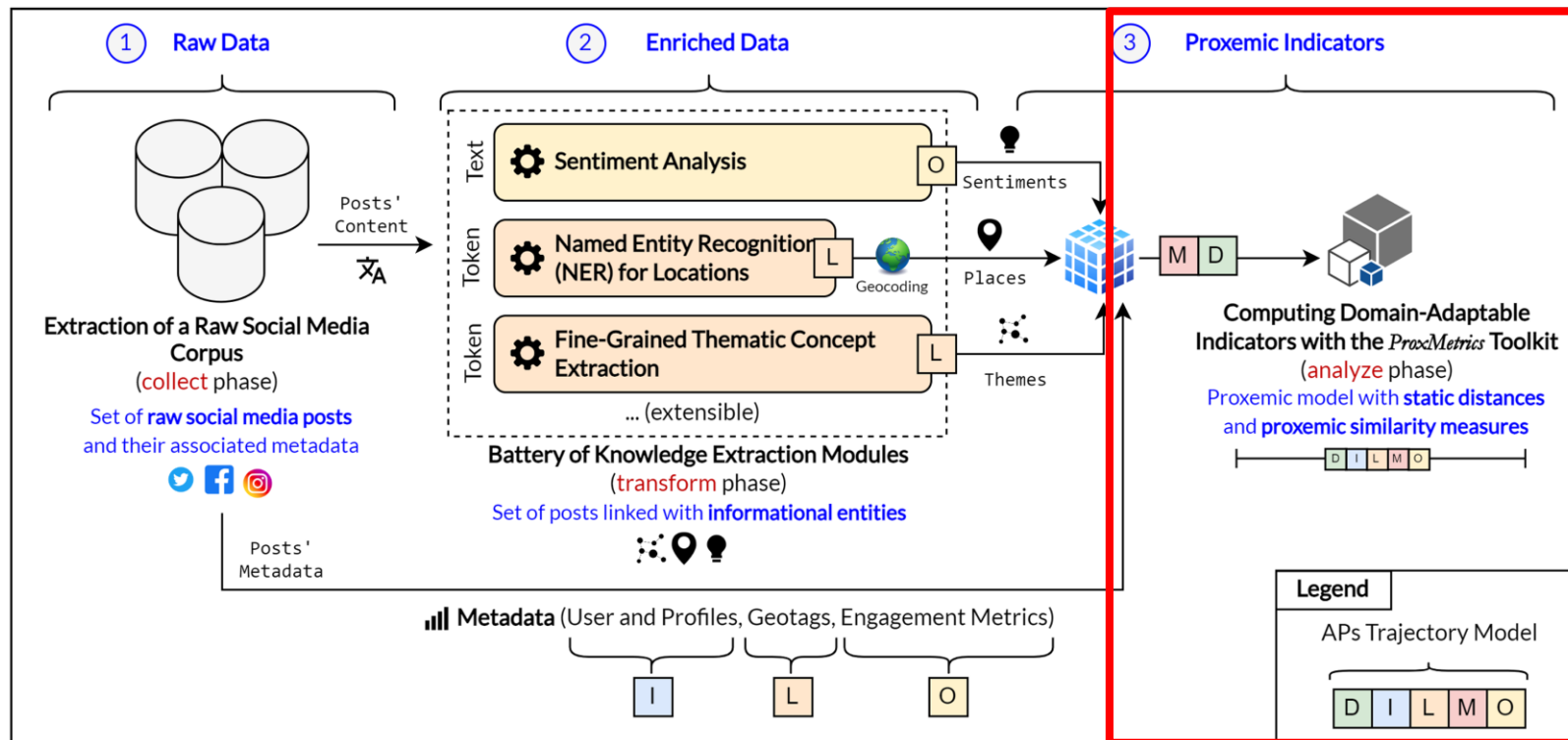
Data Model Overview

Distance	D
Identity	I
Location	L
Movement	M
Orientation	O



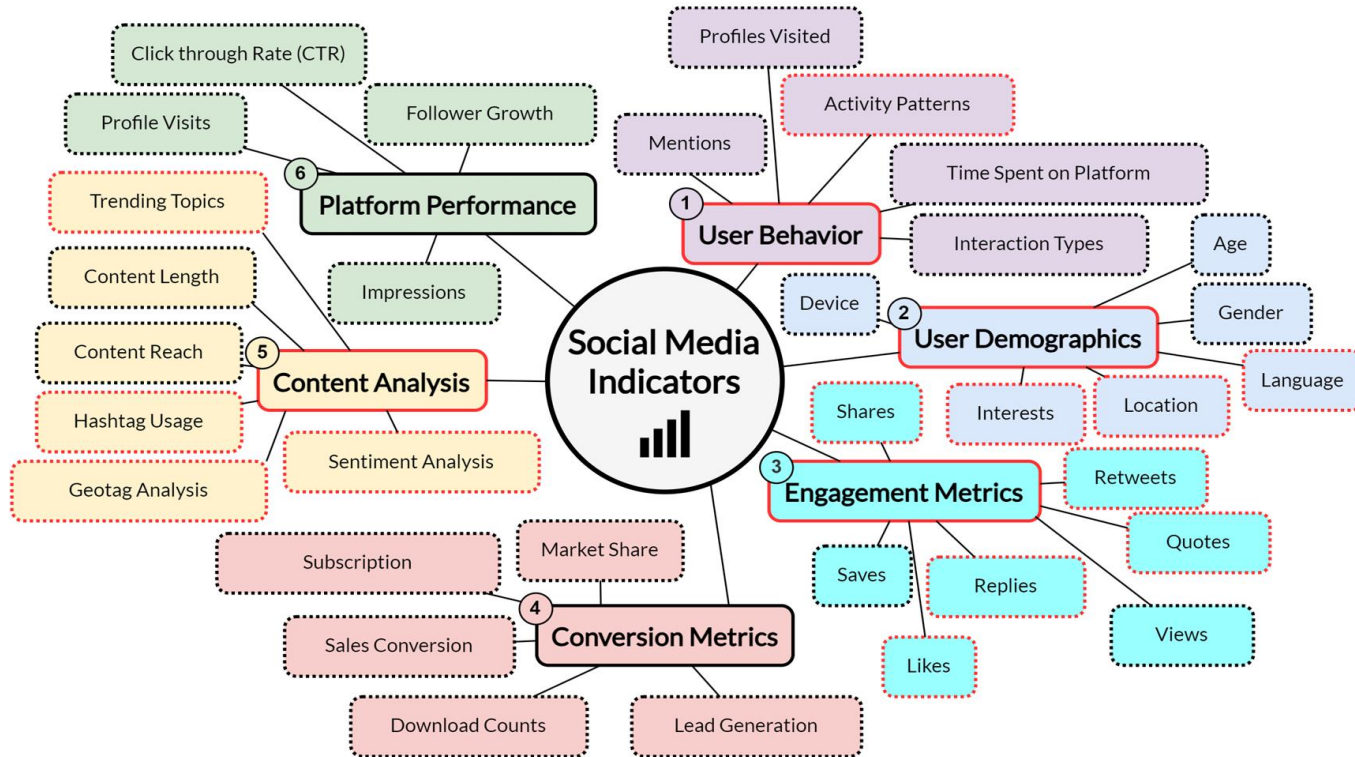
The APs Proxemic Model

Instantiation Process



Related Work

Social Media Indicators



ProxMetrics: Modular Toolkit to Evaluate Proxemic Similarity in Social Media

Social Media Entity Definition and Proxemic Similarity

① Dynamic Entities

Individual User



Group of Users



② Static Entities (Informational Entities)

Places



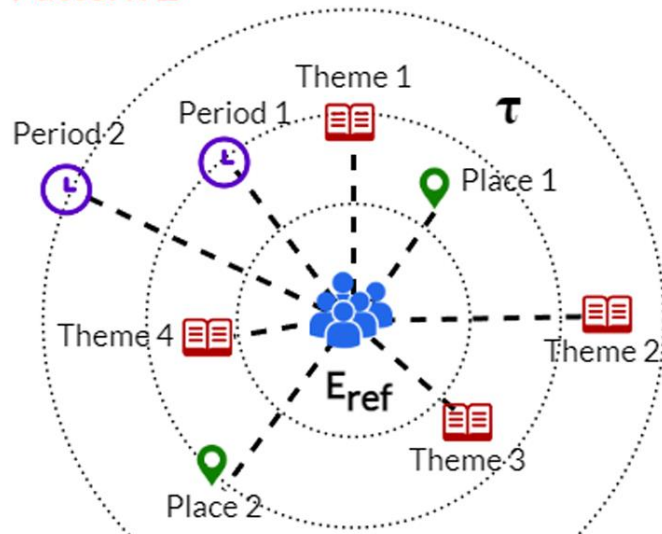
Periods



Themes



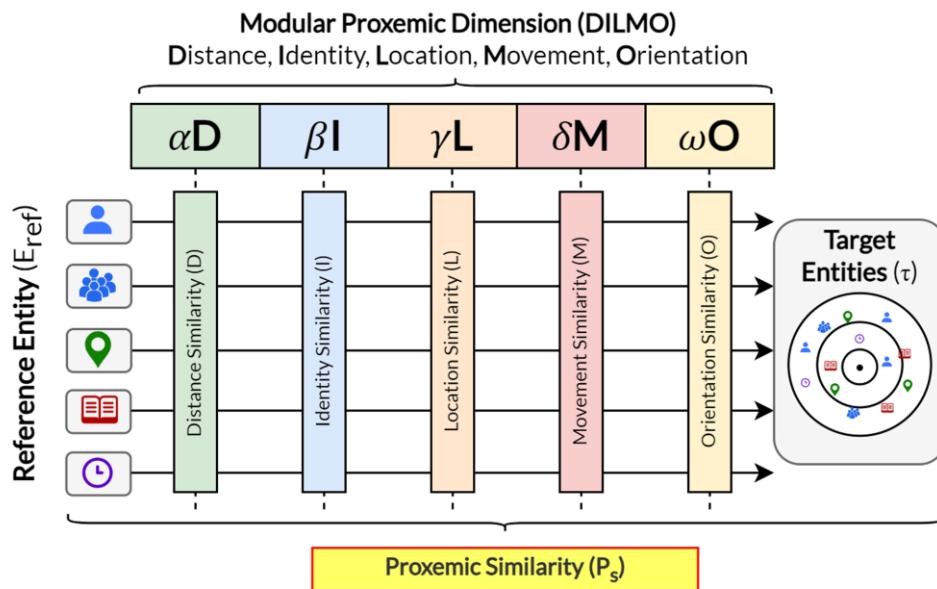
Pattern 2



$$E_{ref} \in \{\text{user, group}\}$$
$$\forall E_{target} \in \tau, E_{target} \in \{\text{place, time, theme}\}$$

ProxMetrics: Modular Toolkit to Evaluate Proxemic Similarity in Social Media





















Proxemic Similarity Definition



$$P_s(E_{ref}, E_{target}) = \alpha D(E_{ref}, E_{target}) + \beta I(E_{ref}, E_{target}) + \gamma L(E_{ref}, E_{target}) \\ + \delta M(E_{ref}, E_{target}) + \omega O(E_{ref}, E_{target}) \\ \text{with } \alpha + \beta + \gamma + \delta + \omega = 1$$

ProxMetrics: Modular Toolkit to Evaluate Proxemic Similarity in Social Media

Proxemic Similarity Definition







	E_{ref}	E_{target}	D	I	L	M	O
Pattern 1	 ∨ 	 ∨ 	$D_{physical}$	$I_{individual}$ I_{group}	$L_{individual}$	$M_{individual}$	$O_{individual}$
Pattern 2	 ∨ 	 ∨  ∨ 	n/a	I_{group}	$L_{occurrences}$	$M_{entropy}$	$O_{occurrences}$
Pattern 3	 ∨  ∨ 	 ∨ 	n/a	I_{group}	$L_{occurrences}$	$M_{entropy}$	$O_{occurrences}$
Pattern 4	 ∨  ∨ 	 ∨  ∨ 	$D_{physical}$ $D_{semantic}$ $D_{interval}$	I_{group}	$L_{co-occurrences}$	$M_{sequencing}$	$O_{co-occurrences}$

① Based on **existing formula**, adapted for social media

- Jaccard ($L_{individual}$)
- Conditional Probability ($M_{sequencing}$)
- Entropy ($M_{entropy}$)

Experimentation

Modelling Tourism Requirements with *ProxMetrics*

	Proxemic Environment		Dimensions				
Req.	Reference (E_{ref})	Targets (τ)	D	I	L	M	O
(A)	 Leisure Activity	 Leisure Activities			•		•
(B)	 Municipality <i>or</i>  POI	 Municipalities, POIs	•			•	
(C)	 User Group	 Municipalities		•	•		•

(A)

What leisure activities do tourists typically engage in together?

(B)

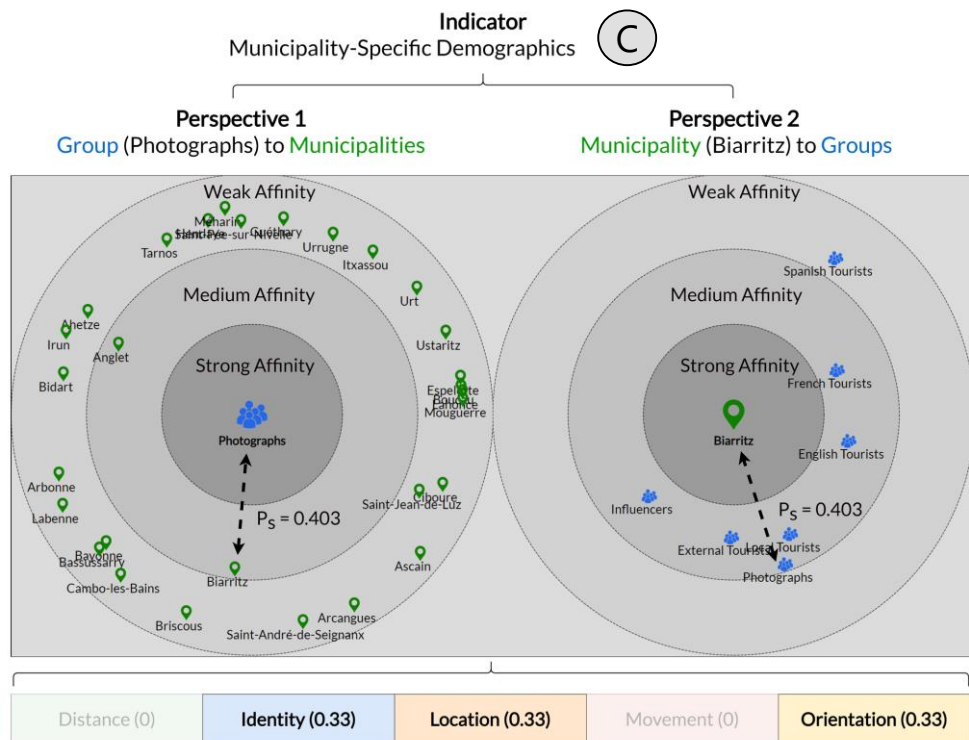
Which municipalities do tourists tend to go to after visiting Bayonne?

(C)

What are the typical demographics of tourists who visit Biarritz?

Experimentation

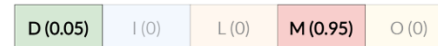
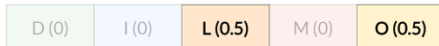
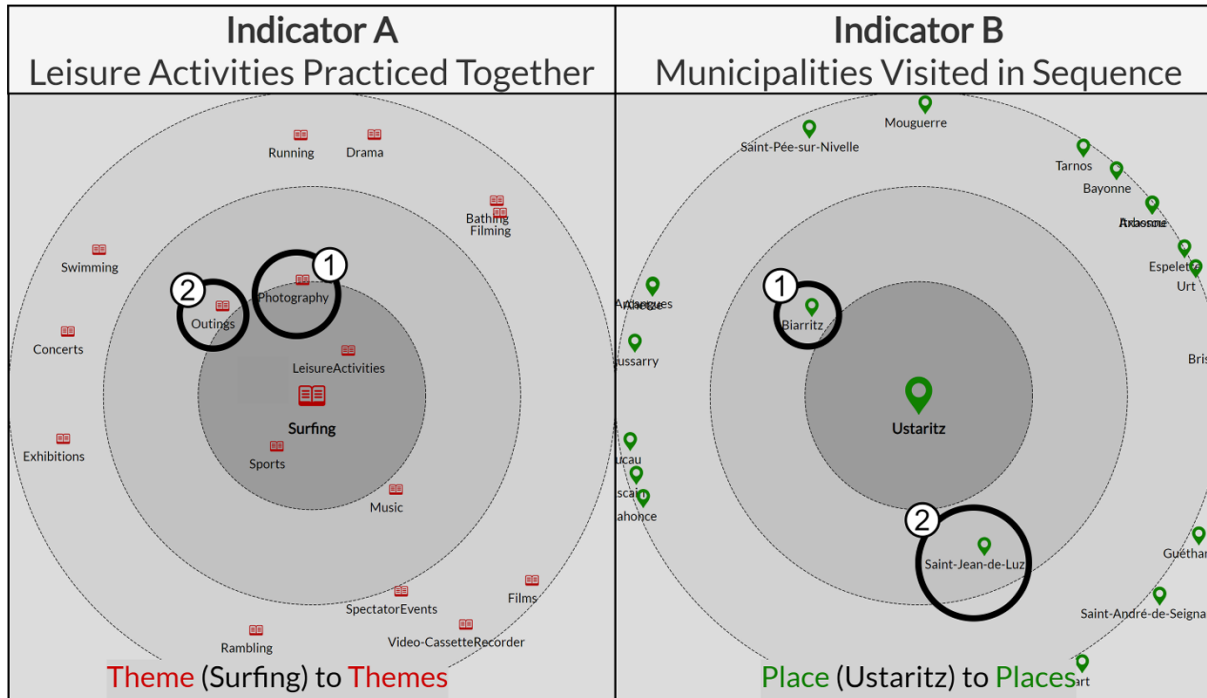
Case Study on Indicator C



$$\begin{aligned}
 P_s(\text{Photographers}, \text{Biarritz}) = & \\
 \frac{1}{3} \times & \underbrace{I_{\text{group}}(\text{Photographers}, \text{getUsersMentioning}(\text{Biarritz}))}_{0.52} + \\
 \frac{1}{3} \times & \underbrace{L_{\text{occurrences}}(\text{Photographers}, \text{Biarritz})}_{0.34} + \\
 \frac{1}{3} \times & \underbrace{O_{\text{occurrences}}(\text{Photographers}, \text{Biarritz})}_{0.36} \\
 & \approx 0.403
 \end{aligned}$$

Experimentation

Examples of Indicators: **A** and **B**

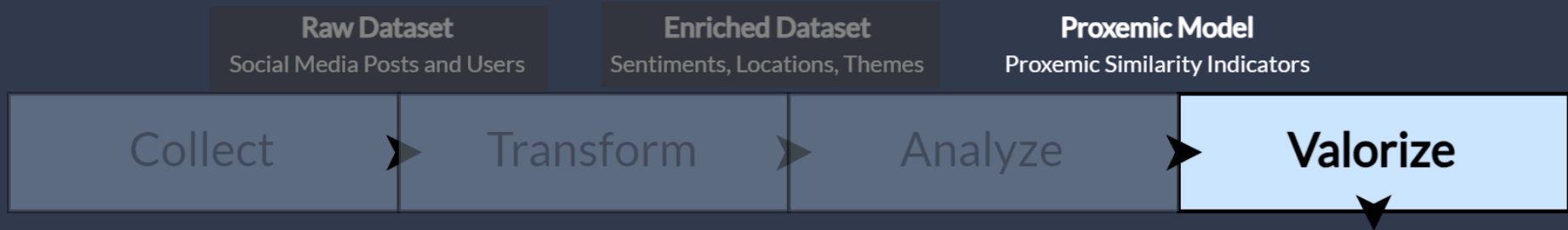


Evaluation

Protocol

Pattern 1 - Dynamic to Dynamic									
Indicator	<i>Connection of Similar Tourists</i>								
Prox. Environment	User (Dominique) to User (Luco)								
	Evaluators					σ	\bar{x}	<i>ProxMetrics</i>	Δ
Distance	5	8	8	5	8	1,47	6,80	8,60	1,80
Identity	3	3	3	2	4	0,63	3,00	6,50	3,50
Location	2	5	2	4	2	1,26	3,00	1,70	1,30
Movement	2	5	3	2	2	1,17	2,80	1,70	1,10
Orientation	1	5	2	2	3	1,36	2,60	2,80	0,20
Combination	DILMO								
	3	5	2	2	4	1,17	3,20	4,26	1,06

Phase 4: Valorize



Insights for End-Users
Non-Computer Scientists 



National Geonumeric Days (GeoDataDays)
Winner of the GeoData Challenge 2023



National Workshop
Workshop "Exploring Traces in an All-Digital World: Challenges and Perspectives" at INFORSID 2023



National Journal
Mappemonde



International Conference (CORE: A)
Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations (EACL 2024)

Research Challenge

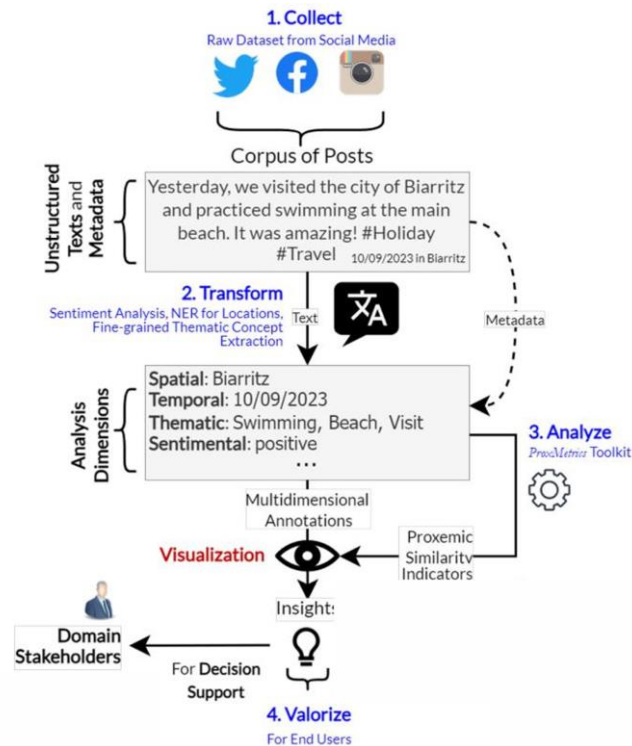
Visualization for Decision Support



Challenge: Presenting multidimensional social media analyses to non-computer scientists in a domain-adaptable manner.



Hypothesis: Extending, integrating, and blending selected features from existing visualization-based decision support tools could allow to build a dashboard addressing our requirements.



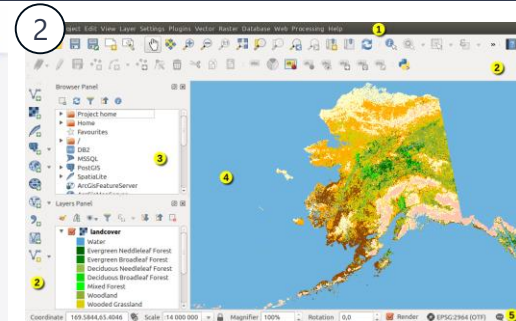
Related Work

Visualization for Decision Support

- 1 Domain-Specific Dashboards (DSD)
- 2 Geographic Information Systems (GIS)
- 3 Business Intelligence Tools (BI)
- 4 Linguistic Information Visualizations (LV)
- 5 Generic Visualization Libraries



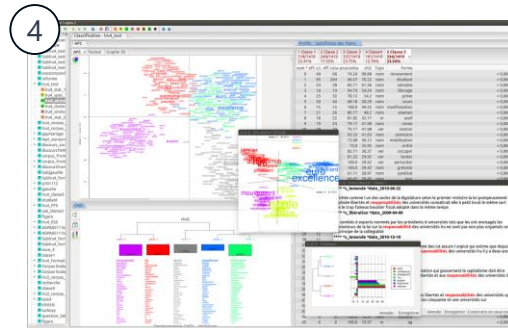
Source: INSEE



Source: QGIS



Source: PowerBI



Source: IRaMuTeQ

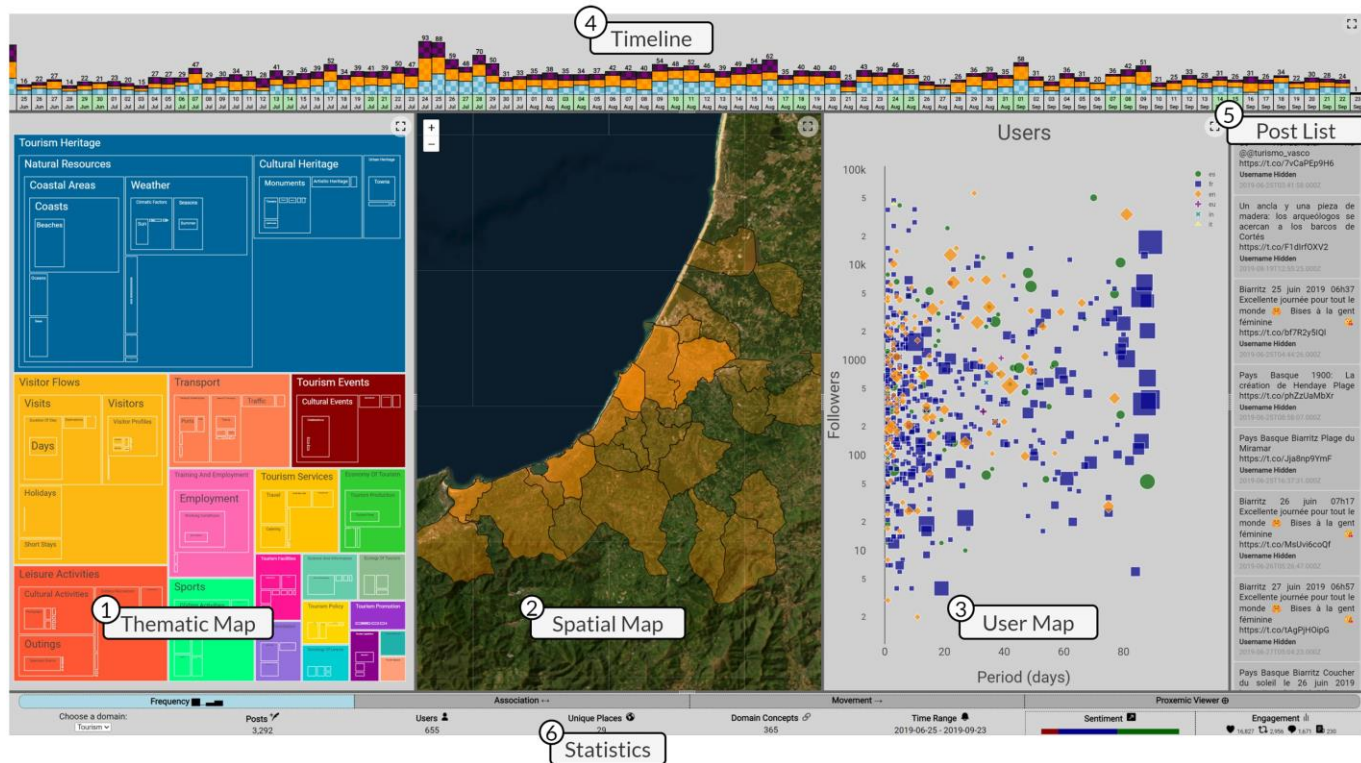


Source: D3JS

The *TextBI* Dashboard

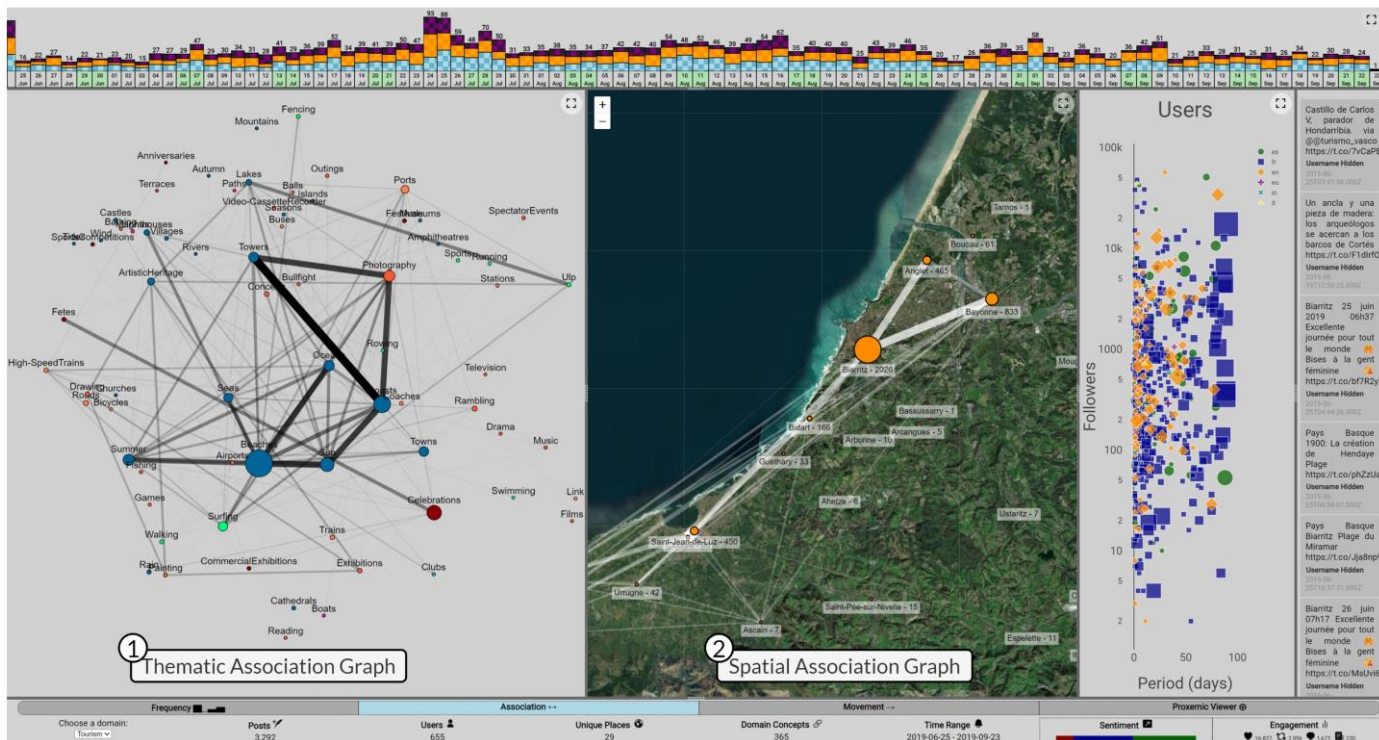
Frequency View

i **Demonstration video** available:
maxime-masson.github.io/TextBI



The *TextBI* Dashboard

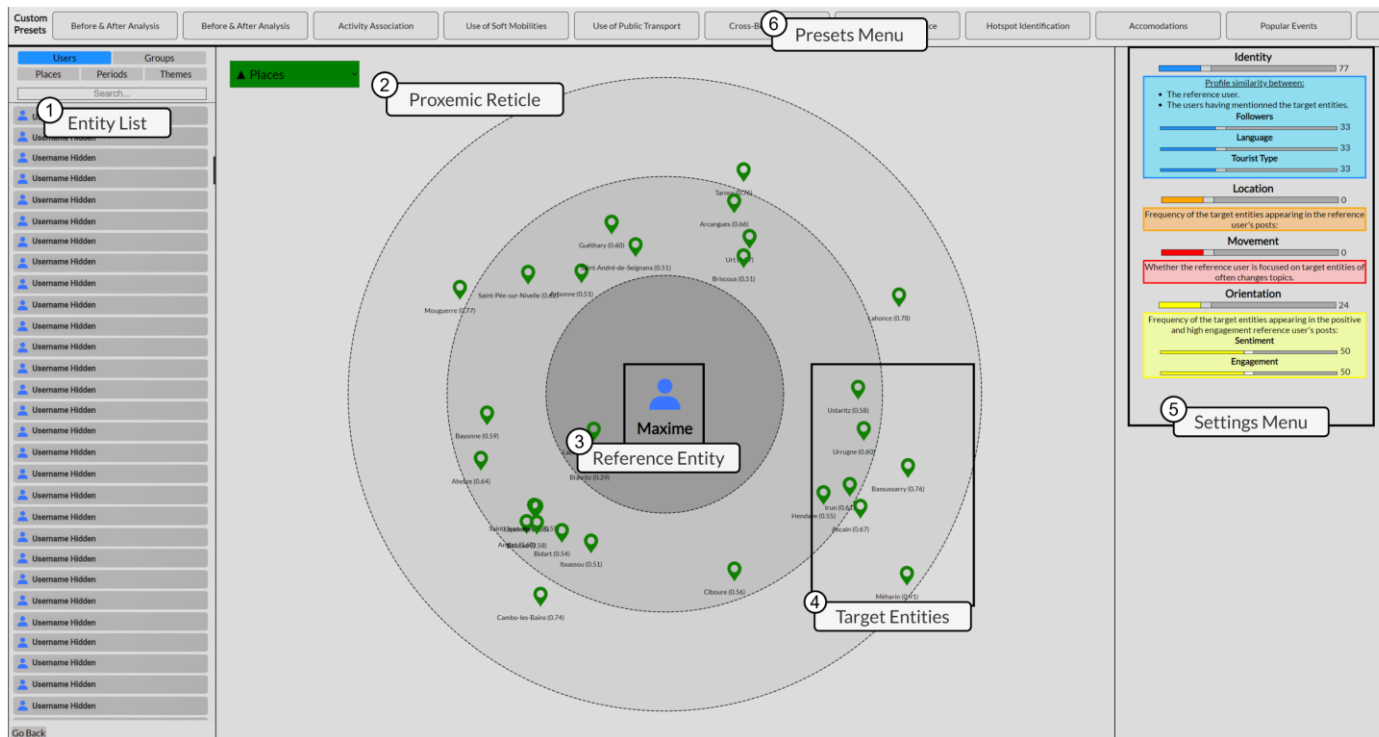
Association View



The *TextBI* Dashboard

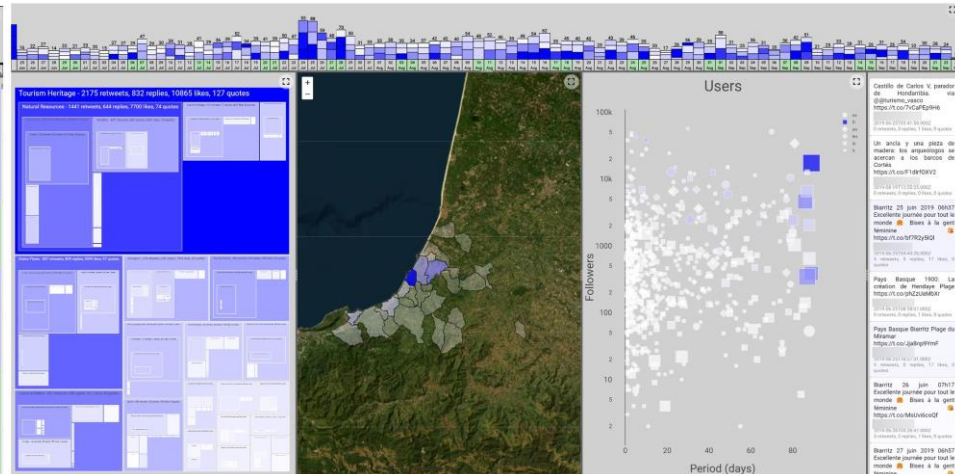
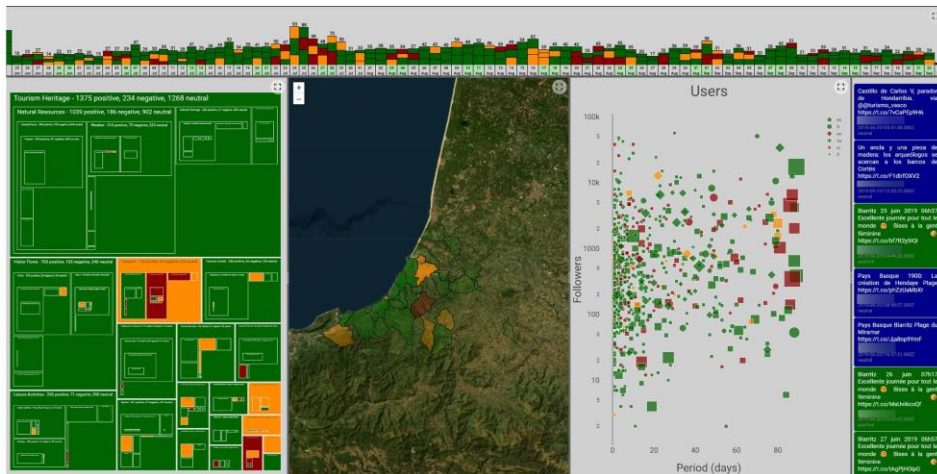
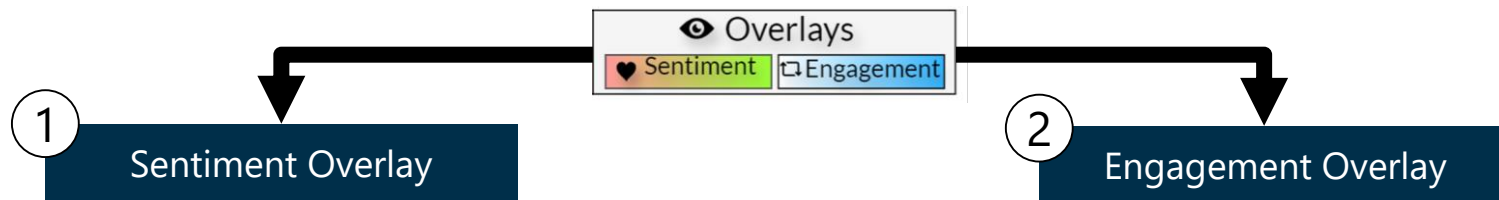
Proxemics View

i Powered by the **ProxMetrics** Toolkit



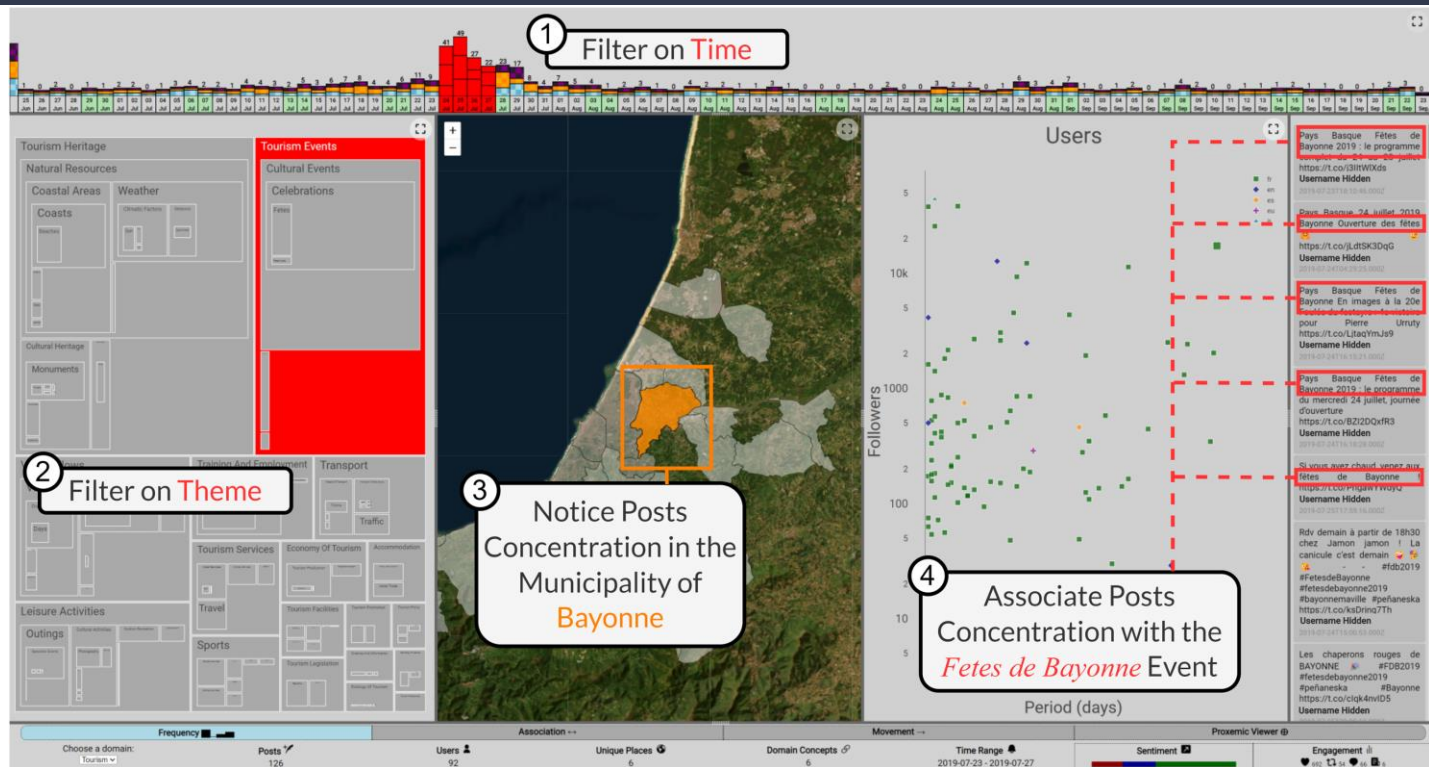
The *TextBI* Dashboard

Overlays



The *TextBI* Dashboard

Interactions

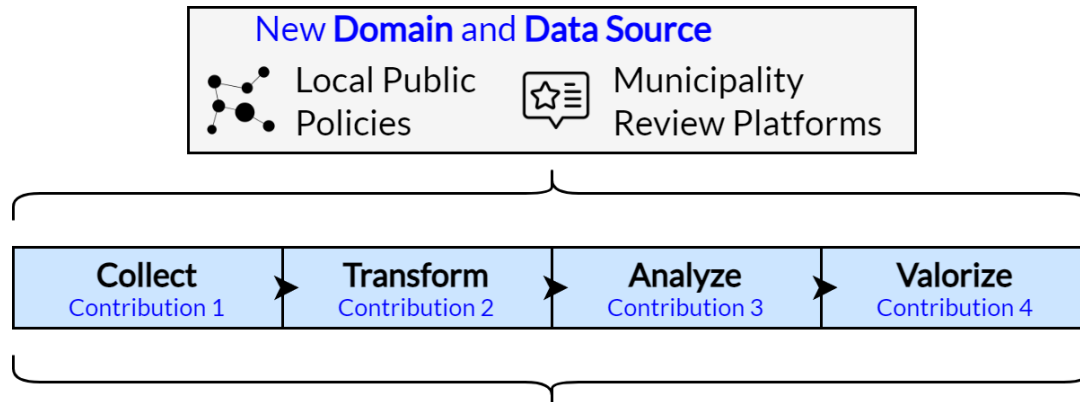


Conclusion

APs Framework

Contributions fitting within the **APs framework**

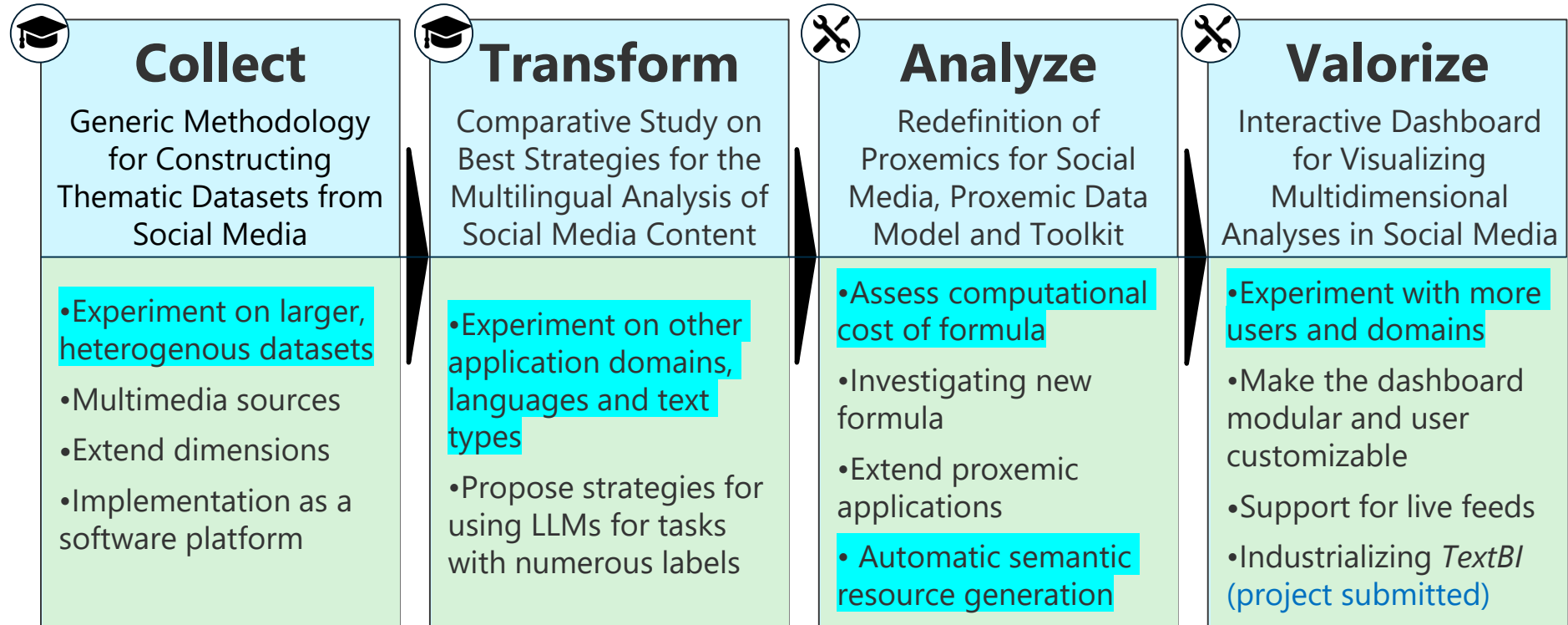
- Framework aiming to **assist in decision-making** based on **social media**
- **Generic:** Domain of application and social media source



Does the **APs Framework** generalizes well to this new **domain of application** and **data source**?

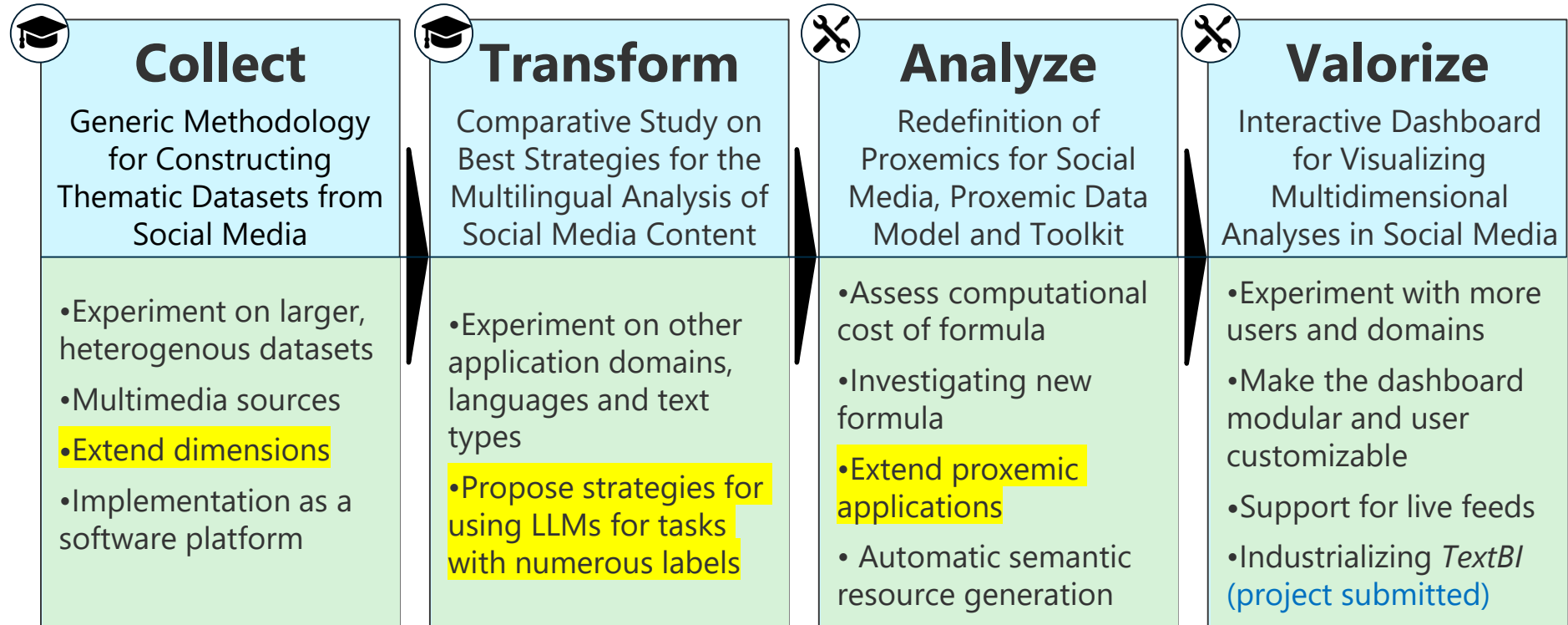
Conclusion

Contributions and Perspectives



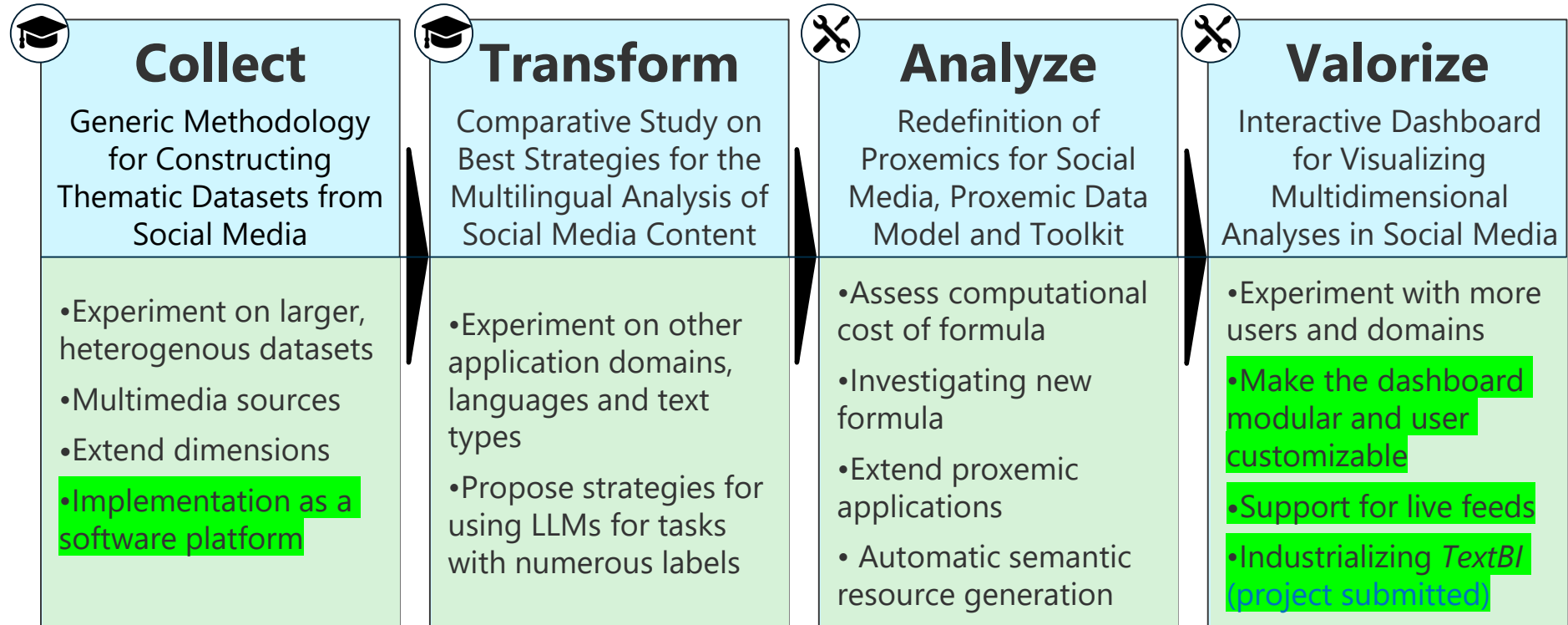
Conclusion

Contributions and Perspectives



Conclusion

Contributions and Perspectives



Thank you for your attention

Any questions?

maxime.masson@univ-pau.fr

[maxime-masson.github.io](https://github.com/maxime-masson)

International Publications



International Conferences

- [M. Masson](#), C. Sallaberry, M. N. Bessagnet, A. Le Parc Lacayrelle, P. Roose, R. Agerri. (2024). TextBI: An Interactive Dashboard for Visualizing Multidimensional NLP Annotations in Social Media Data . In **Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations (EACL 2024)**. (pp. 1-9). Association for Computational Linguistics (ACL).
- [M. Masson](#), P. Roose, C. Sallaberry, R. Agerri, M. N. Bessagnet, A. Le Parc Lacayrelle. (2023). APs: A Proxemic Framework for Social Media Interactions Modeling and Analysis . In **International Symposium on Intelligent Data Analysis (IDA 2023)**. (pp. 287-299). Cham: Springer Nature Switzerland.
- [M. Masson](#), C. Sallaberry, R. Agerri, M. N. Bessagnet, P. Roose, A. Le Parc Lacayrelle. (2022). A Domain-independent Method for Thematic Dataset Building from Social Media: The Case of Tourism on Twitter . In **International Conference on Web Information Systems Engineering (WISE 2022)**. (pp. 11-20). Cham: Springer International Publishing.



International Journals

- [M. Masson](#), P. Roose, C. Sallaberry, M. N. Bessagnet, A. Le Parc Lacayrelle, R. Agerri. (2024). ProxMetrics: Modular Proxemic Similarity Toolkit to Generate Domain-Adaptable Indicators from Social Media . In **Social Network Analysis and Mining (SNAM)**. 14, 124. Springer.
- [M. Masson](#), R. Agerri, C. Sallaberry, M. N. Bessagnet, A. Le Parc Lacayrelle, P. Roose. (2023). Optimal Strategies to Perform Multilingual Analysis of Social Content for a Novel Dataset in the Tourism Domain. Submitted to **Knowledge-Based Systems (KNOSYS)** journal.

National Publications



National Conference

- M. Masson, R. Agerri, C. Sallaberry, M. N. Bessagnet, P. Roose, A. Le Parc Lacayrelle. (2024). Stratégies optimales pour l'analyse multidimensionnelle de contenus multilingues issus des réseaux sociaux . In **Proceedings of the 42nd Conference on Computer Science for Organizations and Information and Decision Systems (INFORSID 2024)**.



National Journal

- M. Masson, C. Sallaberry, M. N. Bessagnet, A. Le Parc Lacayrelle, P. Roose, R. Agerri. (2024). Visualisation de données issues des réseaux sociaux : une plateforme de type Business intelligence . In **Mappemonde**, OpenEdition Journals.



National Workshops

- M. Masson, S. Abdelhedi, C. Sallaberry, R. Agerri, M. N. Bessagnet, A. Le Parc-Lacayrelle, P. Roose. (2023). Visualisation interactive de trajectoires d'activités touristiques: application à des données extraites de twitter. In **Workshop "Exploring traces in an all-digital world: challenges and perspectives" at INFORSID 2023**.
- M. Masson. (2022). Services augmentés pour le tourisme intelligent et l'analyse des pratiques. In **Young Researchers' Forum at INFORSID 2022**.

Awards and Communications



Award

- Ranked 1st at the **Geodata Challenge** of the **National Geonumeric Days 2023 (GeoDataDays 2023)** with the proposal “Visualization of data from social media: a Business intelligence type platform”. This event was organized by the French Association for Geographic Information (Afigéo).



Communications

- [M. Masson](#). (January, 2024). TextBI: An Interactive Platform for Visualizing Multidimensional Data from Social Media. Keynote Speaker: **Webinar on Cartography and Geovisualization of the GdR CNRS MAGIS (CNRS Research Network on Methods and Applications for Geomatics and Spatial Information)** (Online).
- [M. Masson](#). (November, 2023). TextBI: A Generic Dashboard for Interactive Visualization of Multidimensional Data from Social Media. **Workshop on Spatialized Digital Humanities, Annual Meeting of the GdR CNRS MAGIS (CNRS Research Network on Methods and Applications for Geomatics and Spatial Information)**. Maison des Suds (Bordeaux, France).
- [M. Masson](#), P. Roose. (July, 2023). Analyzing Touristic Data in the Basque Country. **Urban community of the Basque Country** (Bayonne, France).
- [M. Masson](#). (June, 2023). A Generic Framework for the Extraction, Processing, Analysis, and Valuation of Social Media content: Application to the Domain of Tourism and the Social Media Twitter. **Ixa Seminar, University of the Basque Country** (EHU/UPV) (San Sebastian, Spain).
- [M. Masson](#), S. Laborie. (June, 2023). A Generic Framework for the Extraction, Processing, Analysis and, Valorization of Social Media Content. **Symposium "Constitution of corpus for the needs of digital marketing in the domain of fashion" (European Cassini Program), Parthenope University of Naples** (Naples, Italy).
- [M. Masson](#). (November, 2022). APs: A Proxemic Approach for Data Analysis on Social Media. **Workshop "Smart city, smart destination: from management to territorial experience", IRGO - Research Institute in Organizational Management, University of Bordeaux** (Bordeaux, France).
- [M. Masson](#). (September, 2022). APs: A Proxemic Approach for Data Analysis on Social Media. **Inter-association Day EGC/INFORSID, IRIT, University of Toulouse III** (Toulouse, France).