

From LLMs to information cascades

Annotate, measure and detect media dynamics at scale

Antoine Lemor, Ph.D.

Postdoctoral researcher

Université de Sherbrooke · RFICS · CIRST · CAPP



Outline

1

Text analysis in social sciences

material, extraction, classification

2

LLM Tool

an end-to-end solution

3

Application – CCF

media cascades on climate change

PART I

Text analysis

Material, extraction, classification

Text as foundational material

MATERIAL

Written

Documents, archives, manuscripts

Oral

Speeches, interviews, conversations

Visual

Images, graphics, videos

Text

malleable

transformable material

EXTRACTION

Information

Facts, data, knowledge

Meaning

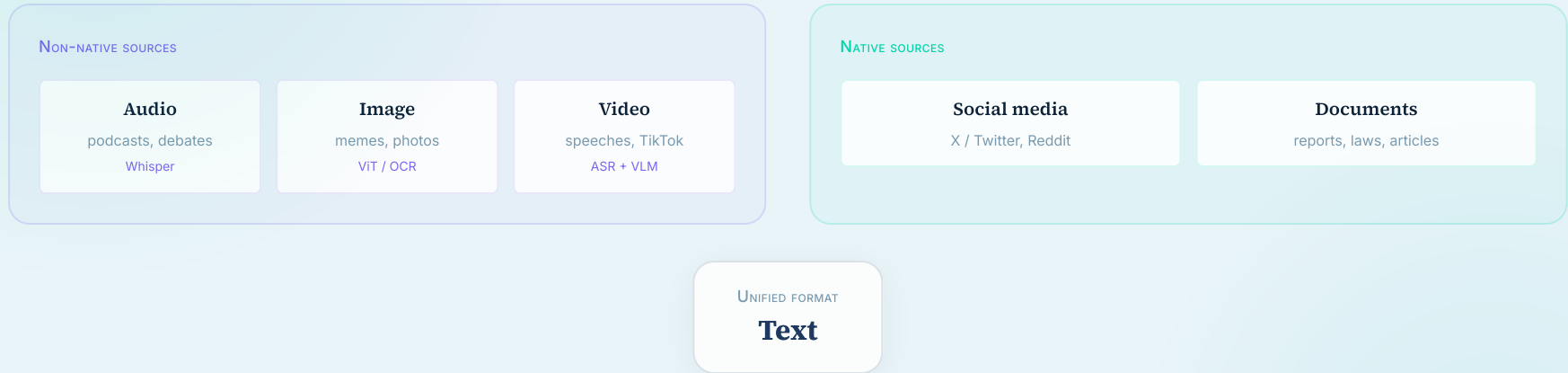
Significations, interpretations

Structures

Recurrences, relations, trends

Text is a malleable material that both humans and machines can read

What is a text in social sciences?



Text unifies native and transformed sources into structured data, ready for computational analysis.

Why classify, and what to extract?

CONTENT

Themes, topics, concepts

Climate change, economy, public health

POSITION

Stance, for or against

Support, criticism, neutrality toward a policy

EMOTION

Sentiment, expressed affect

Anger, hope, fear, joy, sadness

ENTITIES

Actors, places, organizations

UN, Canada, Guterres, WHO

Text analysis turns unstructured language into quantitative, analyzable data

LLM Tool: challenges and answers

CHALLENGE

Scale

Thousands of documents, limited budget and time

Complexity

Different languages, text types, technical skills

Quality

Inter-annotator agreement, reproducibility, validation

Ethics

Sensitive data, commercial APIs, data protection

RESPONSE

Batch processing

High-speed annotation, free and open source

50+ models, 100+ languages

Multilingual, adaptable, LLM consensus

Validation Lab

State-of-the-art metrics, full traceability

100% local

Ollama, no data leaves your machine

LLM Tool: a complete, accessible answer for social science research (quantitative and qualitative)

PART II

LLM Tool

An integrated solution for automated text analysis

What is LLM Tool?

OPEN-SOURCE PLATFORM

Automatic text annotation powered by AI

Turns any text corpus into structured data along your own categories — without writing a line of code.

Accessible

Guided interface, zero programming

Flexible

10 to 10M documents, quanti and quali

Secure

100% local, no external API

Rigorous

Traceability, reproducibility

METHODOLOGICAL VALIDATION

Macro F1 = 66.73

Lemor, Dinan, Gilbert (2025)

CONTINUOUS EVOLUTION

Image, audio (soon)

ViT, Whisper integration

How does LLM Tool work?



• PARAMETERS TRACED FROM START TO FINISH •

A complete protocol: from raw data to a validated classification

Mode 1 — The Annotator



The Annotator: guided interface for LLM-based classification

MAIN FUNCTION

Uses language models (OpenAI, Ollama) to automatically classify your texts along your category schema.

CAPABILITIES

- Advanced prompt engineering features
- Handles diverse text types
- Runs 100% locally with Ollama (free)
- Auto-repairs format errors
- Exports to CSV, Excel, review platforms

IDEAL USE

Quickly annotate a small corpus, or build training data.

Mode 2 — Annotator Factory

MAIN FUNCTION

Chains annotation, data preparation, model training and deployment in a single automated process.

AUTOMATED STEPS

- Annotates a sample with an LLM
- Prepares data (balancing, cleaning)
- Splits train / validation / test
- Selects and deploys the best model
- Applies the model(s) to your whole corpus

IDEAL USE

A custom model for a large corpus in a few hours.



Annotator Factory: end-to-end model creation

Mode 3 — Training Arena



Training Arena: trains and compares model families

MAIN FUNCTION

Trains up to 50 different model types and automatically picks the best one for your data.

CAPABILITIES

- BERT, RoBERTa, DeBERTa, CamemBERT (FR), XLM-R
- Auto-detects the language of your texts
- Compares performance with cross-validation
- Keeps only the best model
- GPU-optimized

IDEAL USE

Get the optimal model for your training corpus.

Mode 4 — BERT Annotation Studio

MAIN FUNCTION

Applies your trained BERT models to very large document collections at high throughput.

PERFORMANCE

- High-speed GPU annotation
- Auto-optimizes available memory
- Confidence score per prediction
- Handles millions of documents
- Auto-resumes if interrupted

IDEAL USE

Apply your model to large volumes of text.



BERT Annotation Studio: applies trained models at scale

Mode 5 — Validation Lab



Validation Lab: agreement metrics and quality control

MAIN FUNCTION

Systematically evaluates annotation quality by computing inter-annotator agreement and surfacing problematic cases.

VALIDATION TOOLS

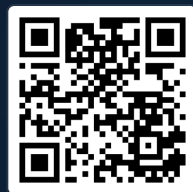
- Cohen's Kappa, Krippendorff's Alpha
- Flags unreliable annotations (score < 0.5)
- Identifies problematic categories
- Stratified samples for human review
- Detailed reports

IDEAL USE

Ensure scientific quality before publication or final analysis.

Demonstration

LLM Tool in action



llm-tool.com

Why use LLM Tool for your research?

ACCESSIBILITY

- **Zero programming** — guided interactive interface
- **Free option** — Ollama, fully local
- **Multilingual** — 100+ languages
- **Multi-platform** — macOS, Linux, Windows

SCIENTIFIC RIGOR

- **Reproducibility** — full traceability
- **Validation** — inter-annotator agreement metrics
- **Transparency** — full parameter export
- **Open source** — public code and docs

Built for social science researchers, by social science researchers

PART III

Application — CCF

Canadian Climate Framing — media cascades

The CCF database: Canadian Climate Framing

266K

Articles

1978 – 2024

20

Newspapers

EN + FR

9.2M

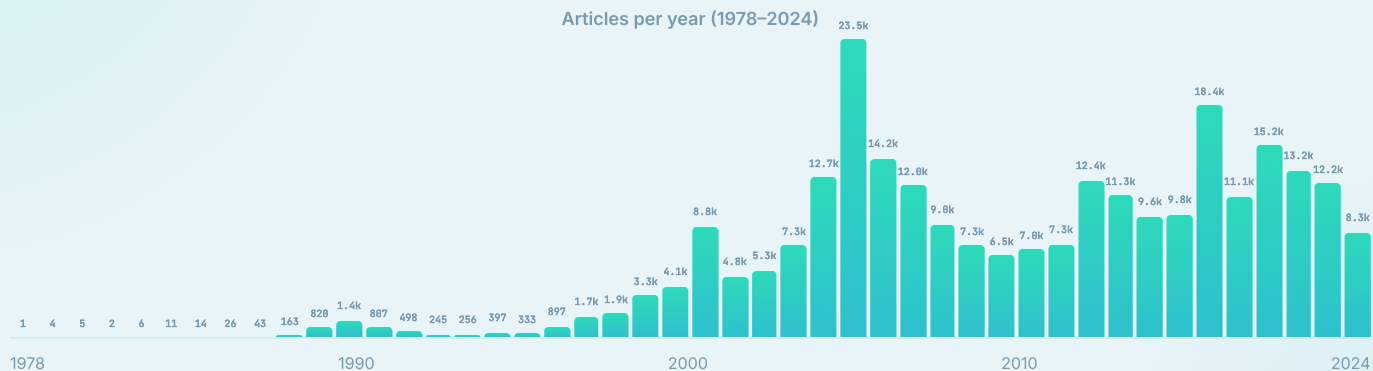
Annotated sentences

BERT / CamemBERT

65

Categories

hierarchical

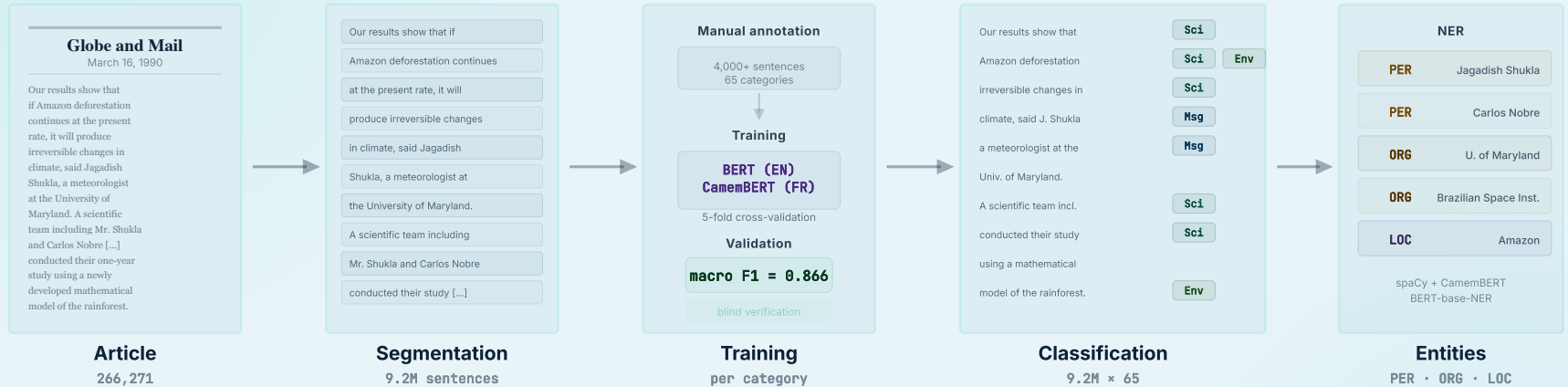


The largest ML-annotated climate discourse corpus in Canada

Macro F1 : 0.866 – Lemor, Pillod, Taylor (2025)

Sentence-level annotation: the pipeline

Each sentence of each article is classified by trained models



67 hierarchical categories

5 thematic frames + complementary dimensions



0.866

MACRO F1

67 BERT / CamemBERT models

1 per category · 63 available

Thematic frames

8 frames → 40 sub-categories

● Economy

Negative impacts
Positive impacts
Costs of action
Benefits of action
Sectoral footprint

● Health

Negative impacts
Positive impacts
Co-benefits
Sector footprint

● Security

Climate refugees
Resource conflicts
Military assistance
Military disruption

● Justice

Winners/losers
Responsibility
Vulnerability
Unequal access
Intergenerational

● Politics

Political action
Political debate
Positioning
Public opinion

● Science

Debate
Discovery
Skepticism
Science advocacy

● Environment

Habitat loss
Species loss

● Culture

Art
Events
Indigenous practices
Cultural footprint

Complementary dimensions

4 dimensions → 27 sub-categories

● Messengers 10

Health
Economy
Security
Legal
Cultural
Nat. sci.
Soc. sci.
Activist
Public official

● Events 9

Extreme weather
Meeting
Publication
Election
Policy
Judicial
Cultural
Protest

● Solutions 3

Mitigation
Adaptation

● Tone + Other 5

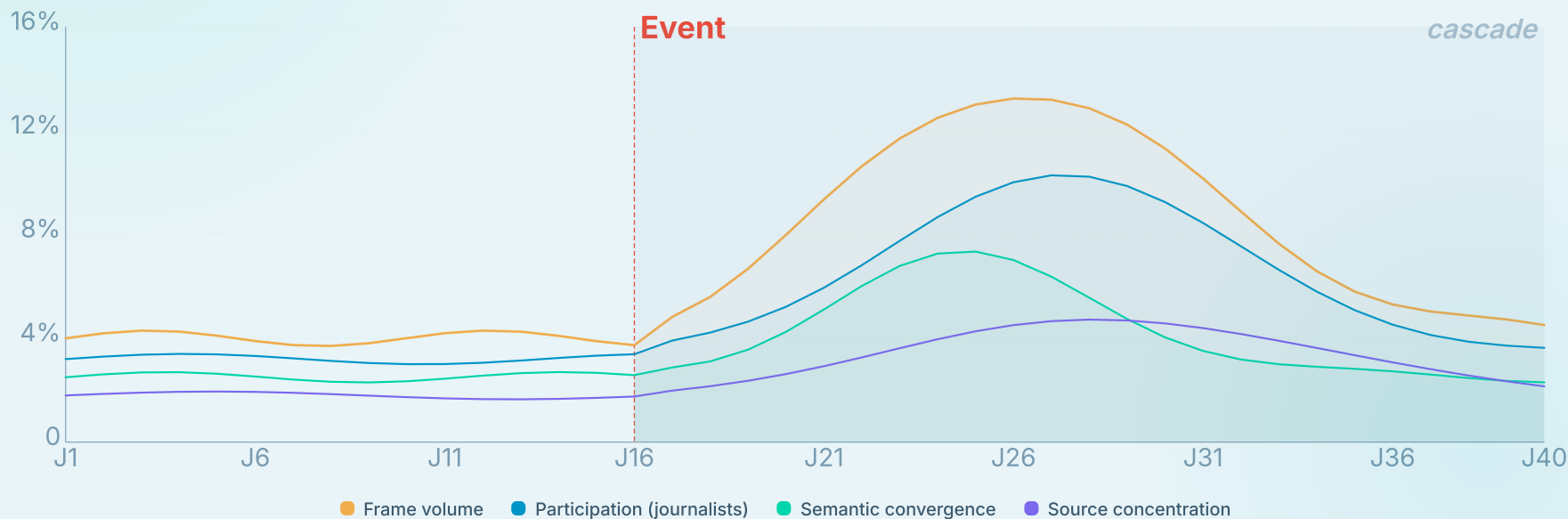
Positive
Negative
Neutral
Canada
Urgency

What is a media cascade?

A frame propagating massively across outlets after a triggering event

Media cascade =

a topic that spreads rapidly across media and journalists, with **convergence** of content, an **increase** in coverage volume, and the arrival of **new actors** in the debate



How to detect and measure a cascade?

Five signals, statistical detection, multi-dimensional score

Detecting cascades

Five signals tracked daily

- **Temporal** *Is the frame rising?*
- **Participation** *New journalists writing?*
- **Convergence** *Is the frame dominant?*
- **Source** *Same messengers across outlets?*
- **Semantic** *Articles look alike?*

Cascade detected

When the signals jointly exceed the background noise

- **Strong**
- **Moderate**
- **Weak**

Detecting events

Group articles into clusters

- **Annotations** *What kind of event? (weather, policy, publication...)*
- **Embeddings** *Are the articles semantically similar?*
- **Temporality** *Published at the same moment?*
- **Entities** *Same people, places, organizations?*

Cluster = event

Each cluster groups articles covering the same concrete event

Health cascade — Summer 2023

— Articles / day — Z-score — % Health • Clusters

The CCF platform

1

Explorer

Map of Canada, temporal trends, media profiles

2

Search

Semantic (embeddings), named entities, cross-tabulations

3

Analysis

Media cascades, event clusters, causal impact

PUBLIC API

data.ccf-project.ca

50+ REST endpoints · bilingual FR/EN · open access

Thank you!

Questions and discussion

Antoine Lemor, Ph.D.

Université de Sherbrooke · RFICS · CIRST · CAPP



antoinelemor.com



llm-tool.com



data.ccf-project.ca