

# *Characterizing healthcare workers vaccine narratives on X (Twitter).*

Leonardo Heyerdahl – Postdoctoral researcher,  
Anthropology and Ecology of Disease Emergence Unit, Global Health  
Department, Institut Pasteur

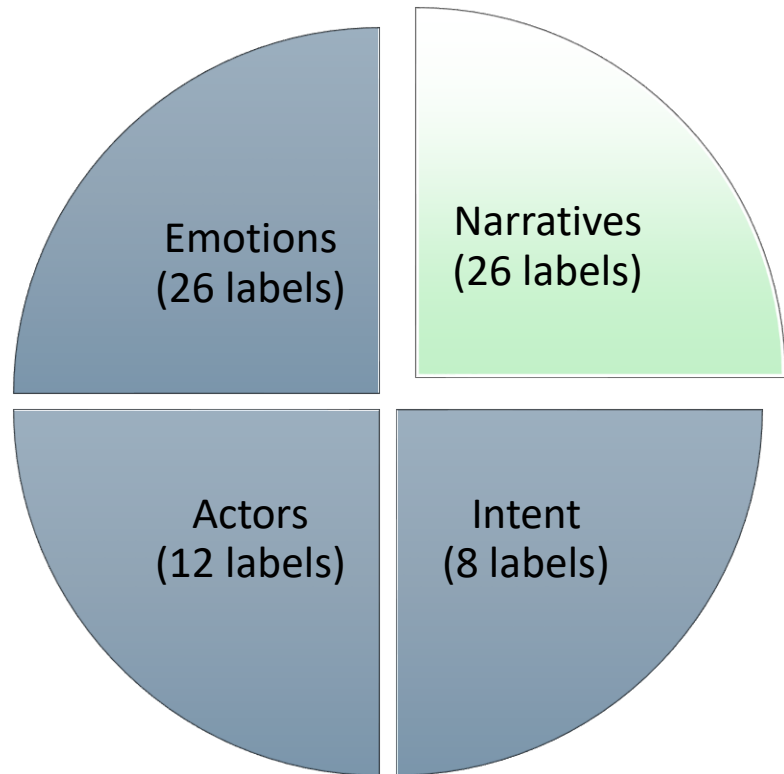
Les professionnels de santé, acteurs et sujets  
de la vaccination

1<sup>ÈRE</sup> JOURNÉE D'ÉTUDES DU RÉSEAU  
SHS-VACCINATION

PARIS 24 JANVIER 2025

# Objective

- To characterize online vaccine discourses among healthcare workers at scale using custom taxonomies.



# Description of the AEE Twitter database

**Daily collection of tweets** (COVID, vaccination, drugs, masks, etc.)

- From February 2020 – present
- Vaccination: >12 million organic tweets

## **Previous uses:**

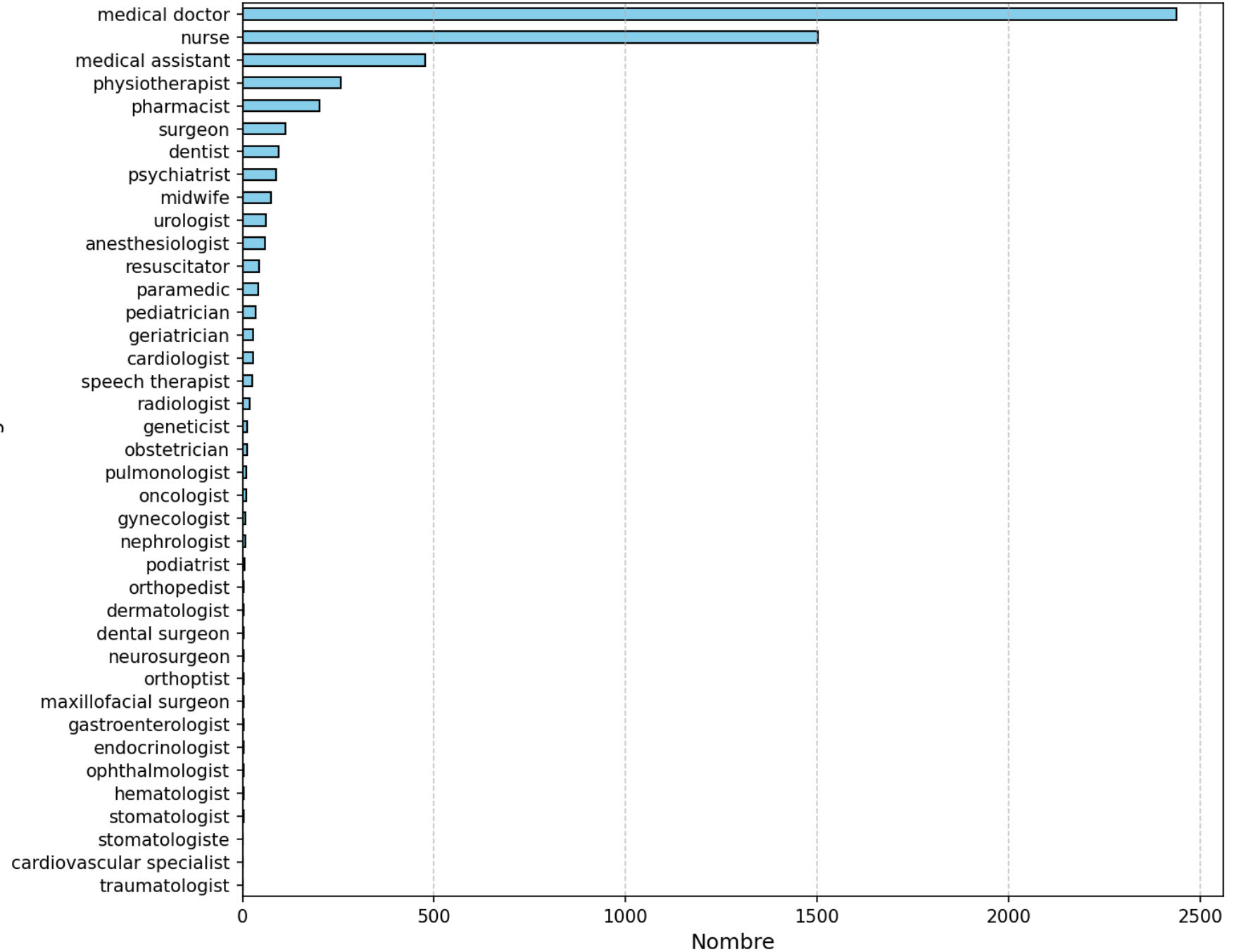
- Hybrid social listening
- Periodic thematic analysis
- Triangulation: interviews, questionnaires (RCCOVID, RECOVER, Transvaxx, VCF Unspoken)

# Distribution of healthcare workers profiles (n= 5676 users)

## Sample discussed today:

- Healthcare profiles: filtered by Twitter bio / medical professions nomenclature
- Period: 2020–2023
- ~ 94,000 tweets

Profil Soignant



# Vaccine narratives taxonomy – ACME

Topics	labels
Safety	safe
	unsure_safe
	not_safe
Usefulness	useful
	unsure_useful
	not_useful
Accessibility	not_accessible
	unsure_accessible
	accessible
	waste
Pathogenic Risk	risky_pathogen
	unsure_risky_pathogen
	not_risky_pathogen
Polarization	criticizing_pro_vax
	criticizing_anti_vax
	unsanitary_other
	de-escalating
Actors' benevolence & competence	trust_actors
	unsure_trust_actors
	not_trust_actors
Probity	vaccine_profitteering
	conspiracy_ideation
Legitimacy	reactance
	coercion_legitimate
information	generic
	candid_question



Hand labeled (qualitatively coded) 600 messages.



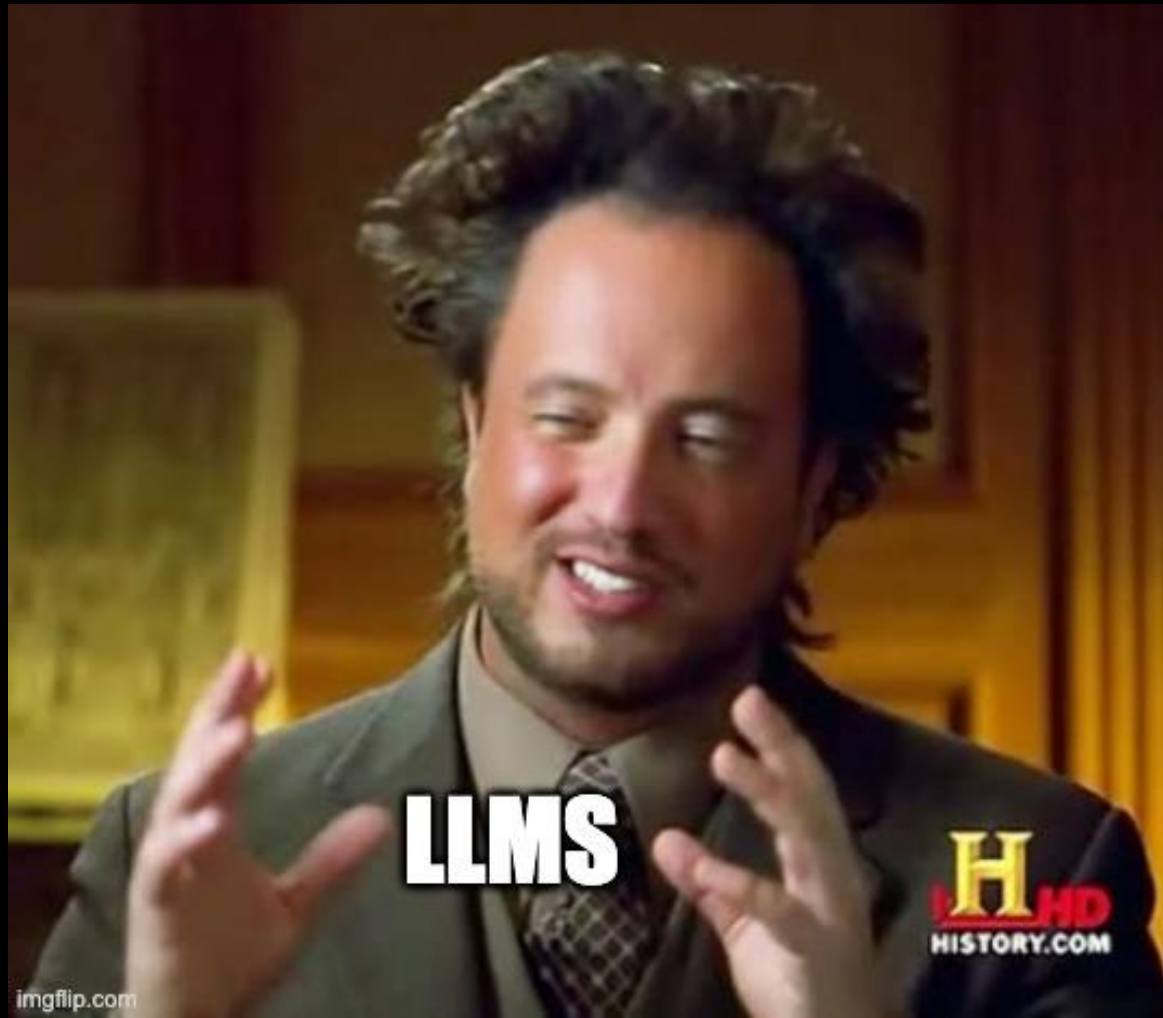
~ 93,400 to go.



How to proceed?

# Evolution of supervised classification options

Era	Techniques	Accuracy	Training/Inference Speed	Volume of Labeled Data Required
<b>Classic Text Mining (2008-2013)</b>	Bag of Words, TF-IDF (-> modelling with RF, XGBOOST using the vectors)	Moderate	Fast	Several thousand labeled texts per label
<b>Word Embeddings (2013-2018)</b>	Word2Vec, GloVe	Improved	Moderate	>500 Hundred to thousands of labeled texts per label
<b>Deep Learning &amp; Transfer Learning (&gt;2018)</b>	LSTM, Transformer models (e.g., BERT, CamemBERT)	Very high	Moderate	>500 Hundred to a few thousand labelled texts
<b>&gt;2020</b>	?	?	?	?



LLMS



imgflip.com

# Evolution of supervised classification options

<b>Era</b>	<b>Techniques</b>	<b>Accuracy</b>	<b>Training/Inference Speed</b>	<b>Volume of Labeled Data Required</b>
<b>Classic Text Mining (2008-2013)</b>	Bag of Words, TF-IDF (-> modelling with RF, XGBOOST using the vectors)	Moderate	Fast	Several thousand labeled texts per class
<b>Word Embeddings (2013-2018)</b>	Word2Vec, GloVe	Improved	Moderate	>500 Hundred to thousands of labeled texts per class
<b>Deep Learning &amp; Transfer Learning (&gt;2018)</b>	LSTM, Transformer models (e.g., BERT)	Very high	Moderate	>500 Hundred to a few thousand labelled texts
<b>Large Language Models (&gt;2020)</b>	Few-shot, zero-shot learning (e.g., GPT-4, Gemini)	Moderate ↔ Very High	Relatively slow	Minimal (dozens to hundreds of texts), or even zero-shot



# Conditions to integrating Large Language Models (LLMs)

## Autonomy

- Augmenting –not replacing– analytical capacity -> deductive coding
- Ability to evaluate different models, prompting strategies
- Keeping data within research group

## Closed models



### Pros

- State of the art, highly efficient models
- Local compute not required

### Cons

- Closed models
- Data & confidentiality ?

## Open source models



[...]

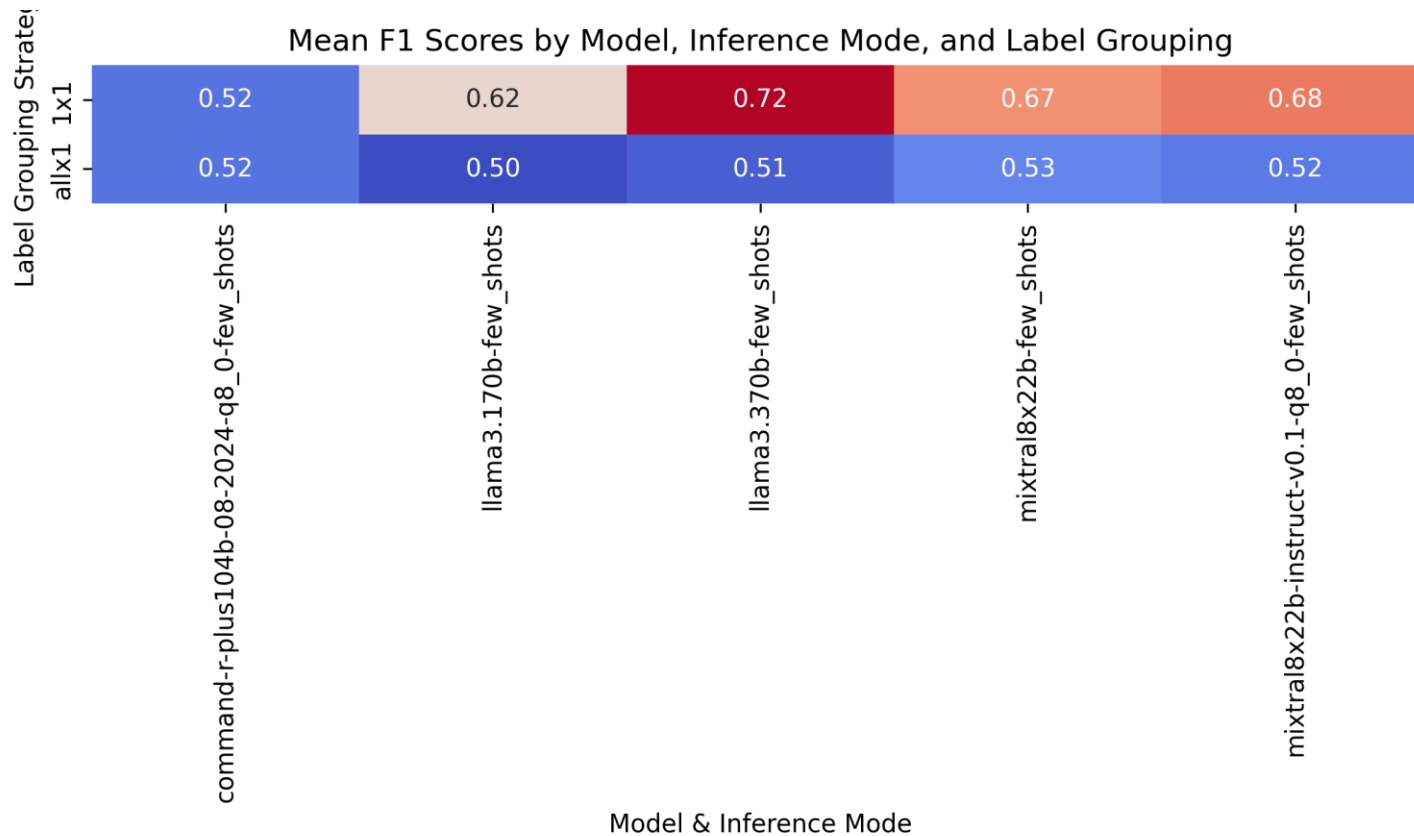
### Pros

- Open source
- Run locally
- Data is not shared

### Cons

- Less accurate
- Requires local compute

# LLMs Evaluation



Mean F1 Score

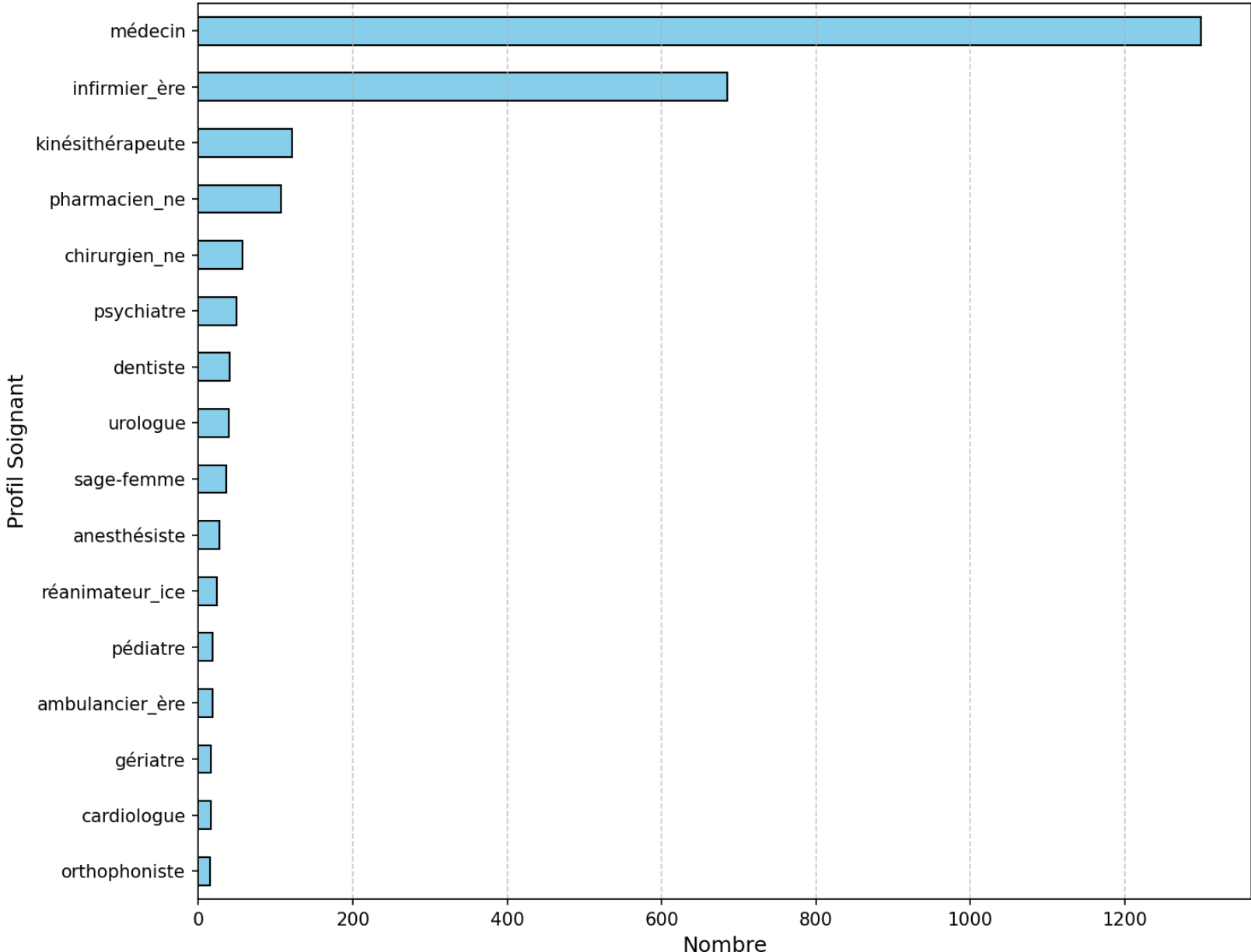
Mixture of models without *uncertain* labels F1 = **0.79**

Topics	labels
Safety	safe
	not_safe
Usefulness	useful
	not_useful
Accessibility	not_accessible
	accessible
	waste
Pathogenic Risk	risky_pathogen
	not_risky_pathogen
Polarization	criticizing_pro_vax
	criticizing_anti_vax
	unsanitary_other
	de-escalating
Actors' benevolence & competence	trust_actors
	not_trust_actors
Probity	vaccine_profitteering
	conspiracy_ideation
Legitimacy	reactance
	coercion_legitimate
information	generic
	candid_question

- Gains through prompt engineering:
  - Contextualizing the task
  - Disambiguating (Collaborative Qual Coding: human / machine) – Ongoing

# Healthcare worker profiles — within the LLM labelled dataset (22,863 tweets)

## Distribution des Profils Soignants (n= 2622 utilisateurs)



**MD over represented:**

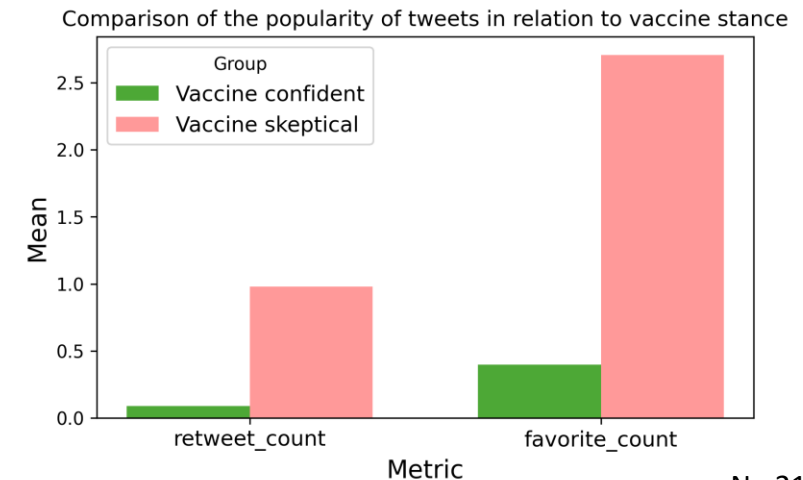
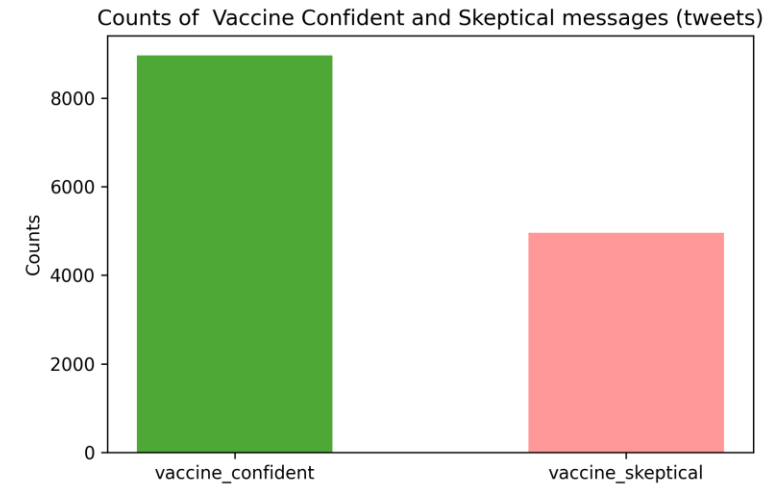
- Status / desirability?
- Specialists not detailing?

# HCW Vaccine skepticism spreads faster

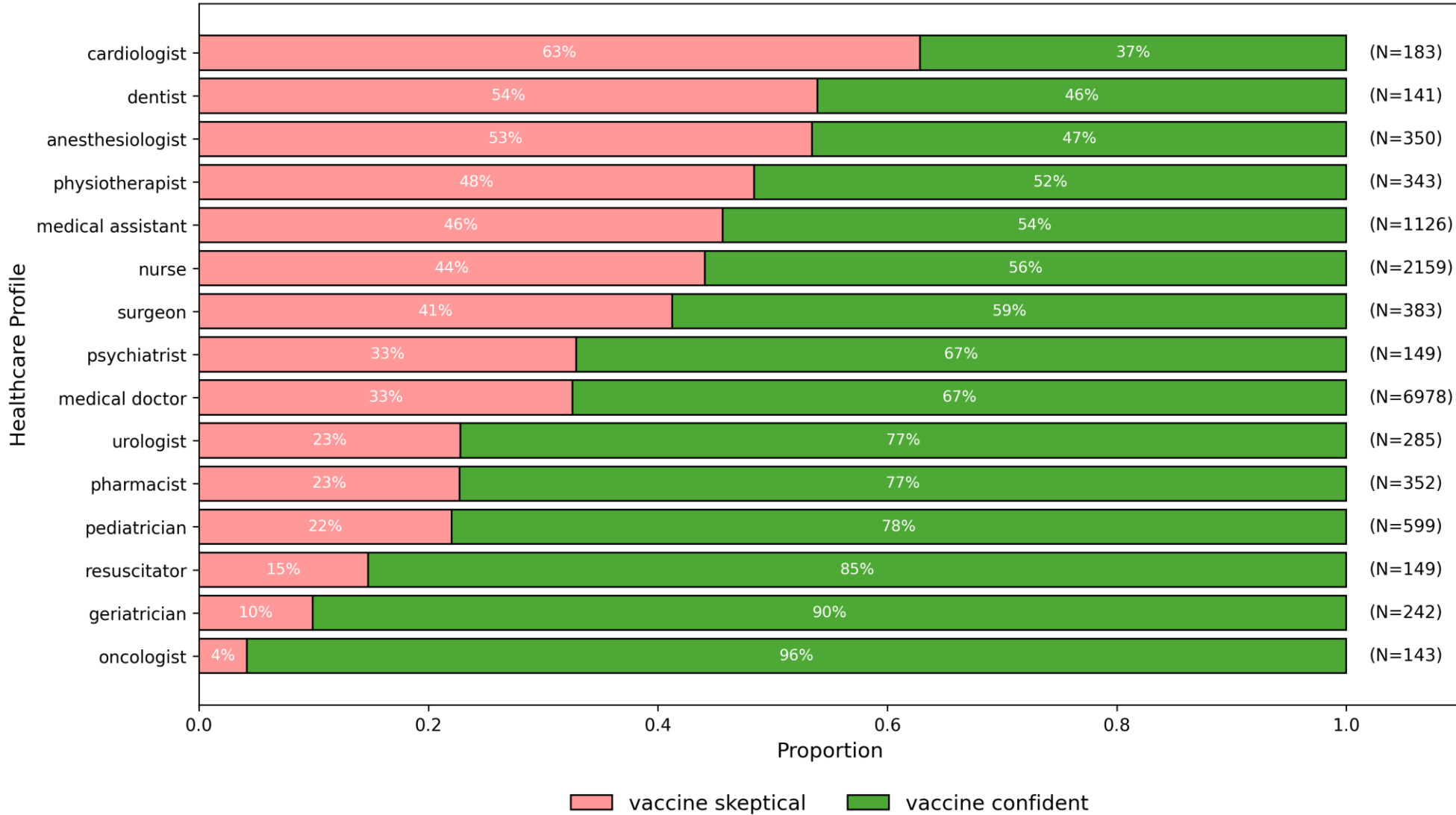
Overarching *vaccine\_skeptical*, *vaccine\_confident* variables created based on key labels (not\_safe, safe..)

Overall, messages are predominantly vaccine confident (nearly x2)

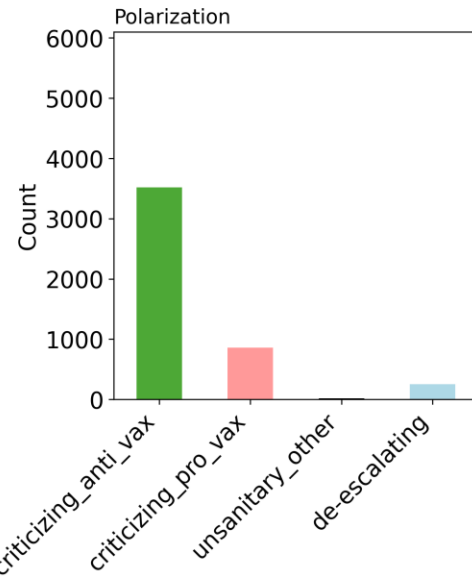
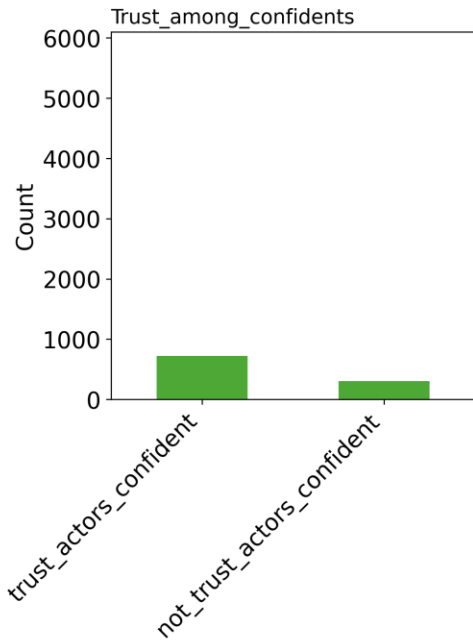
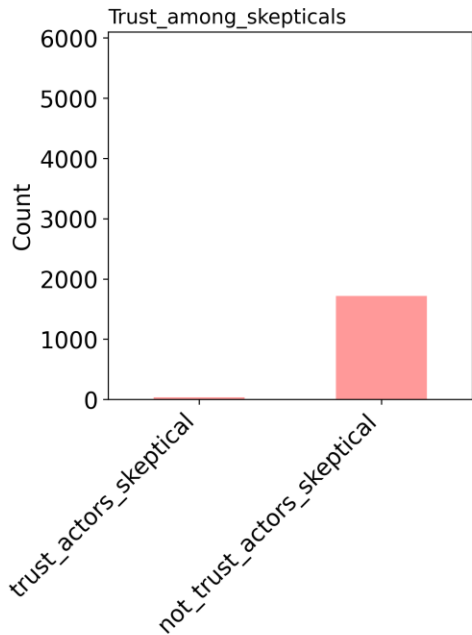
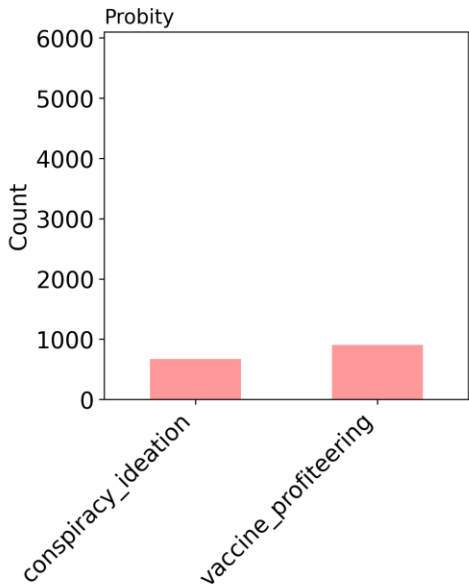
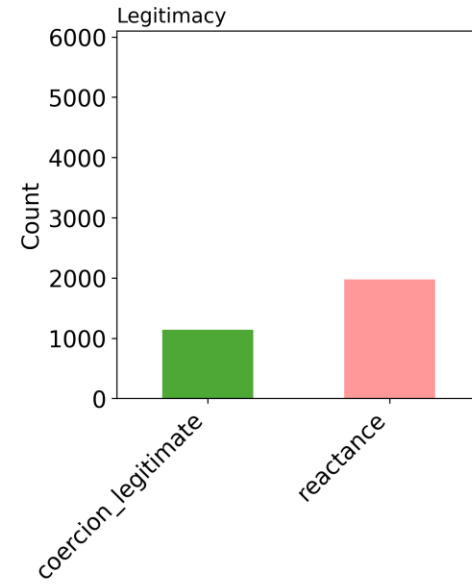
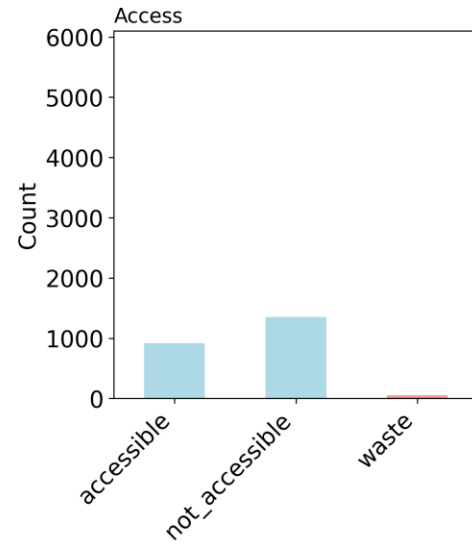
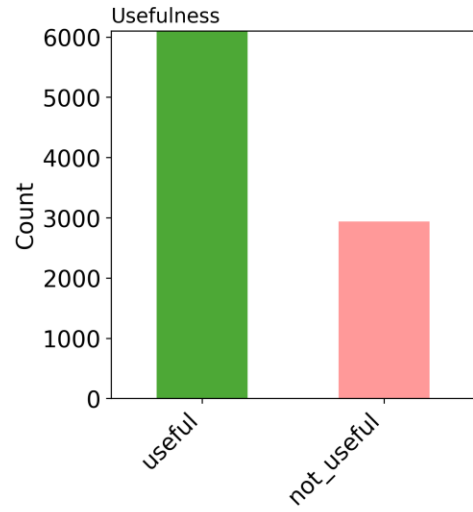
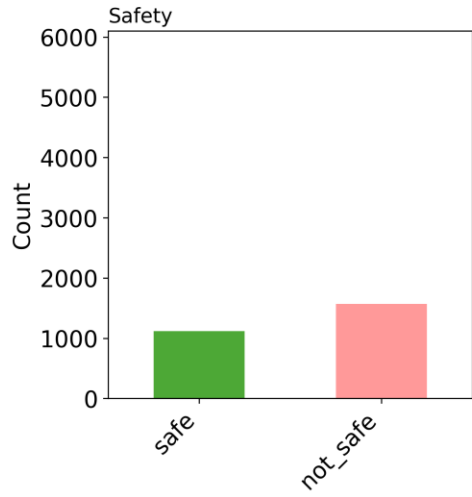
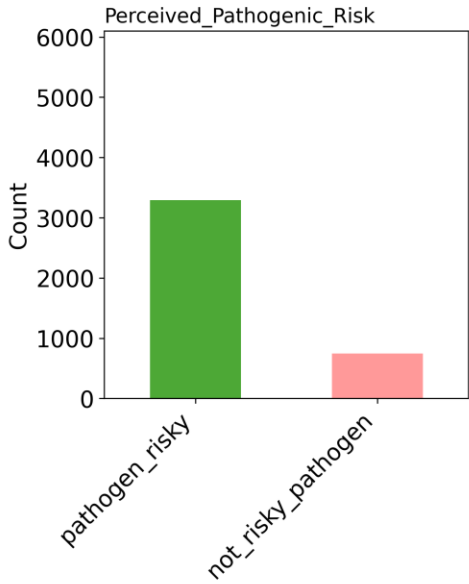
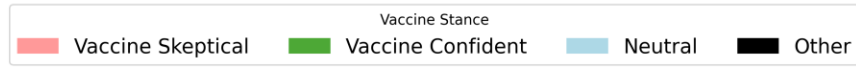
But vaccine skeptical messages are more shared (RT) and receive more favorites.



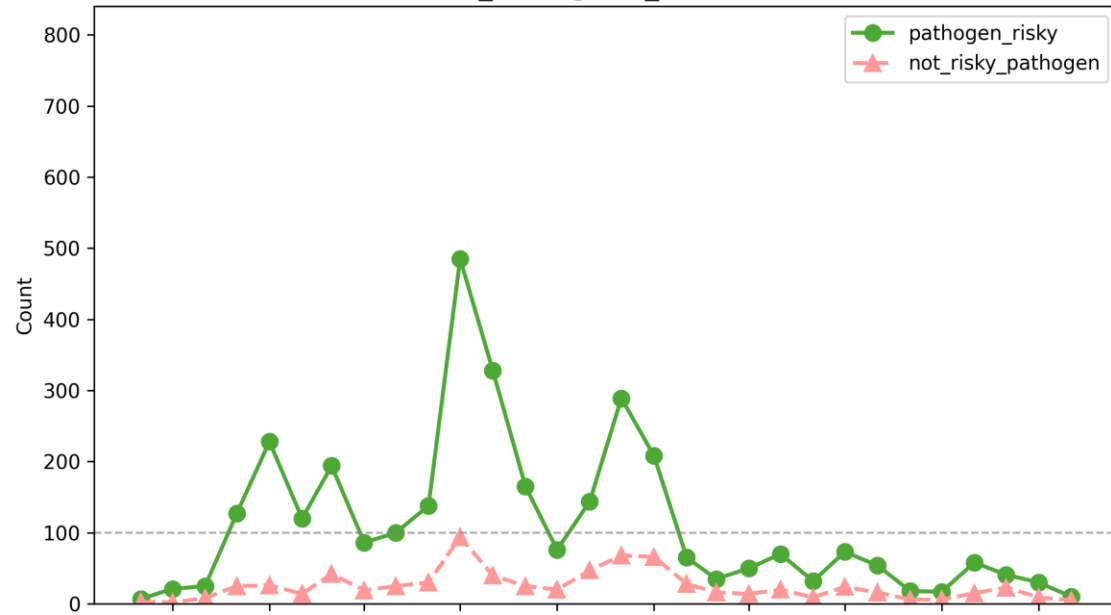
## Proportion of Critics and Supporters by Healthcare Profile (n=22846 messages)



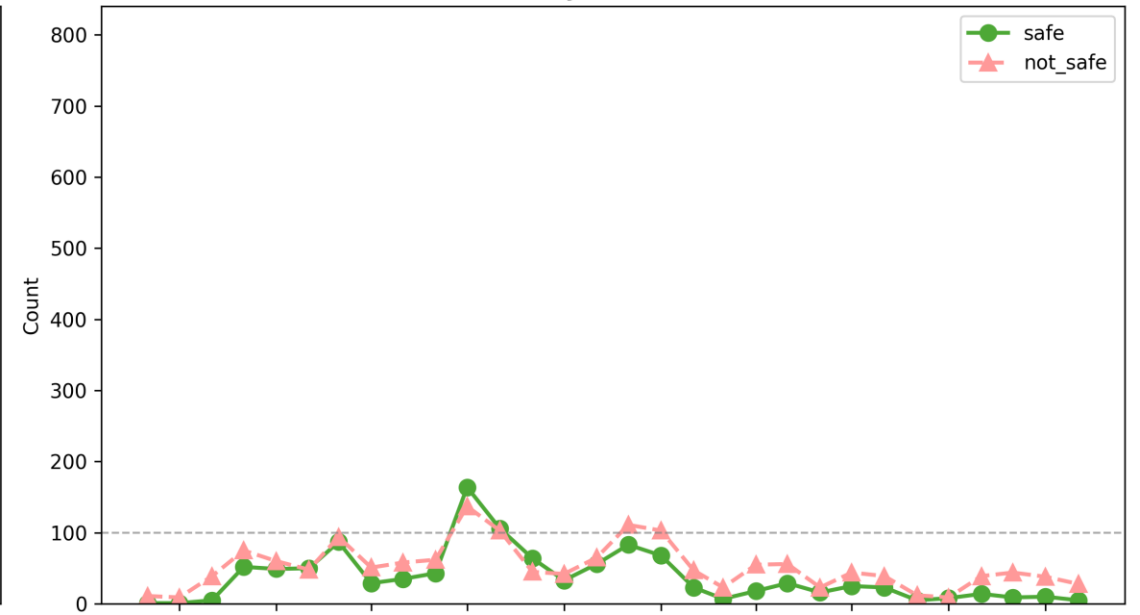
**Label distributions (profile, volume, time)**



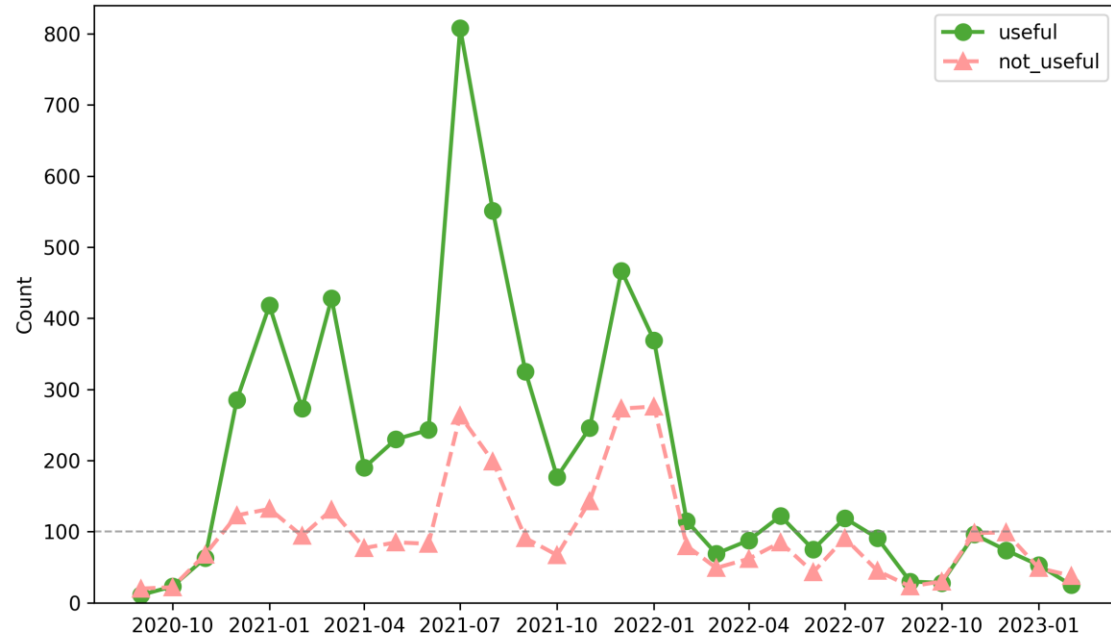
Perceived\_Pathogenic\_Risk Over Time



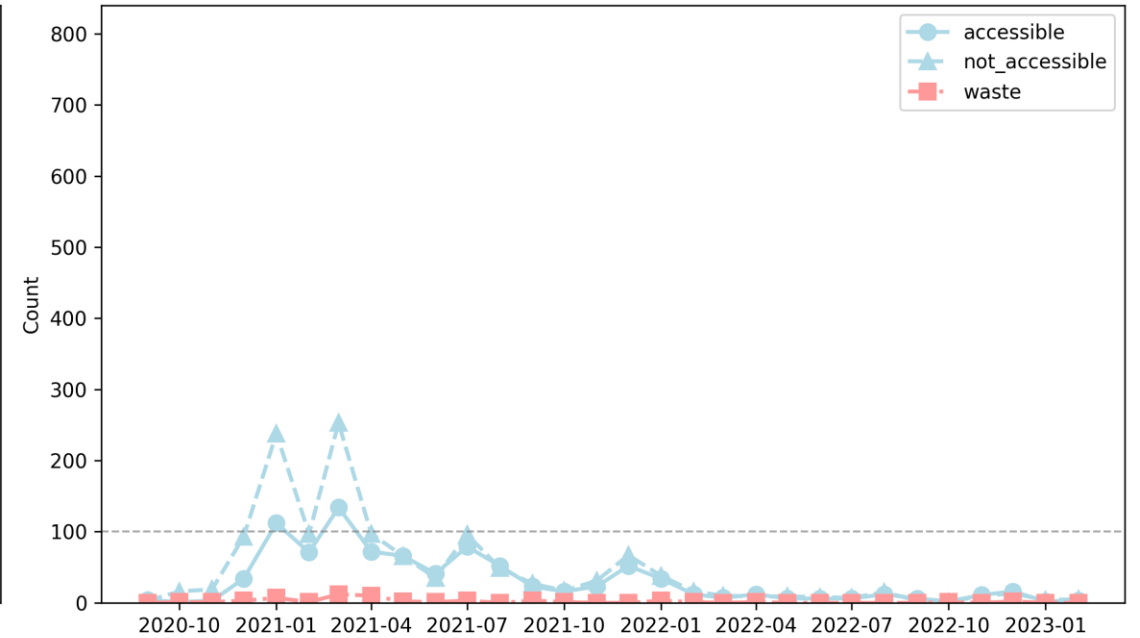
Safety Over Time



Usefulness Over Time



Access Over Time

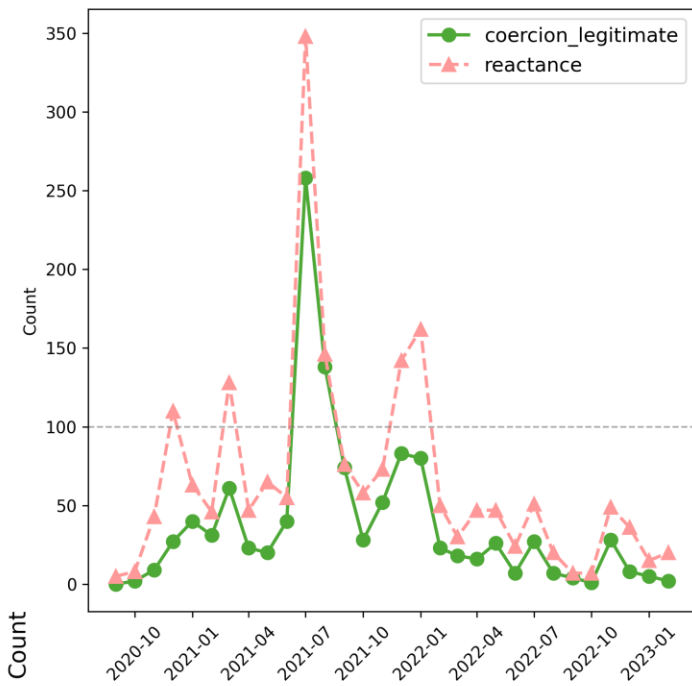


Time (Month)

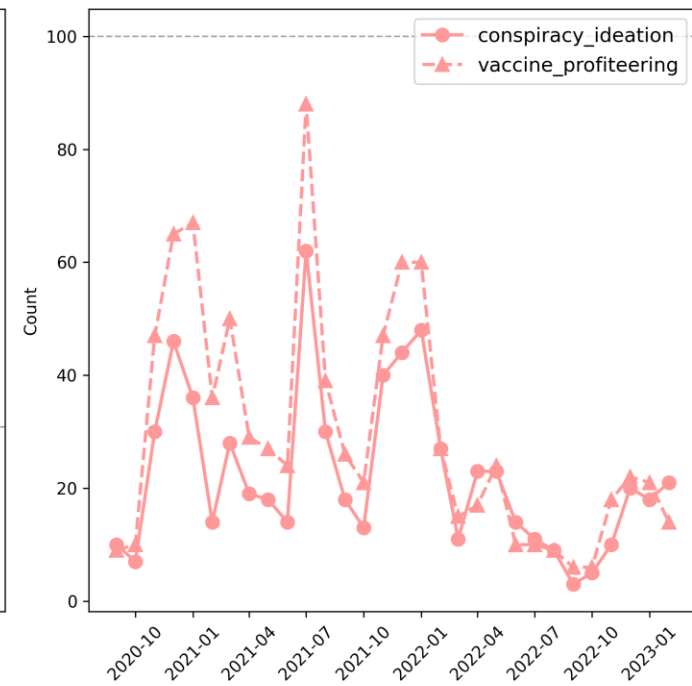
N= 22,846 HCW messages



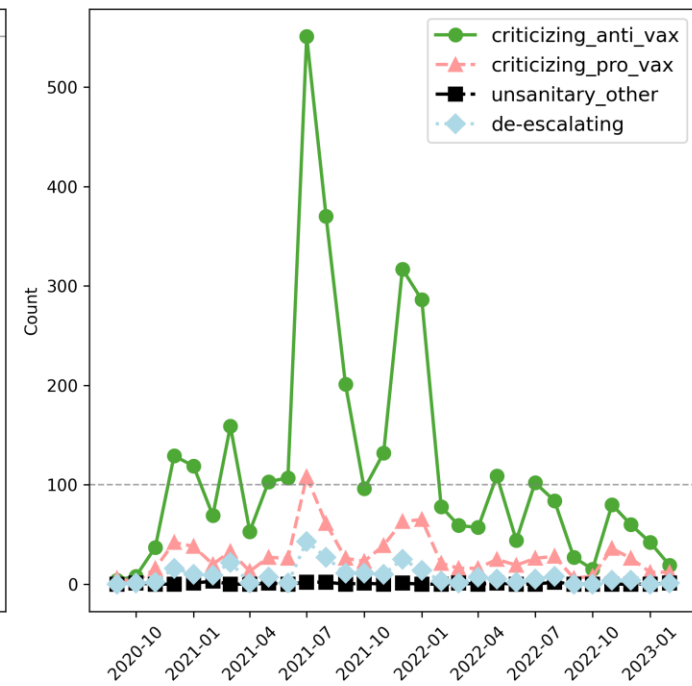
Legitimacy Over Time  
Time (Month)



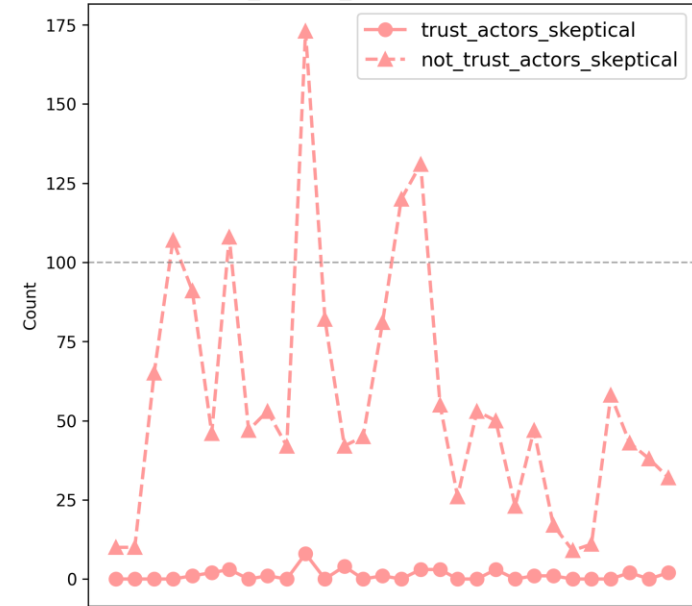
Probity Over Time  
Time (Month)



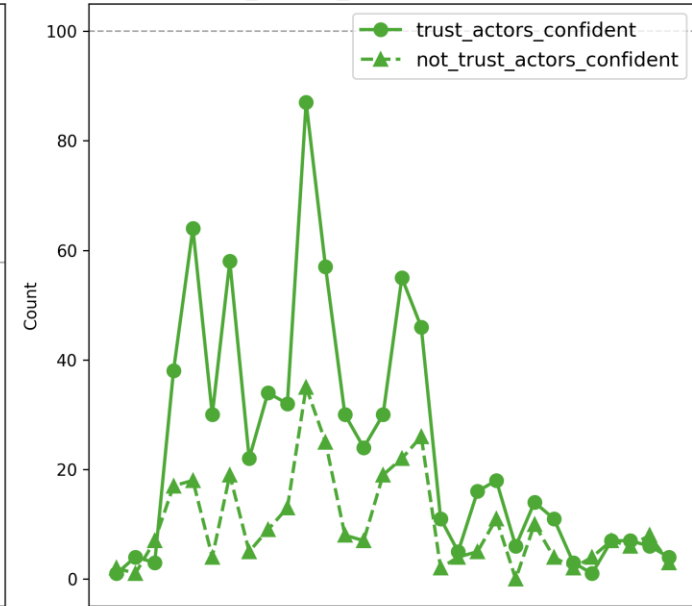
Polarization Over Time  
Time (Month)



Trust\_among\_skepticals Over Time



Trust\_among\_confidants Over Time



# Briefly

- Limitations
  - Preliminary
  - Big data is *thin* data: need to ground / triangulate
- Next steps
  - Enhance human / machine intercoder reliability (disambiguation, prompt engineering, new models).
  - Integration
    - ACME
    - DEV-SHP (ANR, > Feb 25):
      - (Unspoken) vaccine sentiments, climate of vaccine discussion w/ peers & patients.
      - co-construction of vaccine dialogue initiatives.
      - Anthropology, epidemiology, machine learning

# Thank you

- ACME Hybrid Social Listening
  - Gaston Bizel-Bizellot
- ACME WP1 Factors of uptake and adherence, confidence and social equity with regard to epidemic PCMs
  - Judith Mueller
  - Pierre Verger
  - Jocelyn Raude
  - Kathy McColl
- Anthropology and Ecology of Disease Emergence, Institut Pasteur
  - Tamara Giles-Vernick