# DIATOPIT

*A corpus of social media posts*
*for the study of diatopic language variation in Italy*

**Alan Ramponi,**[1] **Camilla Casula**[1,2]

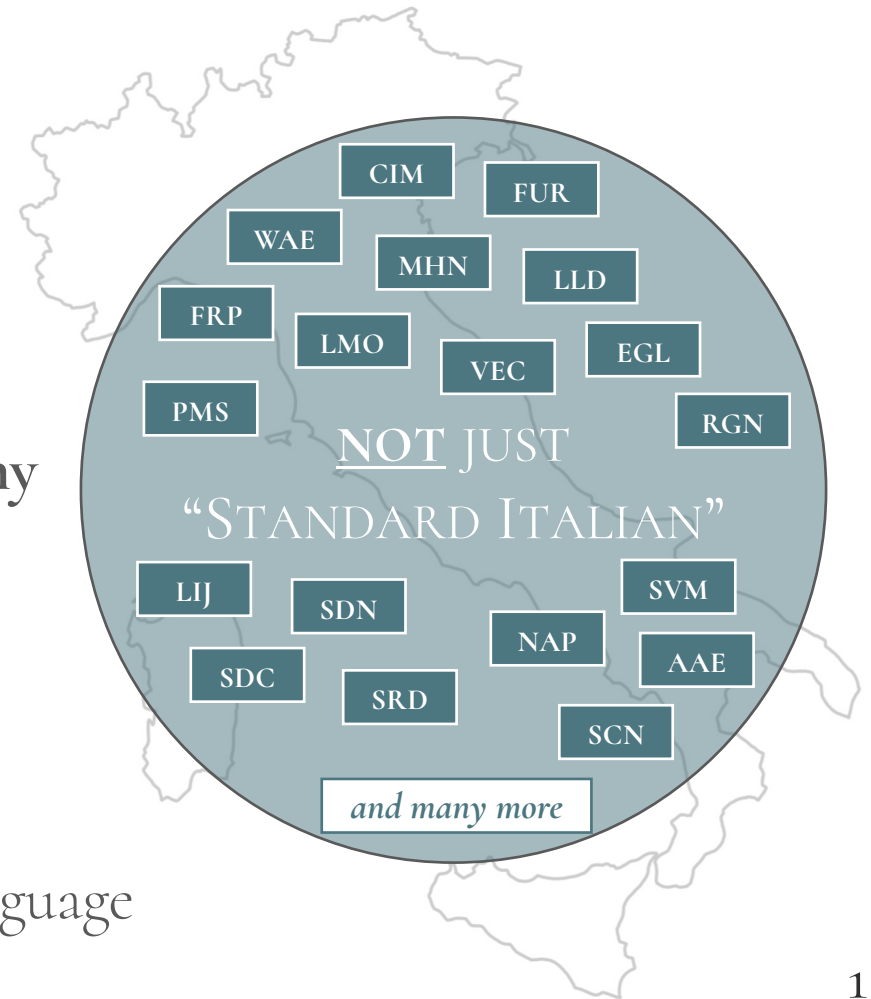[1]*Fondazione Bruno Kessler*   [2]*University of Trento*

# Introduction

**Italy**: linguistically-diverse country

- Many **languages**, **dialects**, and **regional varieties**
- Mostly **oral** and **without established orthography**

**Diatopic language variation** in Italy

- Focal point in **linguistics** (e.g., linguistic atlases)
- **User-generated texts**: informal, spontaneous language

**NOT** JUST "STANDARD ITALIAN"

CIM
FUR
WAE
MHN
LLD
FRP
LMO
EGL
VEC
PMS
RGN
LIJ
SVM
SDN
NAP
AAE
SDC
SRD
SCN
*and many more*

# Contribution

**DIATOPIT**: the first social media corpus focused on **diatopic language variation in Italy** for <u>language varieties other than Standard Italian</u>
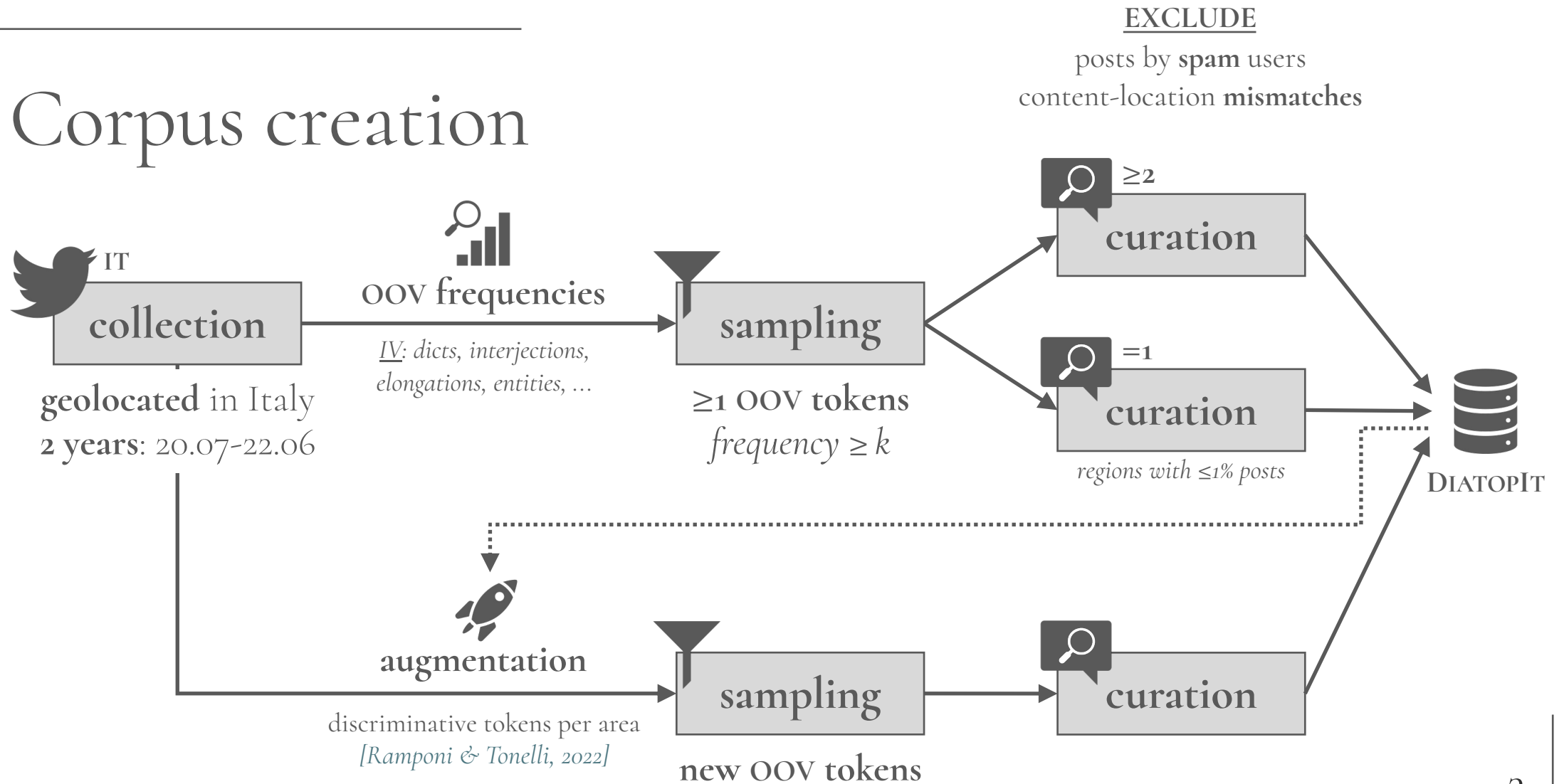
▪ Actual use, orthography choices, code-switching (*language contact* and *vitality*)

**1** **chiov' tutt a jurnat', ce serv' o mbrell'**
 **en.** *it's raining all day, we need an umbrella*

**2** **ho così sonno che me bala l'oeucc**
 **en.** *I'm so sleepy that my eye trembles*

**3** **da caruso anche io ci andavo spesso!**
 **en.** *I used to go there often as a kid too!*

# Corpus creation

**EXCLUDE**
posts by **spam** users
content-location **mismatches**

**IT**

## collection

**geolocated** in Italy
**2 years**: 20.07-22.06

**OOV frequencies**

*IV: dicts, interjections,*
*elongations, entities, …*

## sampling

**≥1 OOV tokens**
*frequency ≥ k*

≥2

## curation

=1

## curation

*regions with ≤1% posts*

**DiatopIt**

**augmentation**

discriminative tokens per area
*[Ramponi & Tonelli, 2022]*

## sampling

**new OOV tokens**

## curation

*Ramponi & Tonelli, 2022. "Features or Spurious Artifacts? Data-centric Baselines for Fair and Robust Hate Speech Detection". NAACL.*
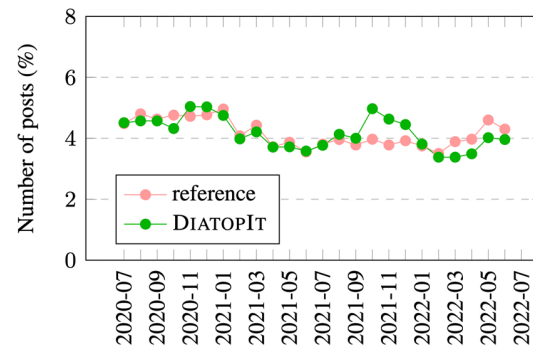
# Corpus analysis: *distribution*

**15K+ posts** by 3,7K authors – *4.1 posts/user*
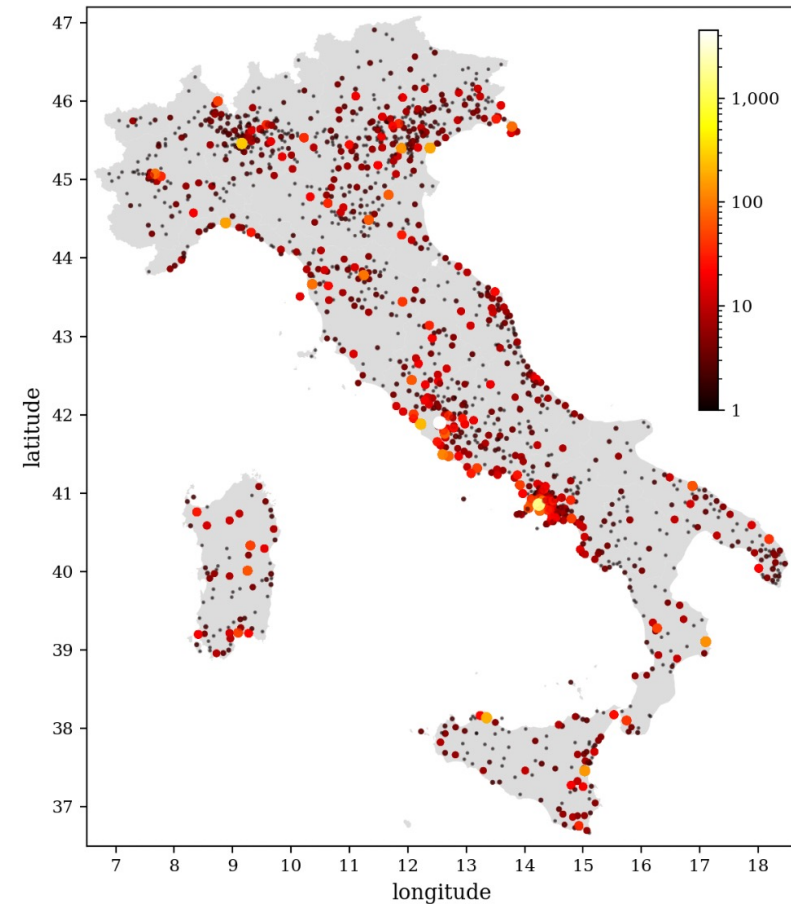
- **55K OOV tokens** – *14.1% avg OOV/post*

- *cities, coastal/lowlands* <u>vs</u> *rural/mountain* areas

*Temporal biases minimized due to the topic-agnostic corpus creation procedure*
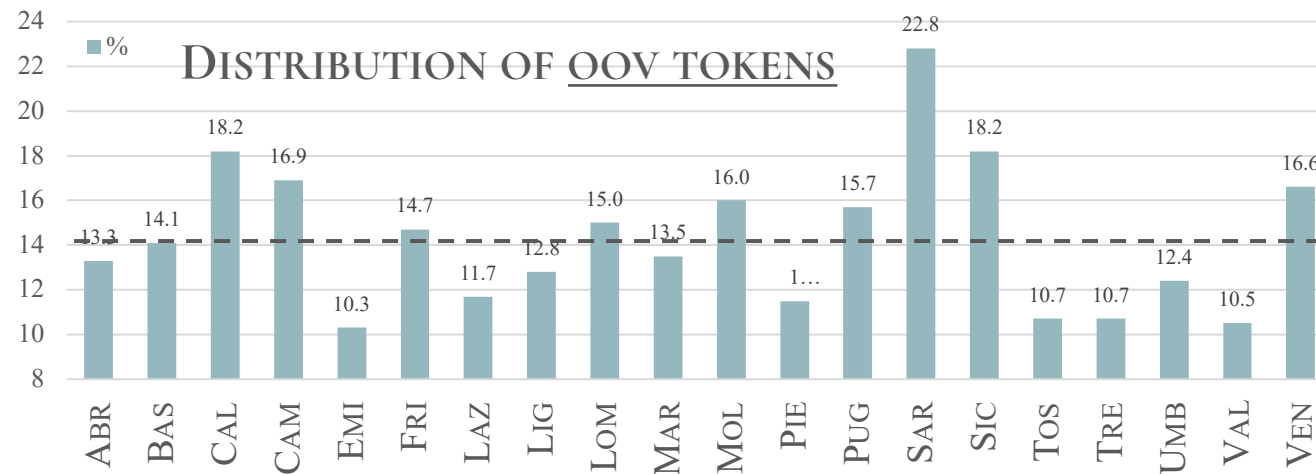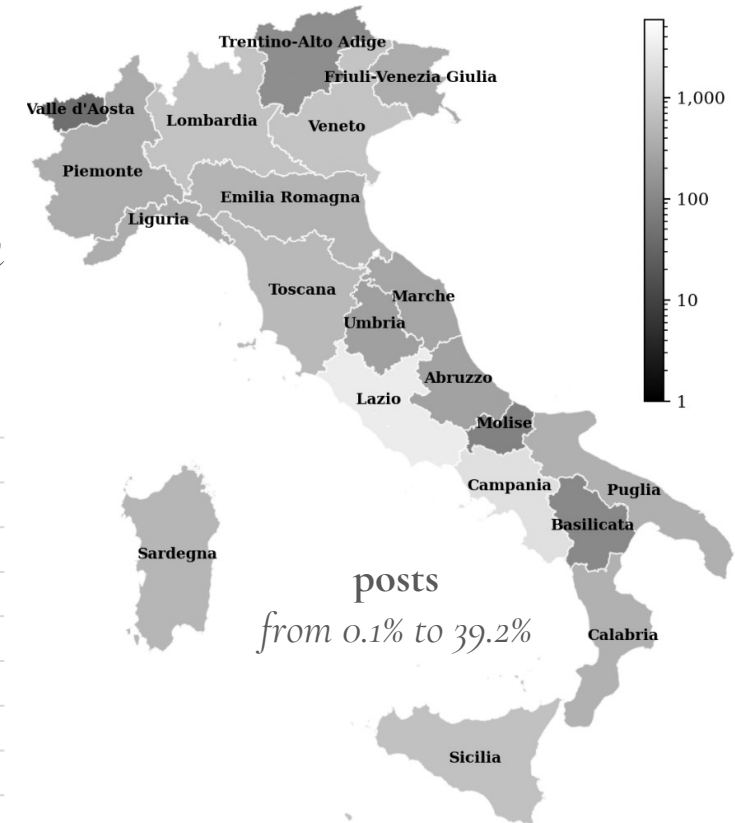
DISTRIBUTION OVER TIME

# Corpus analysis: *distribution*

Regions <u>varies a lot</u> in terms of:

- **posts** – *population & use of varieties other than Italian*

- **OOV tokens** – *non-Standard Italian lexical items*

**posts**
*from 0.1% to 39.2%*
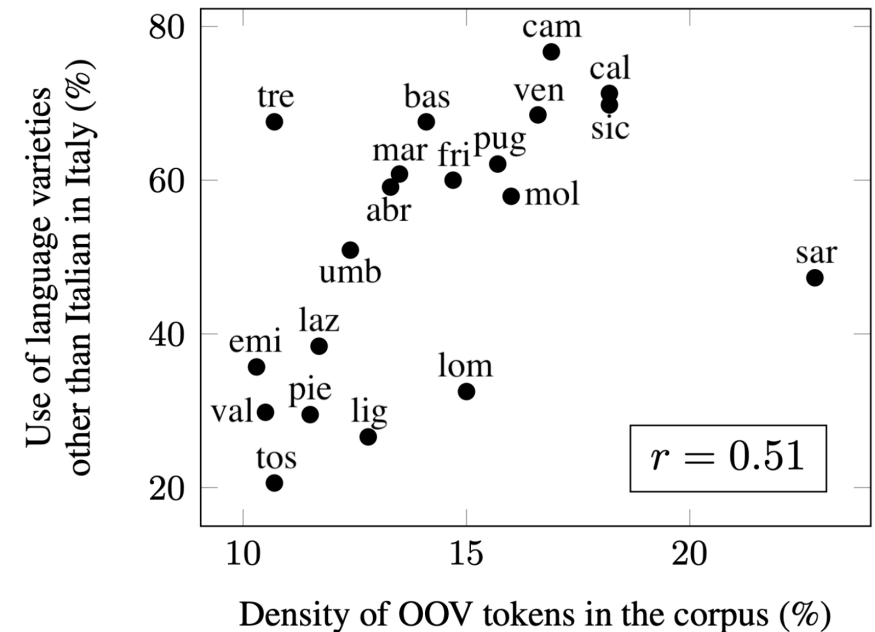


DISTRIBUTION OF OOV TOKENS

5

# Corpus analysis: *OOV density and actual use*

**Degree of mixing** with Standard Italian

- <u>HYP</u>: the *more* a variety is used, the *less* lexical items that belong to Italian are employed

**OOV density** *vs* **use of varieties** *[ISTAT, 2017]*

- Substantial correlation ($r = 0.51$)
- SAR (*speakers' awareness*), TRE (*German varieties*)



*Scatter plot: y-axis "Use of language varieties other than Italian in Italy (%)" ranging 20 to 80; x-axis "Density of OOV tokens in the corpus (%)" ranging 10 to 20. Labeled points: tre, cam, cal, bas, ven, sic, mar, fri, pug, abr, mol, sar, umb, laz, emi, lom, val, pie, lig, tos. Box: $r = 0.51$*

**ISTAT, 2017**. "*L'uso della lingua italiana, dei dialetti e di altre lingue in Italia*". *ISTAT website (accessed 2023-02-01).*

# Corpus analysis: *language use and vitality*

**Most predictive per-region OOV tokens** *[Ramponi & Tonelli, 2022]*

- <u>HYP</u>: the *more* language varieties are spoken in a region, the *higher* the likelihood that non-content OOV tokens (e.g., art, prep, conj) are used

| CAM token | score | SIC token | score | VEN token | score | EMI token | score | TOS token | score |
|-----------|-------|-----------|-------|-----------|-------|-----------|-------|-----------|-------|
| o* | 1.00 | u | 1.00 | ghe | 1.00 | soccia | 1.00 | diaccio | 0.96 |
| e* | 1.00 | bonu | 0.93 | xe | 1.00 | cinno | 0.96 | pigliá | 0.91 |
| tutt | 0.94 | ca | 0.89 | el | 0.96 | maroni | 0.94 | tope | 0.89 |
| nun | 0.90 | cu | 0.88 | no* | 0.83 | cagher | 0.91 | gliè | 0.88 |
| stu | 0.88 | semu | 0.87 | ga | 0.81 | mond | 0.85 | boja | 0.86 |

*confident use of local language varieties*

*restricted function of language varieties*

7

*Ramponi & Tonelli, 2022. "Features or Spurious Artifacts? Data-centric Baselines for Fair and Robust Hate Speech Detection". NAACL.*
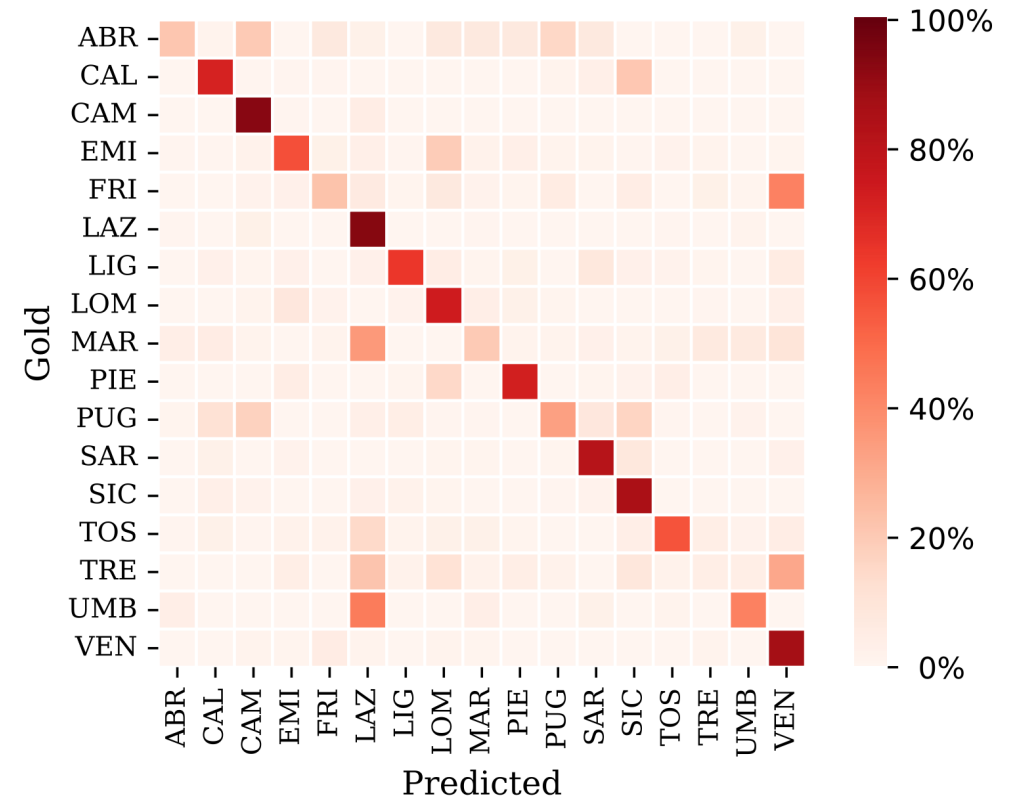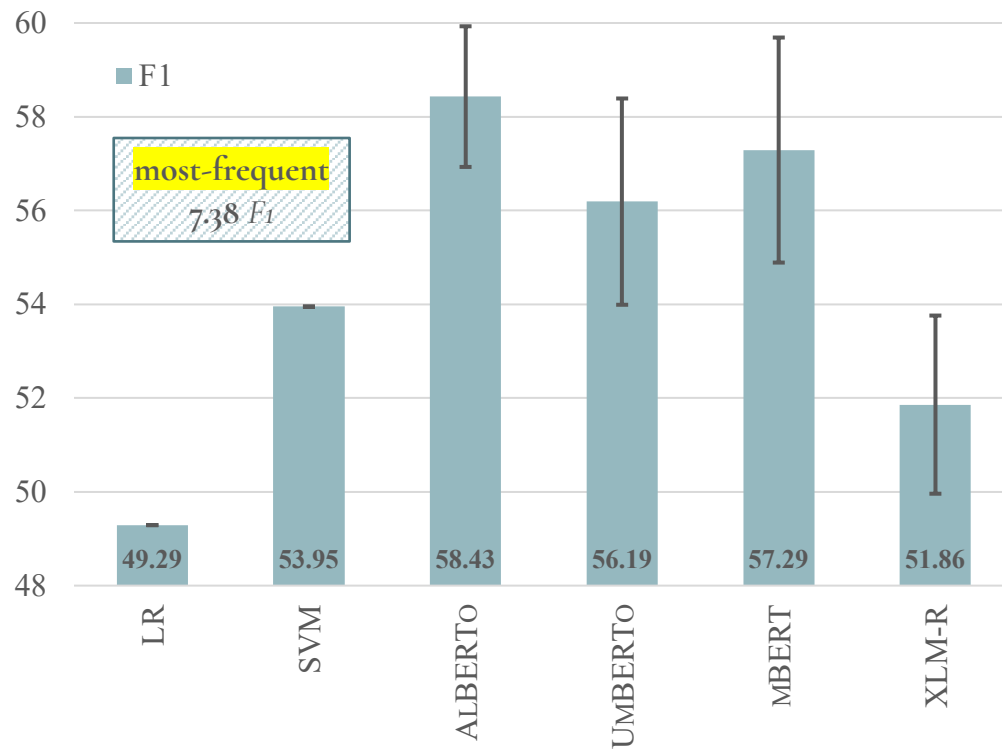
# Experiments

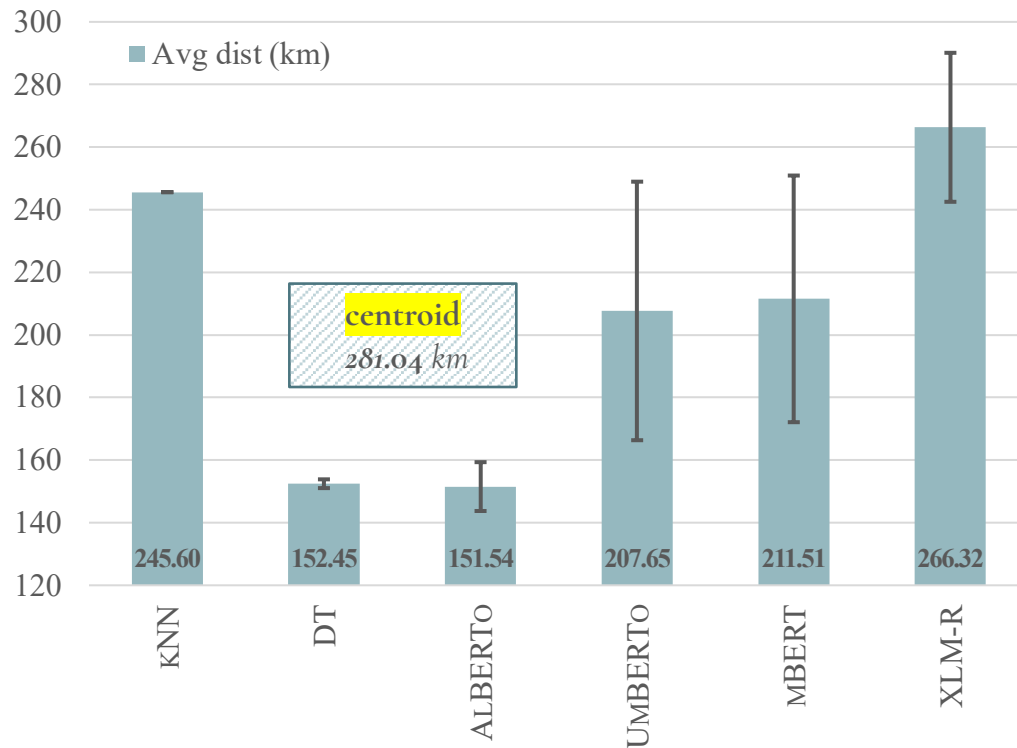*How difficult is it to model diatopic language variation in Italy?*

## Experimental setup

- **Tasks**: *coarse-grained geolocation* [CG] and *fine-grained geolocation* [FG]

- **Evaluation**: *macro F1* [CG] and *avg dist (km)* [FG] *on regions with >50 total posts*

  - *Dev/test* **curation**: wide range of linguistic phenomena and microvariation
  - **Smoothed distribution**: more reliable evaluation for under-represented areas

# Results and discussion: *CG task*

# Results and discussion: *FG task*



Modeling diatopic variation is **hard** for transformer-based models, too

- Limited <u>vocabulary coverage</u>

- <u>Pre-training data</u> impacts stability

# Conclusion

Corpus on *diatopic language variation* in Italy

- (*Partially and even fully*) written in local **languages**, **dialects**, and **regional Italian**

- **Spontaneous**: representative of actual use

- **Varied**: orthography and microvariation

*Useful to study diatopic variation, code-switching and divergences in orthography, in order to assess vitality across local language varieties*

NOT JUST "STANDARD ITALIAN"