# GEOLINGIT AT EVALITA 2023
## *Overview of the Geolocation of Linguistic Variation in Italy Task*

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

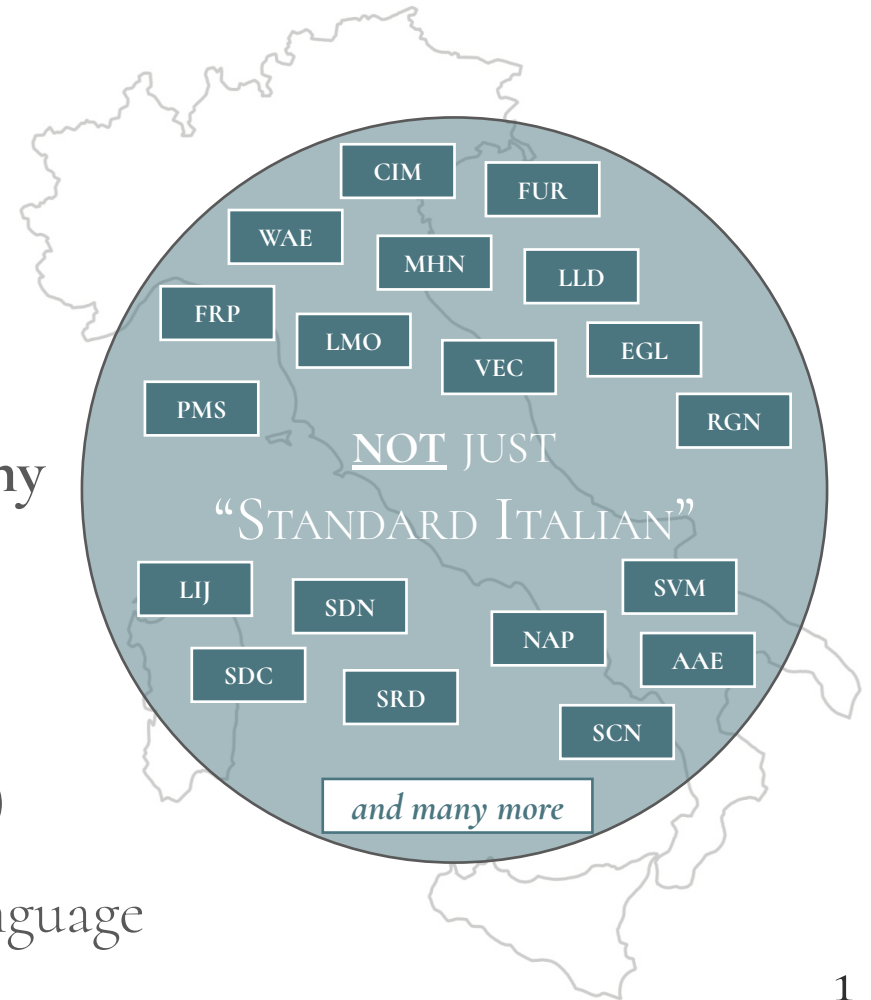[1]*Fondazione Bruno Kessler*   [2]*University of Trento*

# 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



CIM
FUR
WAE
MHN
LLD
FRP
LMO
EGL
VEC
PMS
RGN

**NOT** JUST
"STANDARD ITALIAN"

LIJ
SVM
SDN
NAP
AAE
SDC
SRD
SCN

*and many more*

# Data

**DIATOPIT**: the first social media corpus focused on **diatopic language variation in Italy** for *language varieties other than Standard Italian*

- 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!*

2

*Ramponi & Casula, 2023. "DiatopIt: A Corpus of Social Media Posts for the Study of Diatopic Language Variation in Italy". VarDial@EACL.*

# Data

SOURCE: Twitter, geolocated in Italy
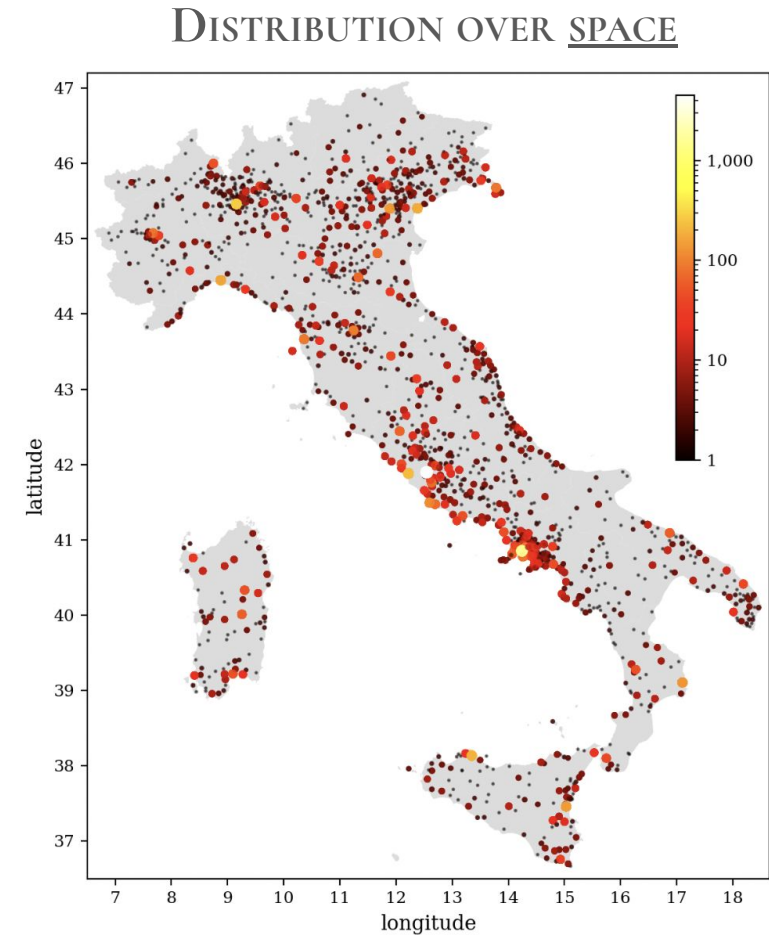TIMEFRAME: 2 years [2020-07 – 2022-06]
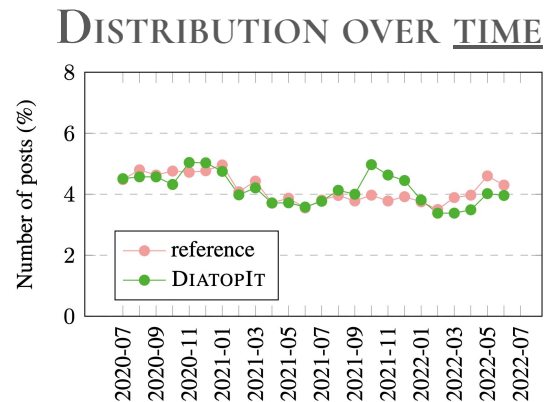SAMPLING: based on (*curated*) OOV tokens
CURATION: manual exclusion of spam/mismatches
AUGMENTATION: for under-represented areas

**15K+ posts** by 3,7K users

- *Thorough corpus analysis in Ramponi & Casula (2023)*

DISTRIBUTION OVER TIME



DISTRIBUTION OVER SPACE



3

*Ramponi & Casula, 2023. "DiatopIt: A Corpus of Social Media Posts for the Study of Diatopic Language Variation in Italy". VarDial@EACL.*

# Task description and evaluation

*Given the text of a post exhibiting regional Italian features or (partially or fully) written in local languages and dialects of Italy,*
**predict the <u>location</u> in which the variety expressed in the post is spoken**

<u>**Tracks**</u>: STANDARD TRACK (*country-level*) or SPECIAL TRACK (*linguistic area of choice*)

▪ **A) Coarse-grained geolocation**. Predict the region; <u>macro F1 score</u>

▪ **B) Fine-grained geolocation**. Predict the lat/lon coordinates; <u>avg dist (km)</u>

*<u>*Dev/test sets</u>: smoothed distribution, further curation to include a wide range of linguistic phenomena / microvariation*

# Participation

**37 registrations** and <mark>35</mark> **submitted runs**

*Heterogeneously composed teams with up to 7 individuals, from master students to senior academic researchers*

| | Subtask A (coarse-grained) | Subtask B (fine-grained) | Total |
|---|---|---|---|
| **Standard track** | **14 (5)** | **12 (5)** | *26 (6)* |
| **Special track** | **6 (2)** | **3 (1)** | *9 (2)* |
| *Total* | *20 (5)* | *15 (5)* | *35 (6)* |

LMU (3)

UniBz (1)

FBK (1)

Maize (1)

PoliTo (10)

UniTo (1)

UniPi (2)

CNR-Isti (1)

UniRoma2 (1)

UniKore (1)

# Subtask A: *Methods and results*

| | Team | Run | P | R | F$_1$ |
|---|---|---|---|---|---|
| 1 | DANTE | (3) | 79.46 | 63.75 | 66.30 |
| 2 | DANTE | (2) | 66.98 | 62.65 | 63.93 |
| 3 | DANTE | (1) | 65.18 | 60.09 | 61.72 |
| 4 | galliz | (1) | 82.94 | 52.25 | 56.20 |
| 5 | ba$\rho$tti | (2) | 67.97 | 51.62 | 53.18 |
| 6 | galliz | (3) | 74.58 | 49.49 | 52.08 |
| 7 | ba$\rho$tti | (3) | 52.93 | 51.75 | 51.74 |
| 8 | ba$\rho$tti | (1) | 56.05 | 51.68 | 51.72 |
| 9 | galliz | (2) | 68.98 | 45.36 | 47.74 |
| | *Log. reg.* | | *62.19* | *42.43* | *46.11* |
| 10 | extremITA | (1) | 72.14 | 38.84 | 39.99 |
| 11 | extremITA | (2) | 65.03 | 37.62 | 38.18 |
| 12 | SCG | (2) | 12.92 | 9.82 | 9.28 |
| 13 | SCG | (1) | 10.15 | 9.97 | 9.04 |
| 14 | SCG | (3) | 10.42 | 6.60 | 7.85 |
| | *Most freq.* | | *1.24* | *5.88* | *2.05* |

| Model & pretraining data | | Extra pretraining methods & data | Fine-tuning methods and data |
|---|---|---|---|
| Ens(15× 👨) | — | — | — |
| 👨 | it | MTL, sources: ≫ W | STL, DiatopIt |
| 👨 | it | MTL, sources: ≫ W | STL, DiatopIt |
| Ens(👨 + 📖) | en | — | STL, *augmented* DiatopIt |
| 👨 | it | CPT, *augmented* DiatopIt w/ ≫ W | STL, DiatopIt |
| Ens(👨 + 📖) | en | — | STL, *augmented* DiatopIt |
| Ens(👨 + lr) | it | — | STL, DiatopIt |
| 👨 | it | — | MTL, DiatopIt |
| Ens(👨 + 📖) | en | — | STL, *augmented* DiatopIt |
| T5 | it | — | MTL, EVALITA 2023 |
| 🦙 | * | LoRA, Alpaca instructions in "**it**" | MTL, EVALITA 2023 |
| svm | — | — | STL, DiatopIt |
| svm | — | — | STL, DiatopIt |
| lr | — | — | STL, DiatopIt |

6

**Models**: **Ens(*)** Ensemble; 👨 BERT-based; 🦙 LLaMa-based; 📖 dictionary-based ( ≫ + DiatopIt); T5 T5-based; svm Support vector machines; lr Logistic regression
**Methods**: **STL**: Single-task learning; **MTL**: Multi-task learning; **CPT**: Contrastive pretraining    **Sources**: ≫ Dialettando; W Wikipedia language editions

# Subtask A: *Analysis*



Results differ <u>a lot</u> between regions

*e.g.*, **Abr, Mar, Tre** and **Umb**
[scarce in `train`, absent in `dev`]
*Very hard to model w/ traditional learning and tuning methods*
– Beyond traditional learning/tuning

*e.g.*, **Cal, Emi, Fri** and **Pug**
[represented in `train`, present in `dev`]
*Easily misclassified w/ regions in which similar varieties are predominantly used*
– Beyond "raw modeling": linguistics!

7

# Subtask B: *Methods and results*

| | Team | Run | Avg dist (km) |
|---|---|---|---|
| 1 | baρtti | (3) | 97.74 |
| 2 | baρtti | (1) | 98.79 |
| 3 | DANTE | (3) | 110.35 |
| 4 | DANTE | (2) | 112.58 |
| 5 | DANTE | (1) | 114.00 |
| 6 | baρtti | (2) | 120.02 |
| 7 | extremITA | (1) | 126.10 |
| 8 | Salogni | (1) | 128.19 |
| 9 | extremITA | (2) | 145.15 |
| | *kNN* | | *263.35* |
| 10 | SCG | (1) | 280.99 |
| | *Centroid* | | *281.04* |
| 11 | SCG | (2) | 281.20 |
| 12 | SCG | (3) | 289.91 |

| Model & pretraining data | | Extra pretraining methods & data | Fine-tuning methods and data |
|---|---|---|---|
| 👨 + **Postproc** | it | – | **MTL**, DiatopIt |
| 👨 | it | **CPT**, *augmented* DiatopIt w/ ≫ W | **MTL**, DiatopIt |
| **Ens(2 × 👨)** | it | – | – |
| 👨 | it | **MTL**, sources: ≫ W | **STL**, DiatopIt |
| 👨 | it | **MTL**, sources: ≫ W | **STL**, DiatopIt |
| 👨 | it | – | **MTL**, DiatopIt |
| T5 | it | – | **MTL**, EVALITA 2023 |
| **Ro🐑a** | – | – | **STL**, DiatopIt |
| 🦙 | * | LoRA, Alpaca instructions in "**it**" | **MTL**, EVALITA 2023 |
| lr | – | – | **STL**, DiatopIt |
| lr | – | – | **STL**, DiatopIt |
| knn | – | – | **STL**, DiatopIt |

8

**Models**: **Ens(*)** Ensemble; 👨 BERT-based; 🦙 LLaMa-based; **Postproc** Geographical postprocessing; T5 T5-based; knn *k*-nearest neighbors; lr Logistic regression
**Methods**: **STL**: Single-task learning; **MTL**: Multi-task learning; **CPT**: Contrastive pretraining    **Sources**: ≫ Dialettando; W Wikipedia language editions
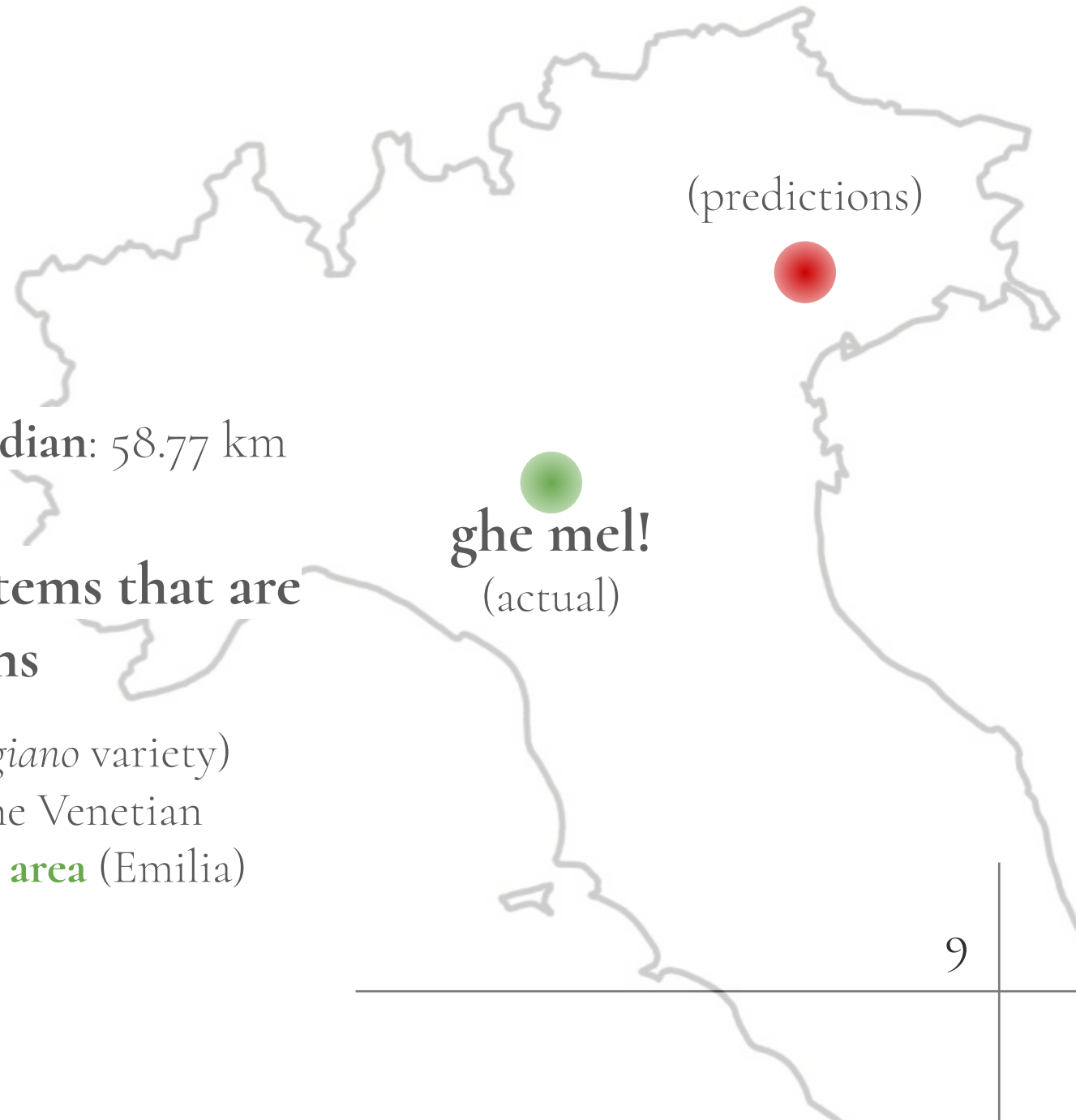
# **Subtask B**: *Analysis*

Test set prediction (avg error)

- **Min**: 0.89 km, **max**: 668.11 km, **median**: 58.77 km

Misclassification due to **lexical items that are highly frequent in other locations**

- e.g., "*ghe* mel!" (en: "of course", *Parmigiano* variety) in the **Treviso area** (Veneto, due to the Venetian ADV/PROP "ghe") instead of the **Parma area** (Emilia)

(predictions)

**ghe mel!**
(actual)

9

# Subtask A and B: *Methods and results*

| Team | Run | P | R | $F_1$ |
|------|-----|-----|-----|-----|
| **TUSCANY-LAZIO AREA** | | | | |
| 1 galliz | (3) | 81.25 | 83.32 | 82.20 |
| 2 galliz | (1) | 72.43 | 80.42 | 73.40 |
| 3 galliz | (2) | 72.43 | 80.42 | 73.40 |
| *Log. reg.* | | *91.79* | *66.67* | *70.53* |
| *Most freq.* | | *38.62* | *50.00* | *43.58* |

*Model & pretraining data*

*Extra pretraining methods & data*

**Ens(👨+📖)** en    STL, *augmented* DiatopIt
**Ens(👨+📖)** en    STL, *augmented* DiatopIt
**Ens(👨+📖)** en    STL, *augmented* DiatopIt

Dialettando

Wikipedia

Vocabolario del Fiorentino Contemporaneo

The Roman Post

The **Gallo-Italic area** (PIE, LOM, LIG and EMI) has been explored, too

- **Subtask A**: SCG team, 3 runs based on Logistic Regression and Support Vector Machines

- **Subtask B**: SCG team, 3 runs based on Logistic Regression and k-Nearest Neighbors

# Discussion and conclusion

- GeoLingIt has attracted **wide interest** from the community

- Modeling diatopic variation in Italy is a **difficult but exciting task**

- Great opportunities for more **linguistically-grounded NLP**

# Poster booster session (~3 min each)

**Baptti team**
*A. Koudounas, F. Giobergia, I. Benedetto, S. Monaco, L. Cagliero, D. Apiletti, and E. Baralis*

**DANTE team**
*G. Gallipoli, M. La Quatra, D. Rege Cambrin, S. Greco, and L. Cagliero*

**Galliz team**
*T. Labruna and S. Gallo*

**Salogni team**
I. Salogni