

# Dynamic Aggregation and Auto-Discovery of Linguistic Features

Yanchen Liu



# DADA: Dialect Adaptation via Dynamic Aggregation of Linguistic Rules

Yanchen Liu, William Held, Diyi Yang

will appear on EMNLP 2023







## Non-Standard Linguistic Features



- □ language usage that deviates from the conventions
- often associated with specific social or cultural groups

#### **Dialect:** A Group of Non-Standard Linguistic Variations in a Language

# Fails on Dialects



- Existing models mainly focus on Standard American English (SAE)
- □ Significant **performance drop**↓when applied to English Dialects



# Previous Work on Dialect Adaptation



#### Mainly focused on targeted adaptation to a specific dialect

- Human Annotation (Blevins et al., 2016, Blodgett et al., 2018)
- Weak Supervision (Jorgensen et al. 2016, Jurgens et al. 2017)
- Alignment (TADA 2023)

### Highly accurate dialect identification systems are required!

#### Incorporating Dialectal Variability Pre-Trained SAE Model for Socially Equitable Language Identification Automatically Processing Tweets from Gang-Involved Youth: Towards **Detecting Loss and Aggression** David Jurgens Yulia Tsyetkov Dan Jurafsky Stanford University Stanford University Stanford University **Terra Blevins** Robert Kwiatkowski Jamie Macbeth {jurgens,tsvetkov,jurafsky}@stanford.edu Department of Department of Department of Electrical and SAE Computer Systems Engineeri Computer Science Computer Science Columbia University Columbia University Fairfield University SAE |CLS<sub>SAF</sub>-CLS<sub>Dial</sub>|<sub>2</sub> - Adv(Dial) New York, NY, USA New York, NY, USA Fairfield, CT, USA 1. @username R u a wizard or wat gan sef: in d mornin Abstract Non-SAE u tweet, afternoon - u tweet, nyt gan u dey tweet. beta t1b2145@columbia\_edu rik21470columbia.edu imacbeth@fairfield.edu get ur IT placement wiv twitter VALUE Language identification (LID) is a criti- Be the lord lantern jaysus me heart after that match!!! Aku hanya mengagumimu dari jauh sekarang . RDK Kathleen McKeown Desmond Patton **Owen Rambow** cal first step for processing multilingual ({}) \* last tweet about you -.- , maybe text. Yet most LID systems are not de-Department of School of Center for Computational Figure 1: Challenges for socially-equitable LID in Twitter signed to handle the linguistic diversity of Computer Science Social Work Learning Systems include dialectal text, shown from Nigeria (#1) and Ireland global platforms like Twitter, where lo-Columbia University Columbia University Columbia University (#2), and multilingual text (Indonesian and English) in #3. cal dialects and rampant code-switching Task-Agnostic New York, NY, USA New York, NY, USA New York, NY, USA lead language classifiers to systematically Dialect Adapters graphic and dialectal variation. As a result, these kathy@cs.columbia.edu dp2787@columbia.edu rambow@ccls.columbia.edu miss minority dialect speakers and mul-Adv(Dial) - Adv(SAE) tilingual speakers. We propose a new systems systematically misclassify texts from populations with millions of sneakers whose local dataset and a character-based sequence-to-Abstract speech differs from the majority dialects (Hovy Gradient Flow sequence model for LID designed to support dialectal and multilingual language and Spruit, 2016; Blodgett et al., 2016). Legend Violence is a serious problems for cities like Chicago and has been exacerbated by the use of so-Token Alianment 1: 1

cial media by gang-involved youths for taunting rival gangs. We present a corpus of tweets from a young and powerful female gang member and her communicators, which we have annotated

varieties. Our model achieves state-of-theart performance on multiple LID benchmarks. Furthermore, in a case study us-

Multiple systems have been proposed for broadcoverage LID at the global level (McCandless, 2010; Lui and Baldwin, 2012; Brown, 2014; Jaech

Frozen

TADA

Critic

Critic Network





Flexible Boundaries

=> no highly accurate dialect identification systems available

□ Vary Depending on Personal and Social Contexts

=> dialects do not neatly fit into predefined categories



Accommodate the diversity of dialects from a Fine-Grained perspective Linguistic Features

# Method

# Dialect Adaptation via Dynamic Aggregation



#### 🏅 Modular and Dynamic

- Multi-Dialectal Robustness
  - Input Dialect-Agnostic
  - Personal- and Social-Contextual

#### 🔆 Within Only 3 Steps

- 1. Synthetic Datasets Construction
- 2. Feature Adapter Training
- 3. Dynamic Aggregation



### Step 1: Synthetic Datasets Construction



🔖 Construct a transformed dataset for each non-standard (morphosyntactic) linguistic feature.



### Step 2: Feature Adapter Training





Train a feature adapter for each non-standard linguistic feature.





# Experiments

# 1/ DADA Can Improve Multi-Dialectal Robustness



Adapt SAE model to multiple dialect variants simultaneously: AppE, ChcE, CollSgE, IndE, AAVE



### 2/ DADA Can Be Task-Agnostic



Adapt instruction-tuned **SAE** model to the dialect variants for multiple tasks



### 3/ DADA Has Great Interpretability!



#### Correlation Coefficients for **AAVE** Adaptation



We use abbreviations for certain terms, such as "nc" for "negative concord."

# **Conclusion and Future Work**

### Dynamic Aggregation of Linguistic Features



- A Fine-grained and Modular Method for Dialect Adaptation
  - □ Improve Multi-Dialect and Multi-Task Robustness
    - □ No need for highly accurate dialect identification systems
    - Taking personal and social context into account
    - Applicable to task-agnostic instruction-tuned LLMs
  - □ Interpretability, reusability and extensibility

#### But!!!



Non-Standard Linguistic Features

- are curated by linguists (<u>eWAVE</u>) and
- □ play a crucial role in a wide range of applications.

However, the manual curation of linguistic rules can be **expensive** and **expertise-intensive**.

