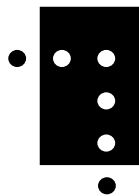


NB Uttale: A Norwegian Pronunciation Lexicon with Dialect Variation

Marie Iversdatter Røsok & Ingerid Løyning Dale

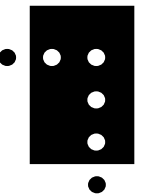
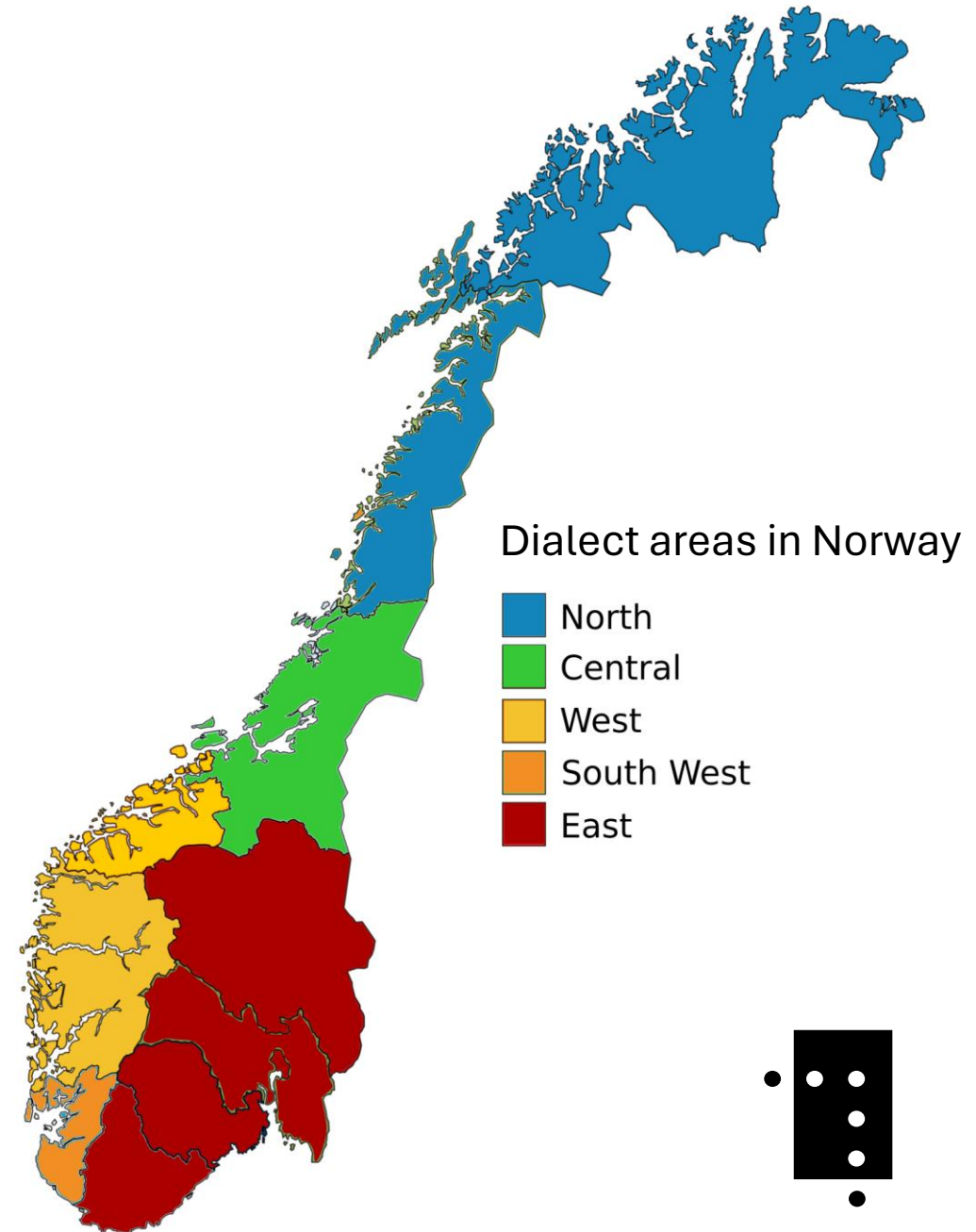
The Language Bank

The National Library of Norway



NB Uttale

- Nearly 800 000 word entries for Norwegian Bokmål
- Phonemic transcriptions for five dialect areas
 - East (e)
 - West (w)
 - Southwest (sw)
 - Central (t)
 - North (n)
- Close-to-spoken and close-to-written style variants



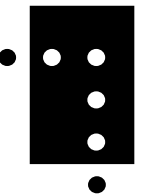
```

# Replace AX0 N AX0 adjective ending with AH0 N AX0 S
# S P EH2 N NX0 AX0 --> S P EH2 N AH0 N AX0 S
dialect_anes_adjective = {
    'name': 'dialect_anes_adjective',
    'areas': [
        'n_spoken',
        'sw_spoken'],
    'rules': [
        {
            'pattern': r'\b(AX0 N|NX0) AX0( S)?$',
            'replacement': r'AH0 N AX0 S',
            'constraints': [
                {
                    "field": 'pos',
                    "pattern": r'JJ|VB',
                    "is_regex": True
                },
                {
                    "field": 'wordform',
                    "pattern": r'endes?$' ,
                    "is_regex": True
                }
            ],
        },
    ],
}

```

Transformation rules

- Based on NST pronunciation lexicon (*Nordisk språkteknologi*)
- Phonemic differences between dialects
- Nofabet transcription standard
- Regular expressions
- Constraints



Error corrections

- Syllabic retroflex nasal

bakeren (the baker), /'bɑ:̥.kæ.ŋ/ → /'bɑ:̥.kæŋ/

Dialect updates

- *Non-retroflexes (west, southwest):*

forsker (scientist), /fɔ̥ʂkər/ → /fɔ̥r**ʂ**kər/

rs, rn, rl, rd, rt:

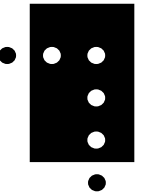
/ʂ ŋ | d t/ → /rs rn rl rd rt/

- *Apocope (central, north):*

elvene (the rivers), /ɛlʋə**nə**/ → /ɛlʋ**an**/

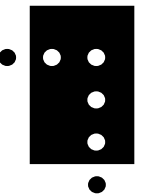
- *Hard d (west):*

kveld (night), /kʋɛl/ → /kʋɛl**d**/



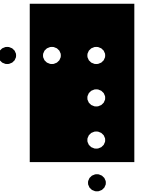
Evaluation

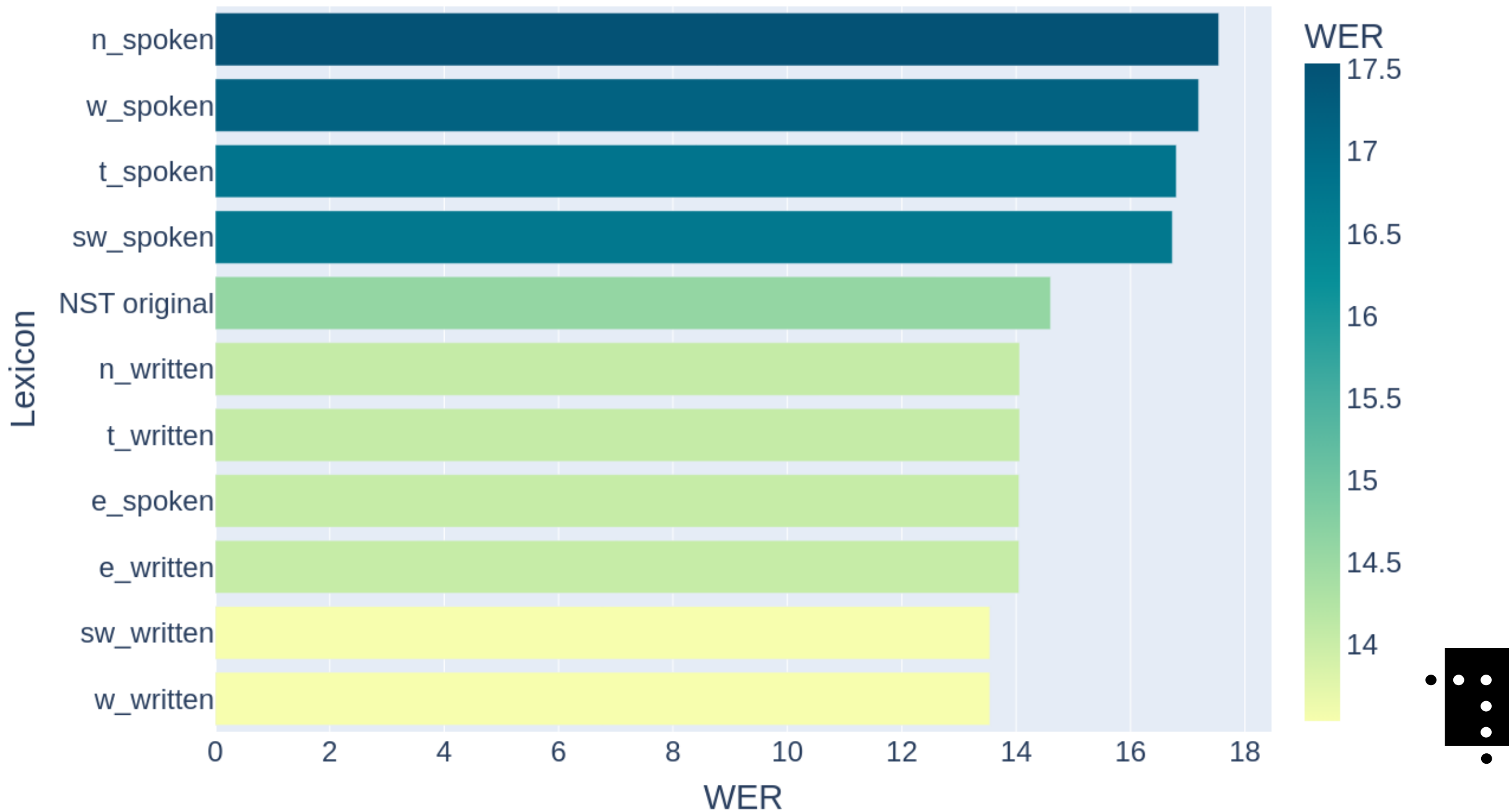
- Consistency
- Correctness
- Statistical predictability as quantitative measure
- Grapheme-to-phoneme models

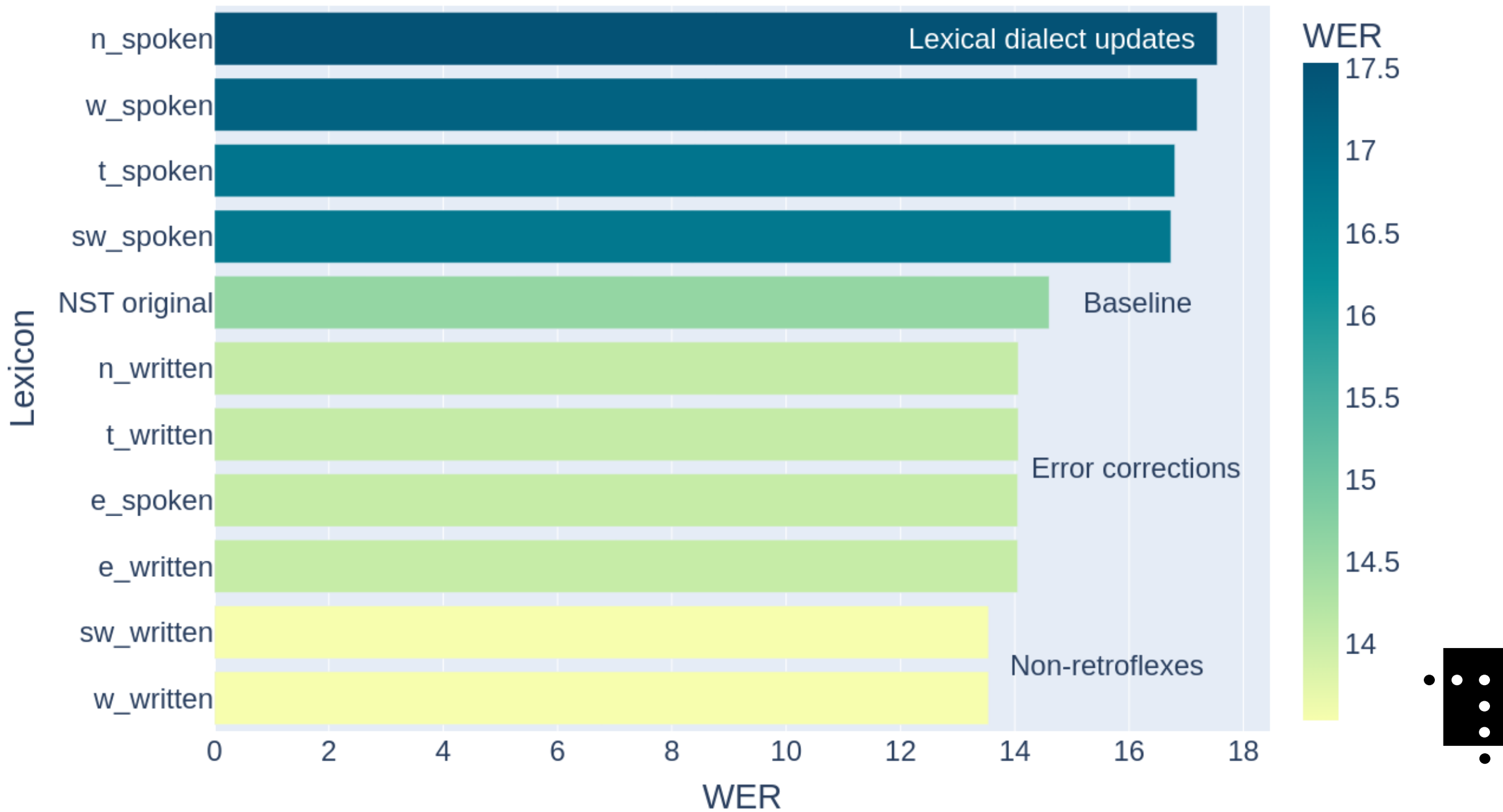


Grapheme-to-phoneme models (G2P)

- 1 per dialect + pronunciation style
- Phonetisaurus + Open FST
 1. Aligned grapheme and phoneme sequences (aligned corpus)
 2. 8-gram models
 3. Finite state transducer models







Thank you!

sprakbanken@nb.no



Funded by
The Research
Council of Norway

This work has in part been supported by
the IKTPLUS grant for the SCRIBE project
financed by the Research Council of
Norway (KSP2IPD).

