# GLAD: Groningen Lightweight Authorship Detection

E. van den Berg, M. Hürlimann, B. Weck, S. Šuster, M. Nissim
e.m.van.den.berg.4@student.rug.nl

**university of groningen**

## Authorship Verification...

Is this an instance of the class **SAME-AUTHOR** or **DIFFERENT-AUTHOR**?

| KNOWN AUTHOR | UNKNOWN AUTHOR |
|---|---|
| Jim, you won't go fightin' in the sun, with the birds all callin'? | You love me, don't you, David? You do, don't you, David? Tell me! |
| I'm not lookin' for it. Daisy, I love you. | I'm your husband, Annie, and you're my wife. Could there be aught but love between us after all these years? |
| And I love you, Jim. I don't want nothin' more than you in all the world. | |

## ... as classification



K1 K2 K3 K4   U → vector1(K1-K4,U),Y

K1 K2   U → vector2(K1-K2,U),N

K1 K2 K3   U → vector3(K1-K3,U),N

training set (#instances = #sets)

## What characterizes authors?

Gentle Tony → Unusual word choice?

Fat Vinny → Shorter sentences?

The Weasel → More complex grammar?

## Sidenote: Individual vs Joint

Individual          Individual          Joint

Vector_K(feat1,feat2) - Vector_U(feat1,feat2) = Vector_Joint(feat1,feat2)

## Sidenote: visual features

... incl. punctuation, letter case, line length
... to detect lay-out differences in plays / poetry

Pa-pa, pa-pa, pa-pa!

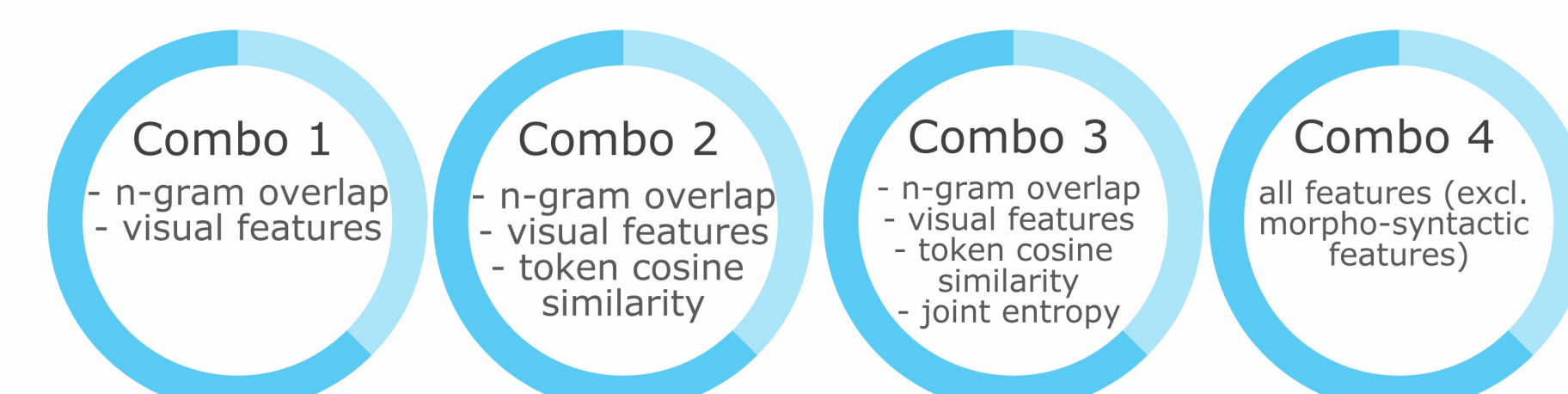Here, stop her. She'll fall down. Here, turn around. Walk this way.

Ma-ma, ma-ma, ma-ma;

Oh, I think you are a darling.

Mer-ry Christ-mas! Mer-ry Christmas.

## All features and combinations of features



**All Individual**
- Sentence length of known
- sentence length of unknown
- entropy of known
- entropy of unknown

entropy of known

entropy of unknown

sentence length of known

sentence length of unknown

**Visual features + Compression dissimilarity**

compression dissimilarity

visual features

joint entropy

n-gram overlap

morpho-syntactic similariy

token cosine similarity

**All Joint**
- n-gram overlap
- token cosine similarity
- joint entropy
- morpho-syntactic similarity
- compression dissimilarity
- visual feature

**Visual features + n-gram overlap + token cosine similarity**

## Features' performance

| Language | Helpful | Harmful |
|---|---|---|
| Dutch | All Individual | All Joint |
| Greek | All Joint | Just Entropy |
| Spanish | All Joint | Just Ngram sim |
| English | Vis+ngram+tok | Just POS sim |

## Best feature combinations

**Combo 1**
- n-gram overlap
- visual features

**Combo 2**
- n-gram overlap
- visual features
- token cosine similarity

**Combo 3**
- n-gram overlap
- visual features
- token cosine similarity
- joint entropy

**Combo 4**
all features (excl. morpho-syntactic features)

## Results: combined scores on PAN 2015 data

| Language | Training | Test | Ranking |
|---|---|---|---|
| Dutch (`full set`) | .55 | .62 | **3**/17 |
| English (`full set`) | .56 | .41 | **8**/17 |
| Greek (`combo2`) | .54 | .60 | **4**/15 |
| Spanish (`full set`) | .90 | .54 | **5**/17 |

**Average runtime:** 1 minute

## Discussion

- Why do complex features not outperform surface-level features?

- Are **individual** vs **joint** features valid ways of classifying texts written by the same of different authors?

## Conclusion: Pros of the system

**Robust**
... works across genres
... works across topics
... works across languages

**Simple**
... easy to run
... uses simple features
... very fast runtime