**Tarek Sakakini**     **Suma Bhat**     **Pramod Viswanath**

# MORSE:
# Semantic-ally Drive-n MORpheme SEgment-er

**Samuel MORSE** minimized the number of on-off clicks for non-verbal communication.

**This MORSE** minimizes the vocabulary size for Natural Language Processing systems.
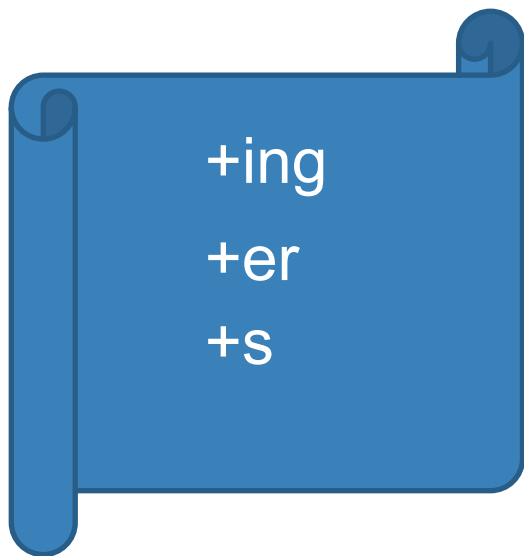
**1** Morpheme Segmentation

Hopefully

## Not a trivial task

+ing
+er
+s

Players

Playing

Beijing

Butterflies

# Applications

**Model:**

| Quickly | Sad |
|---|---|
| •Rapidement | •Triste |

Sadly

↓

**Test:**

???

**Model:**

| Quick | ly | Sad |
|---|---|---|
| • Rapide | • ment | • Triste |

Sadly

↓

**Test:**

Tristement

# 📌 Applications

**Google** | cheap car 🔍

Here at Toyota World, we have the **cheap**est **car**s in town.
We are proudly called the first and last stop.
…
…

# 2 Previous Work

# Helplessly

# Morfessor (Creutz and Lagos, 2005)

Help: 2387

Helping: 1586

Helper: 498

Helps: 2437

Jump: 1847

Jumping: 1664

Jumper: 1290

Jumps: 2987

**Locally Semantic**

**Cosine similarity**

car      🔴      caring

car      🟢      cars

**(Schone and Jurafsky, 2000)**
**(Narasimhan et al., 2015)**
**(Luo et al., 2017)**

# MORSE

3

**Input:**
Word Embeddings

Unsupervised
Morphology Learning

**Segmentation:**
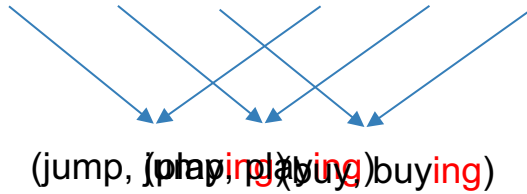Optimization Problem

4 hyperparameters:
Small tuning dataset

# Step 1

Learning Morphology

**Collecting candidate morphological rules**

Vocabulary: jump play buy jumping playing buying jumper player buyer ….. and stand

(jump, jumping) (play, playing) (buy, buying)

(and, stand)

| **(suf, ∅, ing):** | (jump, jumping) | (play, playing) | (buy, buying) |
| **(suf, ∅, er):** | (jump, jumper) | (play, player) | (buy, buyer) |
| **(pre, ∅, st):** | (and, stand) | (ore, store) | (one, stone) |

# Signals

## Orthographic

quick       quick**ly**

beautiful      beautiful**ly**

confident      confident**ly**

wrong       wrong**ly**

## Semantic

**Word Embeddings**

quick → quick**ly**

beautiful → beautiful**ly**

wrong → wrong**ly**

confident → confident**ly**

# What makes a good rule?

**Signal 1: Orthography**

**Rule = (suf, ∅, ly)**      **Size = 8723**

(quick, quickly)     (beautiful, beautifully)

(confident, confidently)

…………………………………………

……………………… (wrong, wrongly)

**Rule = (pre, ∅, st)**      **Size= 16**

(ore, store) ……

(amp, stamp)

# What makes a good rule?

quick → quickly

wrong → wrongly

confident → confidently

beautiful → beautifully

one ↑ stone

amp → stamp

ore → store

and → stand

**What makes a good member of a rule?**

quick → quickly ✔

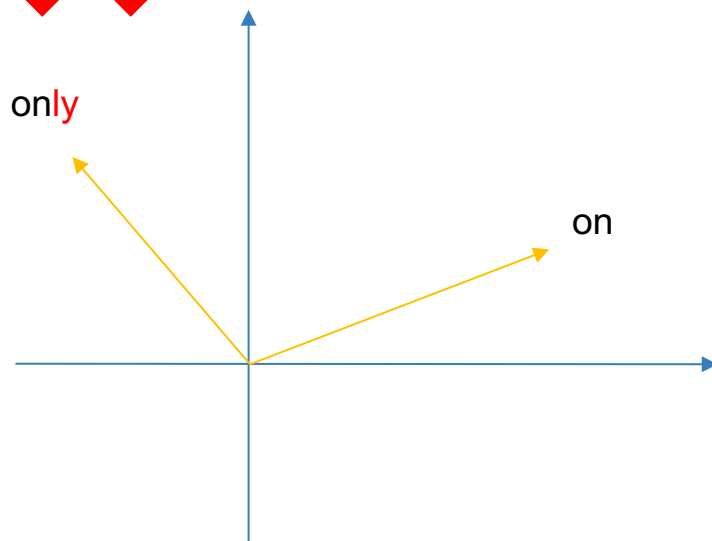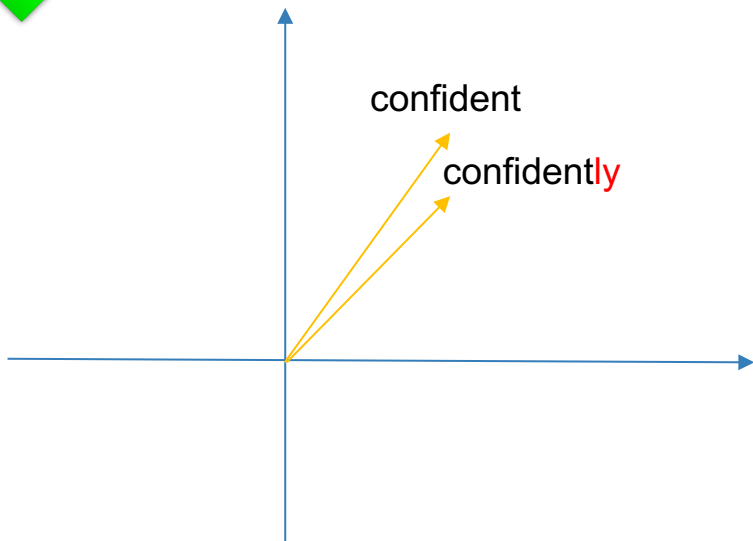confident → confidently ✔

wrong → wrongly ✔

on → only ✘

beautiful → beautifully ✔

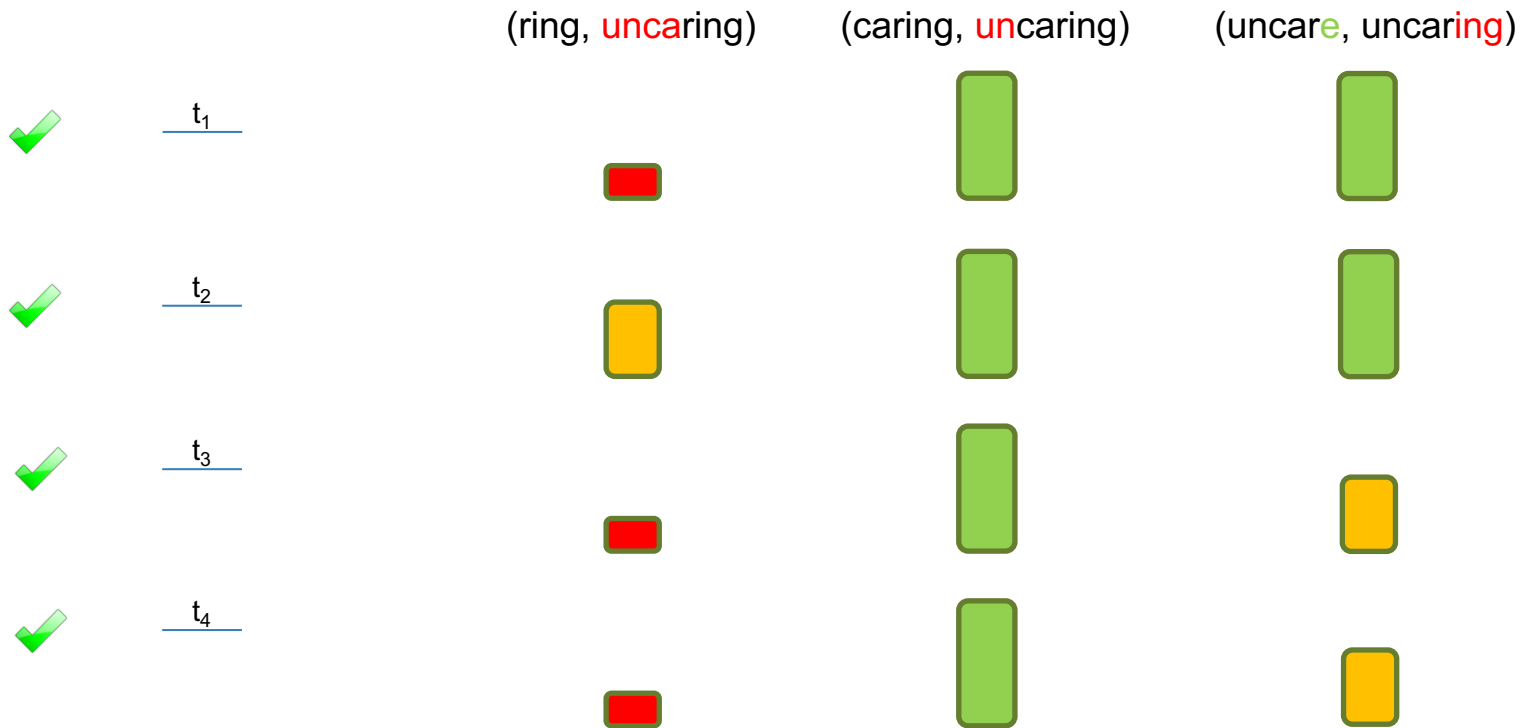# What makes a good member of a rule?

**Scope: Local**

# Step 2

Segmenting

# Linear Optimization Problem

(ring, uncaring)    (caring, uncaring)    (uncare, uncaring)

uncaring

un + care + ing

**Iterate**

care

$t_1$

$t_2$

$t_3$

$t_4$

(car, car**e**)  (ca, ca**re**)  (re, **ca**re)

# 4　Experiments

# Experimental Setup

**Training**

**Languages**

**Gold Datasets**

Morpho Challenge

| jumping | jump | ing |
| playing | play | ing |
| jumps | jump | s |
| calls | call | s |
| rooms | room | s |

# Experiments

## Morpho Challenge downsides

Non-compositional

Trivial instances

Human error

Business

Turning-point        Player's

Turning

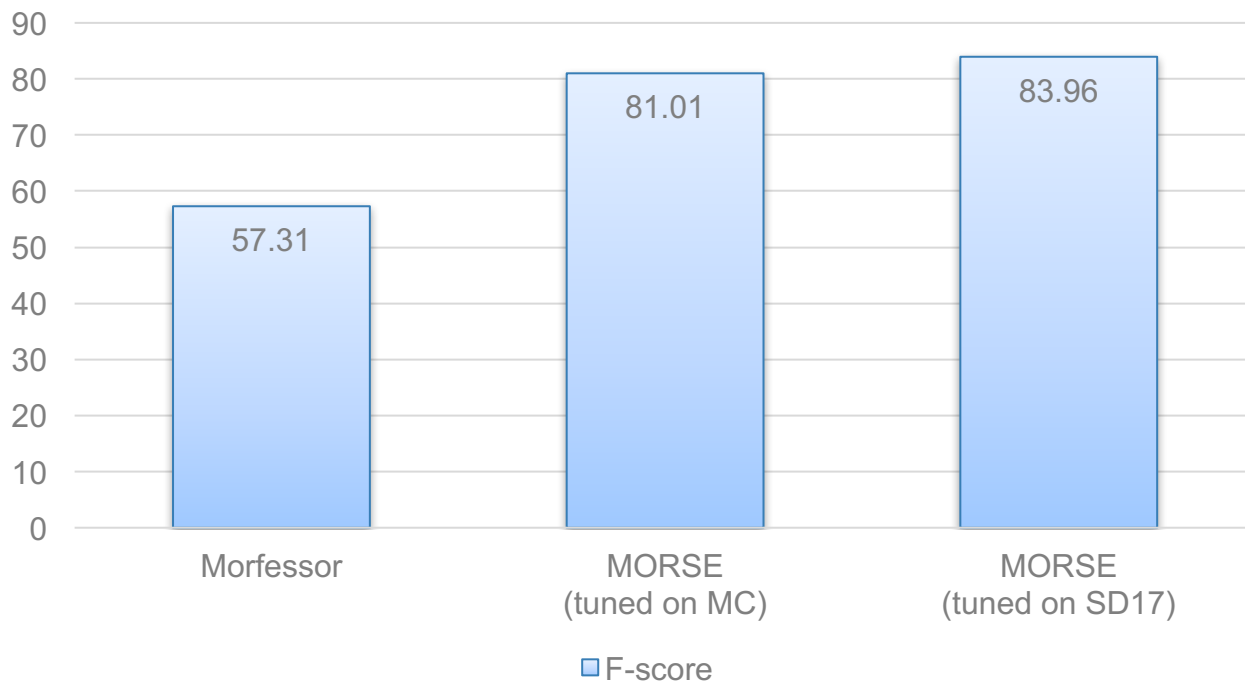# Experiments

⦿ 2000 words

⦿ Compositional

⦿ 91% inter-annotator agreement

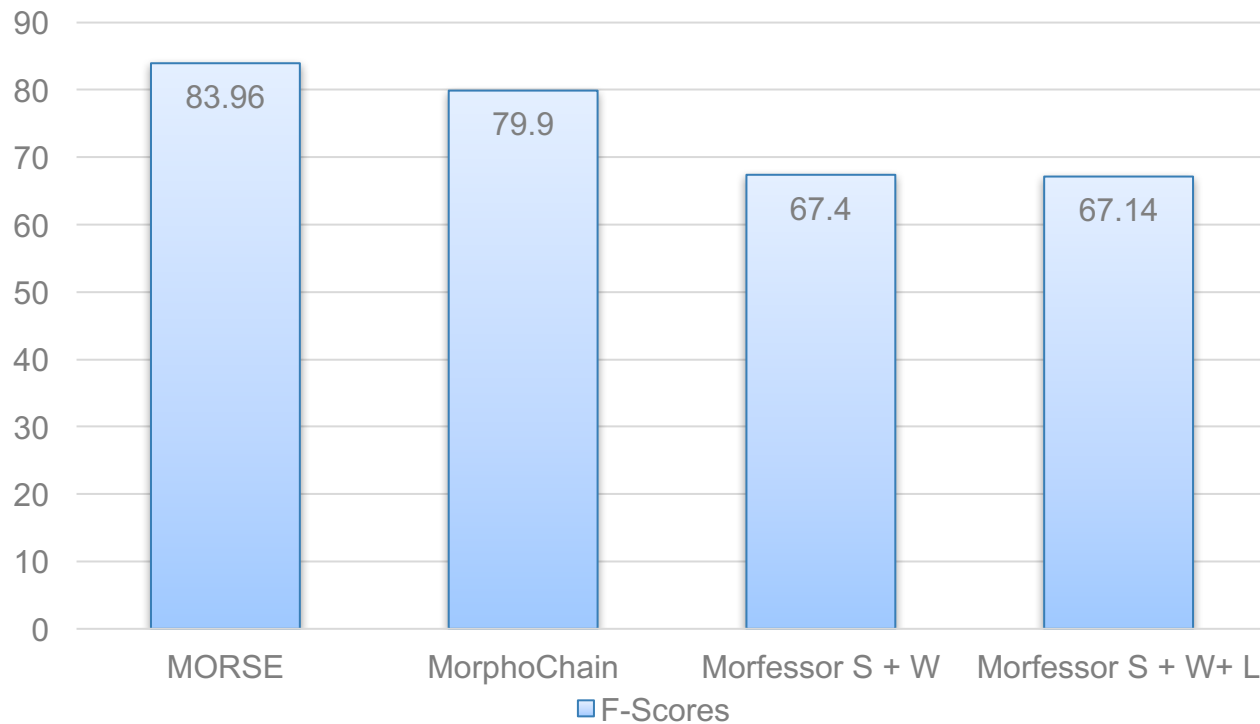⦿ In canonical (butterfly + ies) and non-canonical version (butterfl + ies)

# Results on SD17
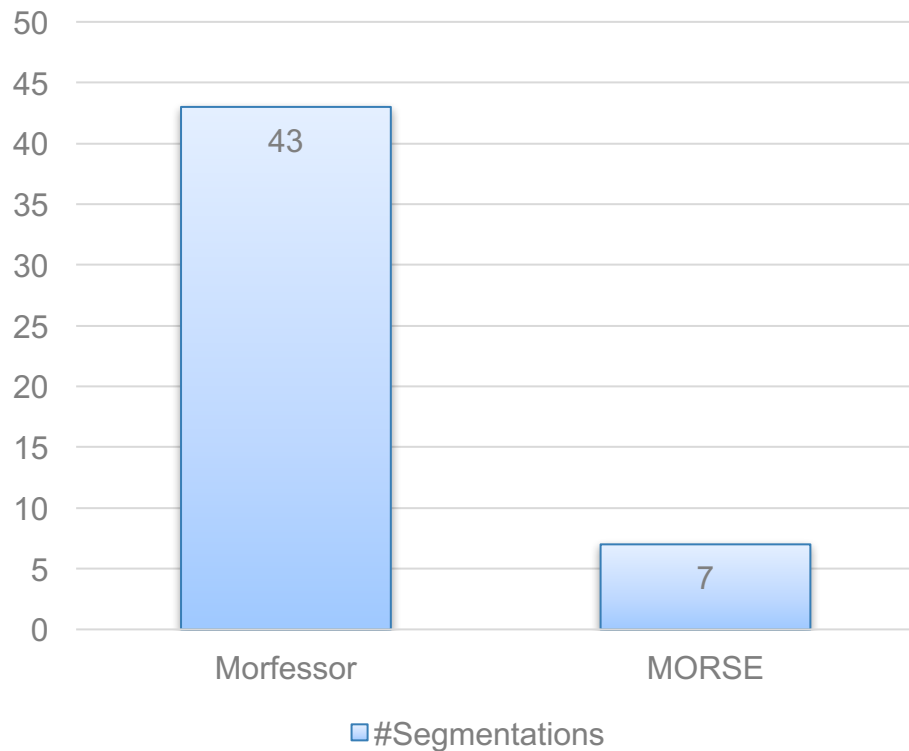
# **Negative Dataset**

◉ 100 words like: honeymoon, passport, outdoors

◉ Checks for robustness

**Looking forward**

- Robustness to highly agglutinative languages

- Extending to other languages (non-concatenative)

katAba

# Looking forward

◉ Morphological mappings across languages



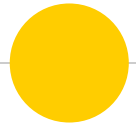| English | | French |
|---------|---|--------|
| (suf, ∅, ly) | → | (suf, ∅, ment) |
| (suf, ∅, s) | → | (suf, ∅, s) |
| | → | (suf, ∅, es) |

# Links



https://morse.mybluemix.net



https://github.com/yoonlee95/morse_segmentation

# Thank you

Questions?

# Effect of Hyperparameters



Recall



Precision

# Prerequisite

**Valid rule** with an **invalid instance**
(suf, ∅, ing)          (s, sing)

**Invalid rule**
(pre, ∅, s)

playing

play

jumping

jump

screaming

scream

sing

s

mile          smile

cream

store

tore

scream

slay

lay

# Demo

4

morse.mybluemix.net