# Unsupervised Morphology Induction using Word Embeddings

Radu Soricut, Franz J. Och*

NAACL 2015

*now at Human Longevity Inc.

# Word Embeddings

- vocabulary $V$, embedding function $e: V \to \mathbf{R}^n$

- vector space encodes semantic similarity

  - $e$(car) $\simeq$ $e$(automobile), $e$(car) $\neq$ $e$(seahorse)

- vector space encodes compositionality

  - semantic: $e$(king) - $e$(man) + $e$(woman) $\simeq$ $e$(queen)

  - syntactic: e(cars) - e(car) + e(fireman) $\simeq$ e(firemen)

- vector space encodes syntactic/semantic transformations
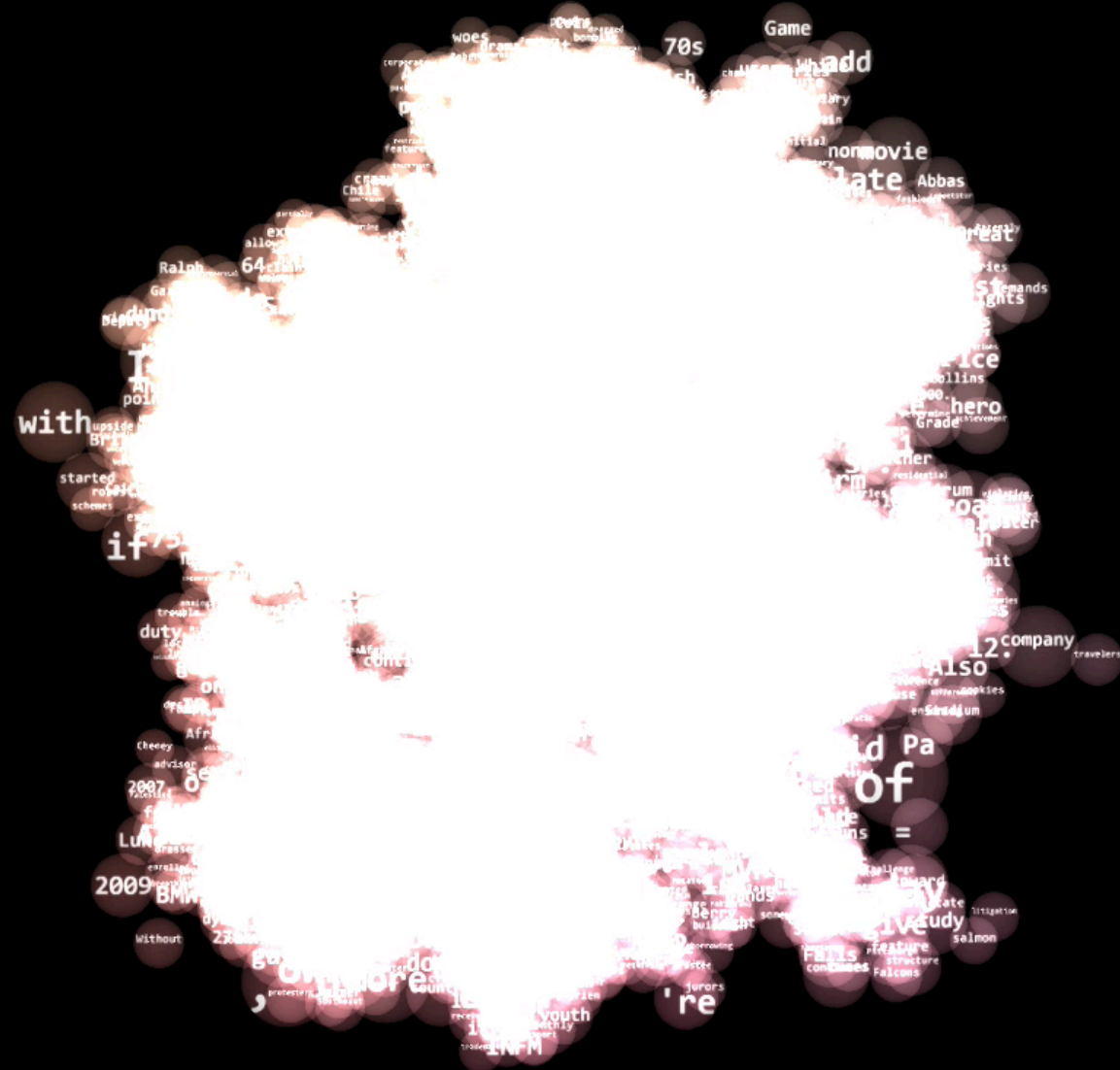
  - anti+ $\simeq$ $e$(anticoruption) - $e$(corruption)

*British scientists recreated Down 's syndrome in mice in order to study the disease and develop new treatments .*

# SkipGram Embeddings [Mikolov et al., 2013]

*British scientists recreated Down 's syndrome in mice in order to study the disease and develop new treatments .*

$$\arg\max_{v_w, v_c}\Big(\sum_{(w,c)\in D_w}\log\sigma(v_w\cdot v_c)+\sum_{(w,c)\in\overline{D_w}}\log\sigma(-v_w\cdot v_c)\Big)$$

$$\sigma(x)=\frac{1}{1+e^x}$$

Train on large monolingual corpus

size: 8,869 tupl

king - man + w

evaluation: acc
(over entire voc

size: 10,675 tuples

cars - car + fireman ≃ firemen

evaluation: acc% closest poin
(over entire vocabulary)

size: 2003 pairs

is a celebrated
an English ?roc

eval: Spearman
(against 10-hun

size: 2034 pairs

belligerence ≃ hostility   8.7
amorphous ≃ inorganic  1.9
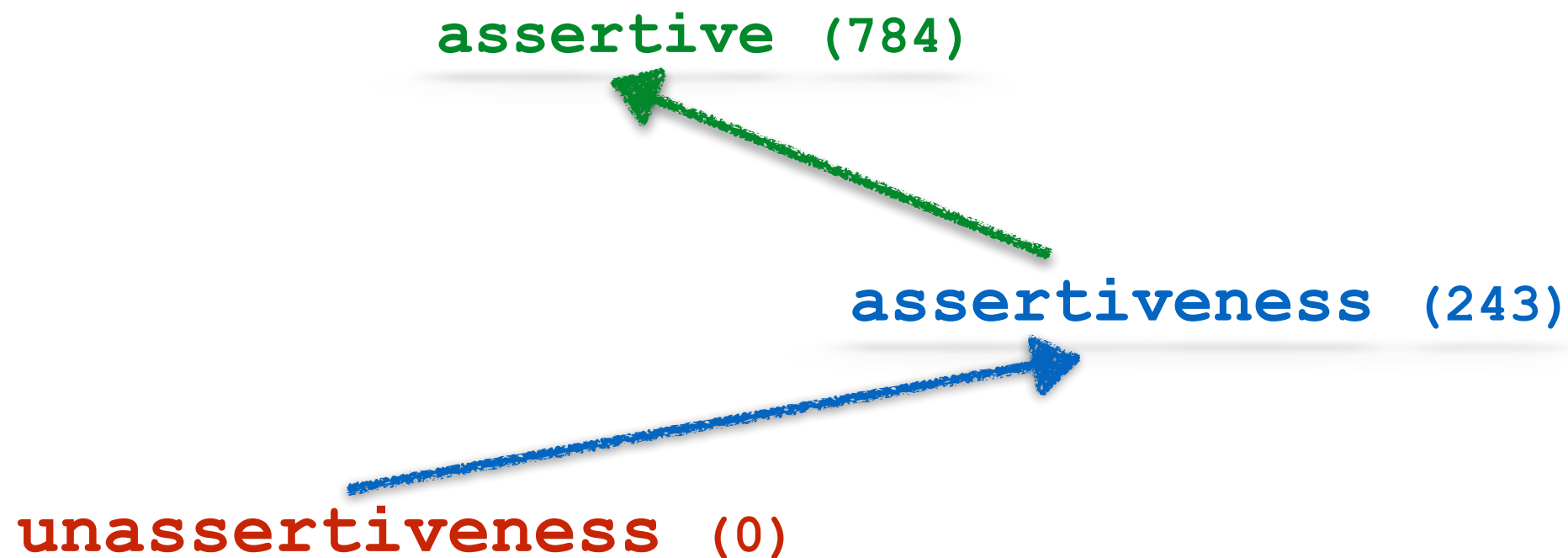
eval: Spearman corr.
(against 10-human score avg.)

| Model | Compositionality (Acc%) | | Similarity (Spearman ρ) | |
|---|---|---|---|---|
| | Semantic | Syntactic | Stanford-C | Stanford-RW |
| SkipGram (Wikipedia 1Bw) | 76.7 | 68.3 | 66.3 | 35.8 |

$$\arg\max_{v_w, v_c}\left( \sum_{(w,c)\in D_w} \log \sigma(v_w \cdot v_c) + \sum_{(w,c)\in \overline{D_w}} \log \sigma(-v_w \cdot v_c) \right)$$
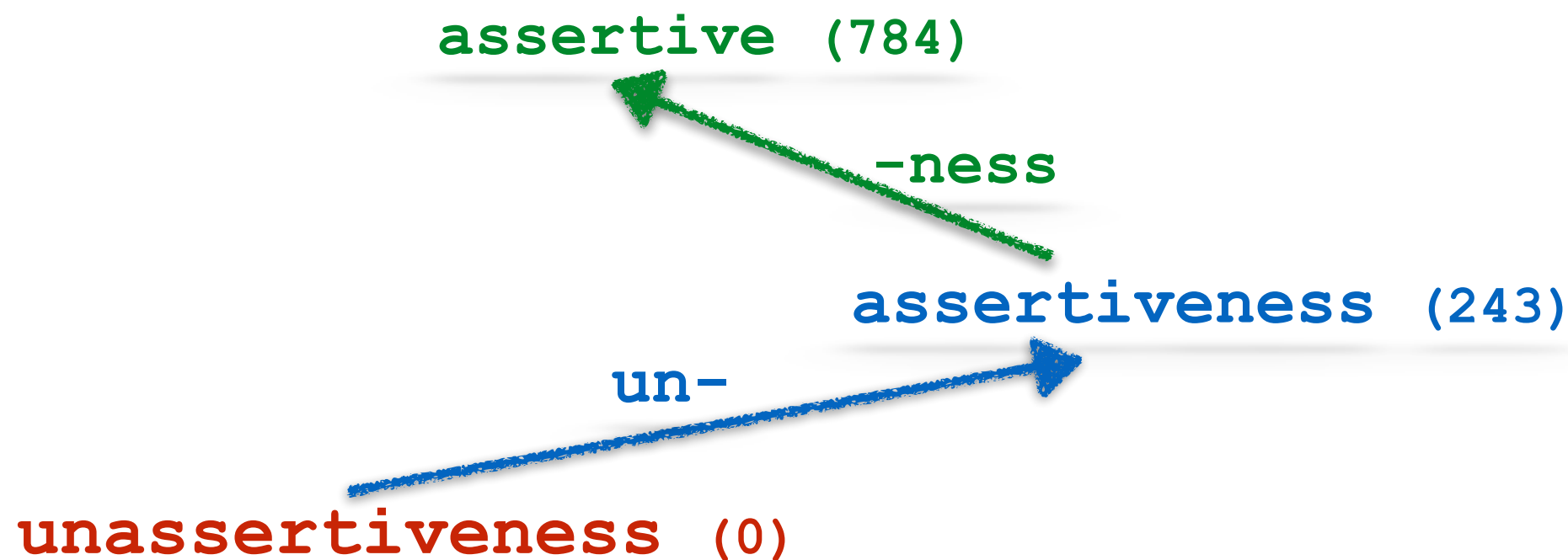
$$\sigma(x) = \frac{1}{1+e^x}$$

Q: What do we want?

A: We want *high-quality* embeddings for all words (even ones outside *V*)

**assertive** **(784)**

**assertiveness** **(243)**

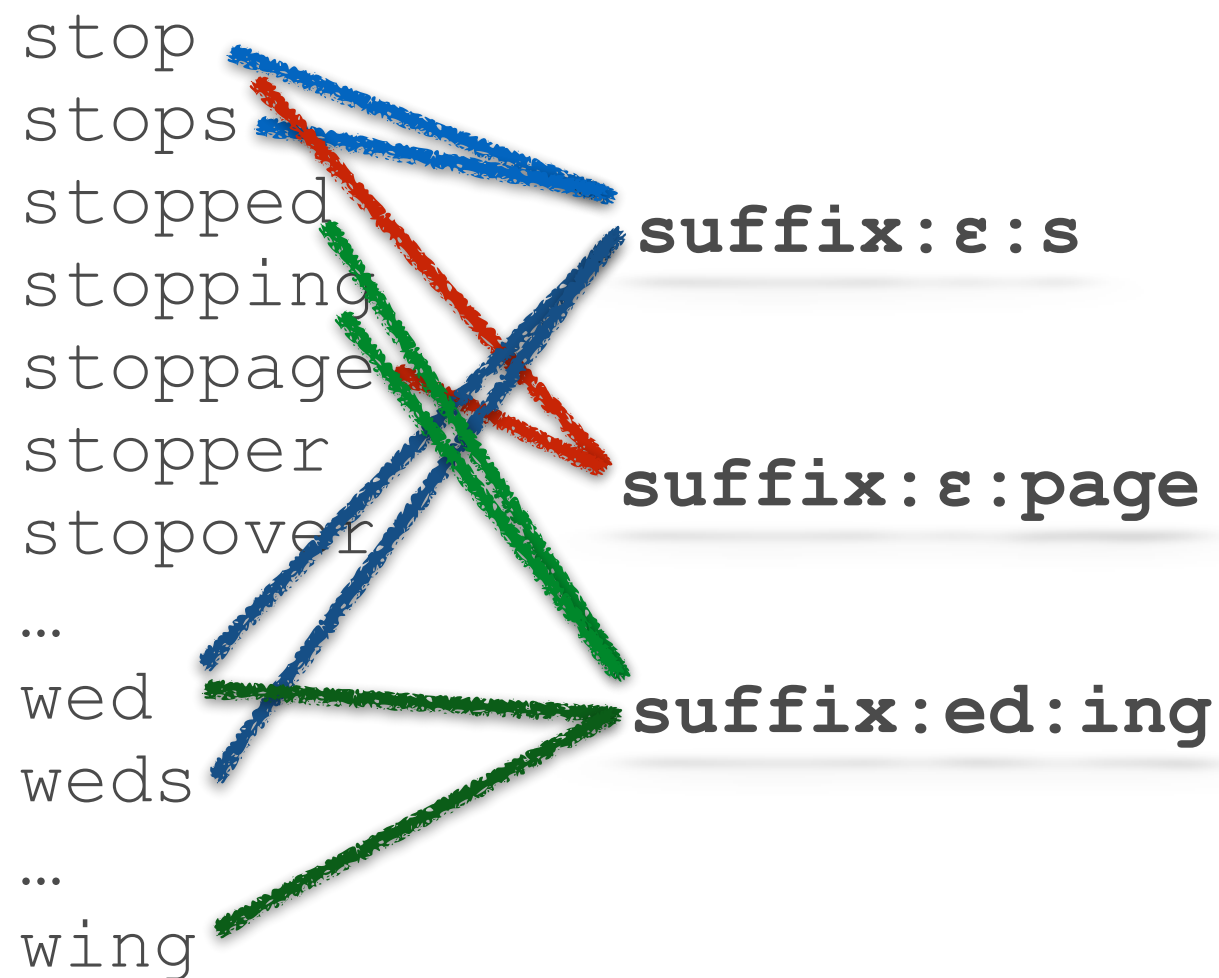**unassertiveness** **(0)**

# Unsupervised Morphology Induction

Q: What do we want?

A: We want *morphology-based transformations* that can accurately analyze words (even ones unseen at training time)

**assertive** (784)

**-ness**

**assertiveness** (243)

**un–**

**unassertiveness** (0)

Steps:

1. From *V*, extract candidates for morphological rules (prefix & suffix only)

```
stop
stops
stopped
stopping
stoppage
stopper
stopover
…
wed
weds
…
wing
```

**suffix:ε:s**

**suffix:ε:page**

**suffix:ed:ing**

Steps:

1. From *V*, extract candidates for morphological rules (prefix & suffix only)



stop
stops
stopped
stopping
stoppage
stopper
stopover
…
wed
weds
…
wing

**suffix:ε:s**

**suffix:ε:page**

**suffix:ed:ing**

Screams
Screw
…
aware
creams
crew
done
…
truthful
unaware
undone
untruthful

**prefix:S:ε**

**prefix:un:ε**

Steps:

2. Query against embedding space: *morphology does not shift meaning*

**suffix:ed:ing**

adored adorned affected …
blamed blitzed blogged …
stayed stepped stopped …
weaned wed wedged whirled

$rank($blamed → blam**ing**$) = 1$
$rank($stopped → stopp**ing**$) = 2$
$rank($wed → w**ing**$) = 28609$

**prefix:ε:S**

aura aux ave …
canned cans car care …
crape cream creams …
miles mitten mothers …

$rank($care → Scare$) = 57778$
$rank($cream → Scream$) = 9434$
$rank($miles → Smiles$) = 18800$

Steps:

2. Query against embedding space: *morphology does not shift meaning*

**prefix:un:ε**

unabated unable unabridged…
unaware unbalance unbeaten…
undoing undone undoubted…
untrusted untrustworthy…

*rank*(unaware → aware) = 1
*rank*(undone → done) = 129

Steps:

2. Query against embedding space: *morphology does not shift meaning*

*morphology shifts meaning consistently*

**prefix:un:ε**

unabated unable unabridged…
unaware unbalance unbeaten…
undoing undone undoubted…
untrusted untrustworthy…

**↑un–**

clear – unclear
delivered – undelivered
truthful – untruthful

*rank*(unaware → aware) = 0
*rank*(undone → done) = 129

*rank*(undone + ↑un– → done) = 4

Steps:

3. Extract candidate rules using embedding-based stats

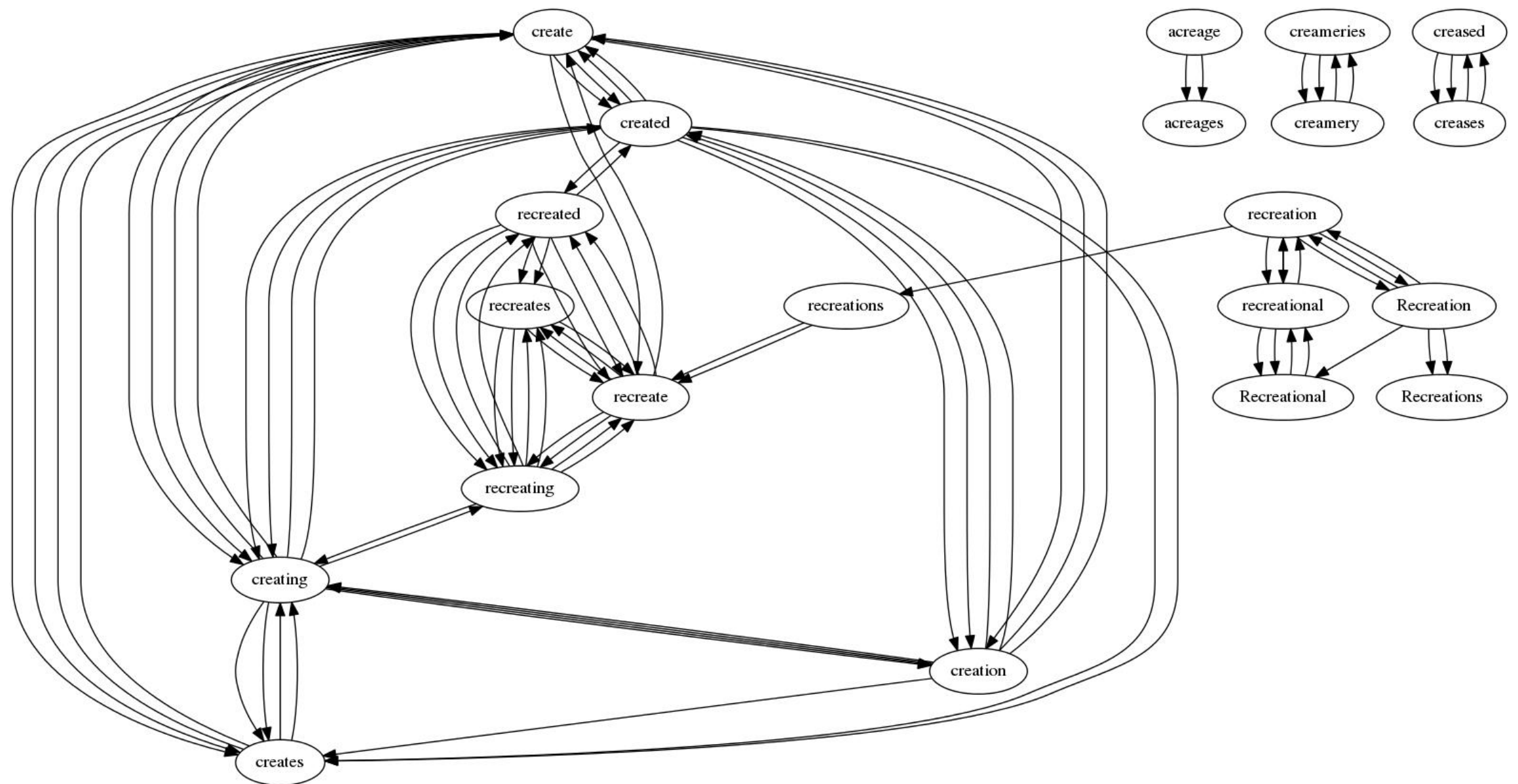| Candidate Rule | Direction | #Correct | #Total | Acc10 |
|---|---|---|---|---|
| **suffix:h:a** | ↑Teh | 1 | 449 | 0.4% |
| **suffix:o:es** | ↑Tono | 7 | 688 | 1.0% |
| **prefix:D:W** | ↑Daring | 9 | 675 | 1.3% |
| … | | | | |
| **prefix:un:ε** | ↑undelivered | 166 | 994 | 23.3% |
| **suffix:ed:ing** | ↑procured | 2138 | 4714 | 56.2% |
| … | | | | |
| **suffix:ating:ate** | ↑formulating | 255 | 395 | 74.7% |
| **suffix:sed:zed** | ↑victimised | 153 | 186 | 90.9% |

**Bad**

**Good**

Steps:

4. Use rules to extract lexicalized, weighted morphological transformations

| Start | Rule + Direction = Transformation | End | Cosine | Rank |
|---|---|---|---|---|
| ... | | | | |
| recreat**ions** | **suffix:ions:e + ↑investigations** | recreat**e** | 0.69 | 1 |
| recrea**tions** | **suffix:tions:te + ↑investigations** | recrea**te** | 0.70 | 1 |
| recreat**ions** | **suffix:ions:ed + ↑delineations** | recreat**ed** | 0.51 | 29 |
| recreat**ions** | **suffix:ions:ing + ↑reconstruction** | recreat**ing** | 0.72 | 1 |
| ... | | | | |
| **un**aware | **prefix:un:ε + ↑uncivilized** | aware | 0.77 | 1 |
| **un**aware | **prefix:un:ε + ↑undelivered** | aware | 0.63 | 7 |

Output (I): labeled, weighted, cyclic, directed multigraph $G^V_{Morph}$

- words are nodes, morphological transformations are (weighted) edges

Output (II): labeled, weighted, acyclic, directed graph $D^V_{Morph}$

- words are nodes, morphological mappings are weighted edges



fix-point node

Q: What do we want?

A: We want *morphology-based transformations* that can accurately analyze words (even ones unseen at training time)

**assertive** (784)

**-ness**

**assertiveness** (243)

**un-**

**unassertiveness** (0)

Basic algorithm: embedding words outside *V*

Outside Wikipedia (1B tokens, $|V| = 4.3M$)

**animalize** (0)

**balminess** (0)

**caesarism** (0)

**containerful** (0)

**nonindulgent** (0)

**unassertiveness** (0)

# Unsupervised Morphology Induction: Algorithm

Analyze words outside *V*

1. Train time: extract and count all paths ending in a "fix-point" from the directed acyclic graph $D^V_{Morph}$

   - each path is called a "rule sequence"



fix-point node

| rule sequence | count |
|---:|:---|
| **suffix:s:ε** | 3119 |
| **suffix:ed:ε** | 687 |
| **suffix:ing:ed** | 412 |
| **prefix:un:ε** | 207 |
| **suffix:ness:ε** | 162 |
| **suffix:ness:ly** | 25 |
| **suffix:y:ier,suffix:er:ness** | 10 |
| **prefix:un:ε,suffix:ed:ing** | 5 |

Analyze words outside *V*

2. Run time: apply each rule sequence in descending order of counts
   - if rule fires, check that result has count > 0 and in-degree > 0
   - stop at first winner

| | rule sequence | count | |
|---|---|---|---|
| **unassertiveness(0)** → | **suffix:s:ε** | 3119 | → **unassertivenes(0)** |
| | **suffix:ed:ε** | 687 | |
| | **suffix:ing:ed** | 412 | |
| **unassertiveness(0)** → | **prefix:un:ε** | 207 | → **assertiveness(243)** |
| | **suffix:ness:ε** | 162 | |
| | **suffix:ness:ly** | 25 | |
| | **suffix:y:ier,suffix:er:ness** | 10 | |
| | **prefix:un:ε,suffix:ed:ing** | 5 | |

**unassertiveness = assertiveness + ↑un+**

A: We want *morphology-based transformations* that can accurately analyze words unseen at training time

**assertiveness (243)**

**un-**

**unassertiveness (0)**

| Language | \|Tokens\| | \|V\| | $\|G^V_{Morph}\|$ | $\|D^V_{Morph}\|$ |
|---|---|---|---|---|
| EN | 1.1b | 1.2m | 780k | 75,823 |
| DE | 1.2b | 2.9m | 3.7m | 169,017 |
| FR | 1.5b | 1.2m | 1.8m | 92,145 |
| ES | 566m | 941k | 2.2m | 82,379 |
| RO | 1.7b | 963k | 3.8m | 141,642 |
| AR | 453m | 624k | 2.4m | 114,246 |
| UZ | 850m | 2.0m | 5.6m | 194,717 |

Evaluate OOV analysis using low-count words



~~**unassertiveness (0)**~~

**unassertive (240)**

**assertive (1840)**

**missive (472)**

**assert (15535)**

**asserted (19352)**

**miss (18113)**

**asserting (3087)**

Evaluate OOV analysis using rare words



| Language | $\|V_{[1000,2000)}\|$ | | Accuracy | |
|---|---|---|---|---|
| | Have analysis | Don't have analysis | Have analysis | Don't have analysis |
| EN | 3421 | 10617 | 89.7% | 89.6% |
| DE | 10778 | 21234 | 90.8% | 93.1% |
| FR | 6435 | 9807 | 90.3% | 90.4% |
| ES | 5724 | 7412 | 91.1% | 90.3% |
| RO | 11905 | 9254 | 86.5% | 85.3% |
| AR | 7913 | 5202 | 92.4% | 69.0% |
| UZ | 11772 | 9027 | 81.3% | 84.1% |

Q: What do we want?

A: We want *morphology-based transformations* that can accurately analyze words (even ones unseen at training time)

Key result

Improved word similarity judgment: unknown, low-count, high-count words

- evaluation on Stanford Rare Word similarity dataset (RW-EN)
- evaluation of similarity datasets on various languages (RG-DE)

Training Setup

| | Language | Train Set | \|Tokens\| | \|V\| | $\|G^V_{Morph}\|$ | $\|D^V_{Morph}\|$ |
|---|---|---|---|---|---|---|
| **Small** | EN | Wiki-EN | 1.1b | 1.2m | 780k | 75,823 |
| | DE | WMT-DE | 1.2b | 2.9m | 3.7m | 169,017 |
| **Large** | EN | News-EN | 120b | 1.0m | 2.9m | 98,268 |
| | DE | News-DE | 20b | 1.8m | 6.7m | 351,980 |

Evaluation on similarity datasets (RG-DE, RW-EN)

| Language | Train Set | \|Tokens\| | \|V\| | $\|G^V_{Morph}\|$ | $\|D^V_{Morph}\|$ |
|---|---|---|---|---|---|
| EN | Wiki-EN | 1.1b | 1.2m | 780k | 75,823 |
| DE | WMT-DE | 1.2b | 2.9m | 3.7m | 169,017 |
| EN | News-EN | 120b | 1.0m | 2.9m | 98,268 |
| DE | News-DE | 20b | 1.8m | 6.7m | 351,980 |

size: 2034 pairs

| | | |
|---|---|---|
| impossibilities | unattainableness | 8.8 |
| deregulating | liberation | 8.0 |
| baseness | unworthiness | 4.0 |
| transmigrating | born | 1.1 |

| | RW-EN Testset | | | |
|---|---|---|---|---|
| | \|Unembedded\| | | Spearman ρ | |
| System | Wiki-EN | News-EN | Wiki-EN | News-EN |
| SkipGram | 78 | 177 | 35.8 | 44.7 |
| SkipGram+Morph | 1 | 0 | 41.8 | 52.0 |

+9   +7

# Unsupervised Morphology Induction: Evaluation

Evaluation on similarity datasets (RG-DE, RW-EN)

size: 65 pairs

| | | |
|---|---|---|
| Edelstein | Juwel | 3.8 |
| Autogramm | Unterschrift | 3.5 |
| Irrenhaus | Friedhof | 0.3 |
| Kraftfahrzeug | Magier | 0.0 |

| Language | Train Set | \|Tokens\| | \|V\| | $\|G^V_{Morph}\|$ | $\|D^V_{Morph}\|$ |
|---|---|---|---|---|---|
| EN | Wiki-EN | 1.1b | 1.2m | 780k | 75,823 |
| DE | WMT-DE | 1.2b | 2.9m | 3.7m | 169,017 |
| EN | News-EN | 120b | 1.0m | 2.9m | 98,268 |
| DE | News-DE | 20b | 1.8m | 6.7m | 351,980 |

### RW-EN Testset

| | \|Unembedded\| | | Spearman ρ | |
|---|---|---|---|---|
| System | Wiki-EN | News-EN | Wiki-EN | News-EN |
| SkipGram | 80 | 177 | 35.8 | 44.7 |
| SkipGram+Morph | 1 | 0 | 41.8 | 52.0 |

+9   +7

### RG-DE Testset

| | \|Unembedded\| | | Spearman ρ | |
|---|---|---|---|---|
| System | WMT-DE | News-DE | WMT-DE | News-DE |
| SkipGram | 0 | 20 | 62.4 | 62.1 |
| SkipGram+Morph | 0 | 0 | 64.1 | 69.1 |

+0   +7

# Conslusions

1. Method for inducing morphological transformations between words

     • from scratch, unsupervised, language agnostic

2. Provides morphology-based structure over embedding spaces

3. Provides high-quality embeddings for out-of-vocabulary and low-count morphological variants

Google™

- Going beyond suffix & prefix morphology
  - nothing in the approach prevents from extending it

- Use it for improved Machine Translation
  - quick and painless morphological analysis on source side
  - generate morphological variants on target side (even new ones!)

- Use it for improved Information Retrieval

# Thank you