

# ASSESSING THE LIMITS OF THE DISTRIBUTIONAL HYPOTHESIS IN SEMANTIC SPACES

---

## TRAIT-BASED RELATIONAL KNOWLEDGE AND THE IMPACT OF CO-OCCURRENCES

Mark Anderson and Jose Camacho Collados



# TRAIT-BASED RELATIONAL KNOWLEDGE

## Traits/attributes

- Traits commonly associated with concepts.
- Split into types (colour, shape&size, material, components, tactile)
- E.g. The ubiquitous yellow banana.



## Datasets

- McRae<sup>1</sup>
- Norms<sup>2</sup>
- McRae (Spanish)

---

<sup>1</sup> Ken McRae, George S. Cree, Mark S. Seidenberg, and Chris McNorgan. 2005. *Semantic feature production norms for a large set of living and nonliving things*. Behavior Research Methods, 37:547–559.

<sup>2</sup>Barry Devereux, Lorraine K. Tyler, Jeroen Geertzen, and Billi Randall. 2014. *The centre for speech, language and the brain (CSLB) concept property norms*. Behavior Research Methods, 46:1119 – 1127

# EXAMPLE TABLE — McRAE

| trait type |              | N <sub>C</sub> |  | N <sub>T</sub> |  |
|------------|--------------|----------------|--|----------------|--|
|            | colour       | 148            |  | 7              | green (32), brown (32), black (24), white (21), red (16), yellow (13), orange (10) |
|            | components   | 110            |  | 6              | handle (39), legs (19), wheels (14), leaves (14), seeds (13), doors (11)           |
|            | materials    | 144            |  | 4              | metal (79), wood (43), cotton (11), leather (11)                                   |
|            | size & shape | 234            |  | 4              | small (83), large (70), long (44), round (37)                                      |
|            | tactile      | 117            |  | 7              | heavy (21), soft (19), furry (18), sharp (17), hard (16), juicy (16), slimy (10)   |

# Co-OCCURENCES

- Sentence — e.g. *The **bananas** were ubiquitous in that town, and **yellow** jackets were the outerwear of choice for the inhabitants.*

---
- Window ( $\pm 5$  tokens of concept) — e.g. *He wore a [**yellow** t-shirt while eating a **banana** in the garage of a] kind stranger.*

---
- Syntactic — e.g. *If **bananas** aren't **yellow**, nobody will be interested in eating them, except for perhaps some questionable folk.*

# EXPERIMENTAL DETAILS 1

## Corpora

- **UMBC** 3.4B Tokens (English)
- **Wikipedia** Dump 2.5B/0.6B Tokens (English/Spanish)
- **ES1B** 1.4B Tokens (Spanish)

## Removing Co-occurrences

- Use 80% of corpus as training data and keep 20% as reserves.
- Sentences with co-occurrences removed from training data.
- Replaced with random sentence from reserves.

# EXPERIMENTAL DETAILS 2

## Vector Space

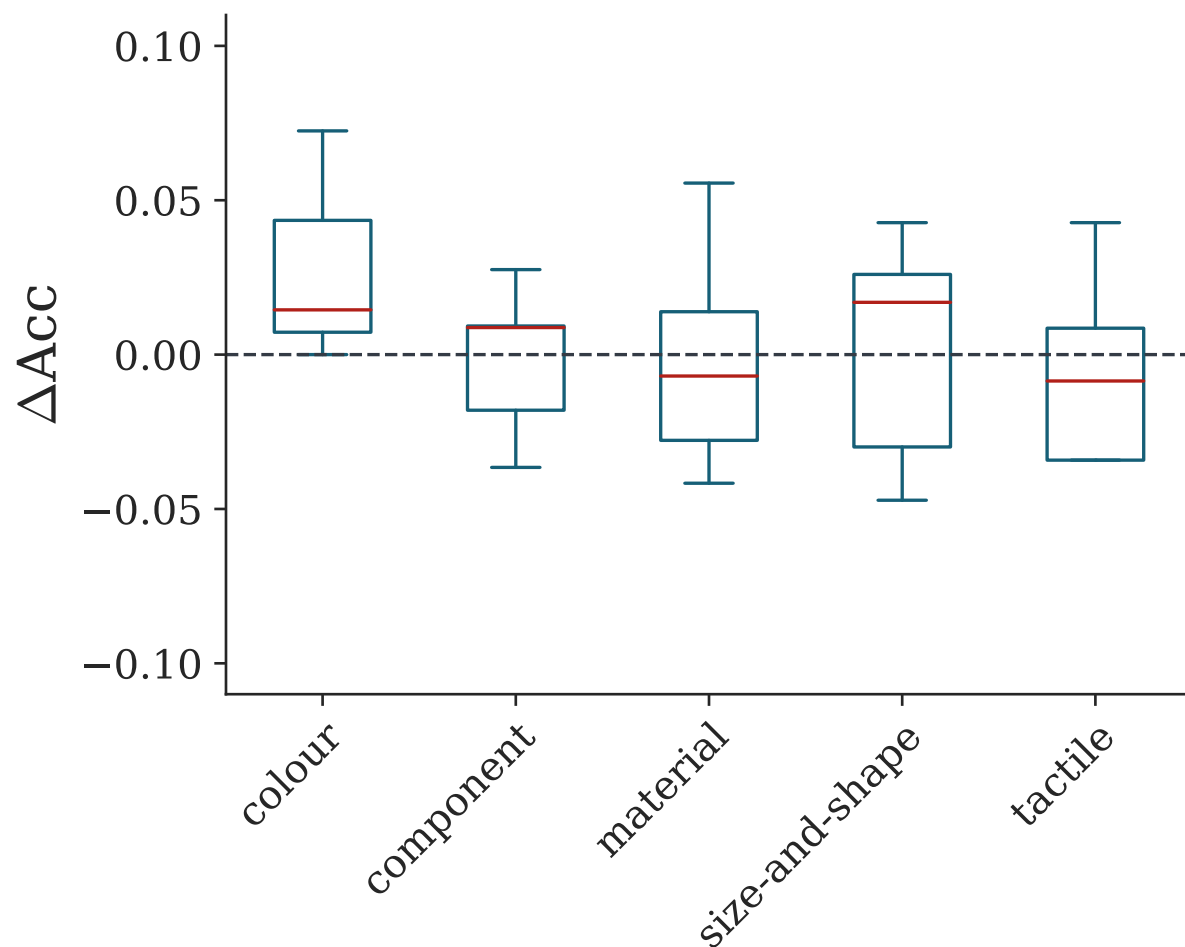
- **CBOW Word2Vec**
- Faster than skip-gram.
- Trained models with and without co-occurrences.
- For each separate trait type (e.g. colour)

## Classifiers

- Multi-class (predict yellow given banana as input, as embedding from model).
- Binary-class (predict related or not related with  $e_{\text{concept}} - e_{\text{relation}}$  as input).
- Used SVMs.

**RESULTS?**

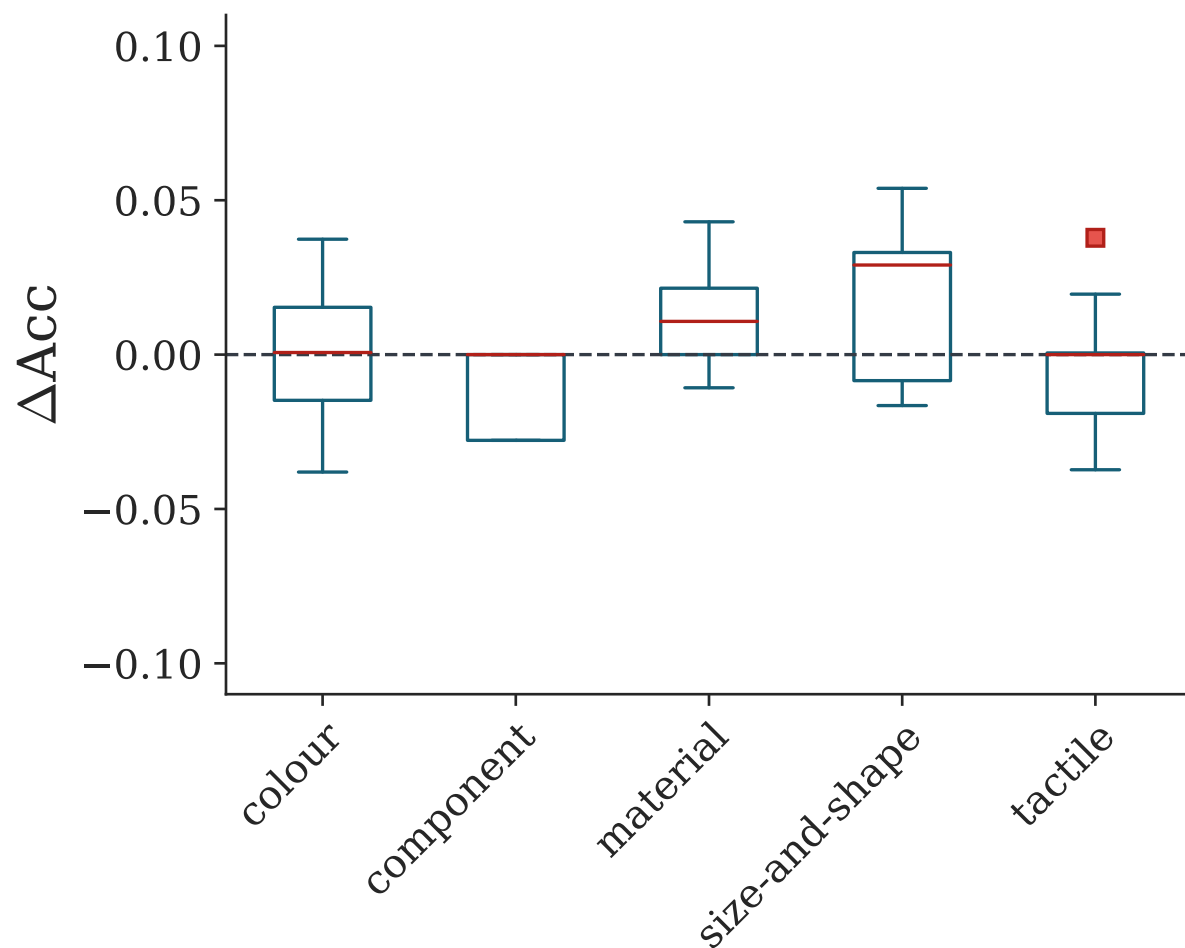
# McRAE (ENGLISH) — MULTI-CLASS



**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on trait type. For McRae dataset and English corpora.

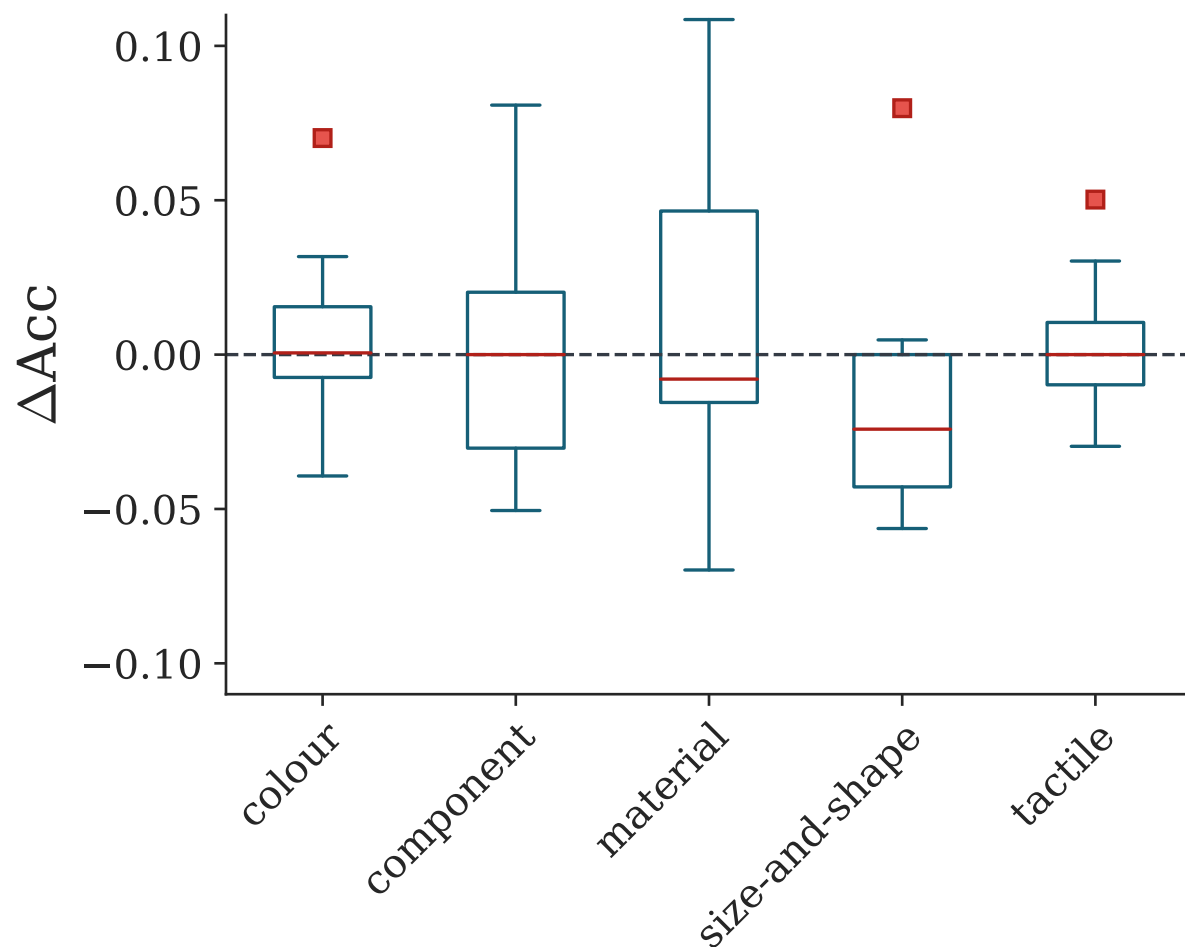


# NORMS (ENGLISH) — MULTI-CLASS



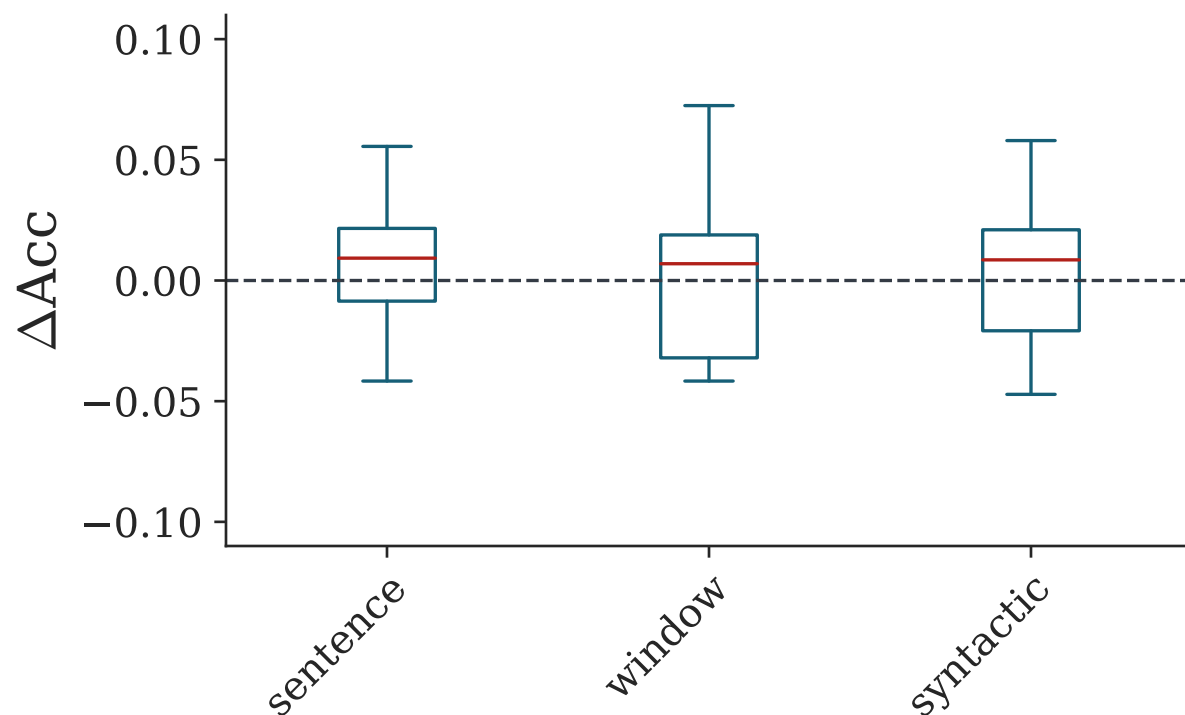
**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on trait type. For Norms dataset and English corpora.

# McRAE (SPANISH) — MULTI-CLASS



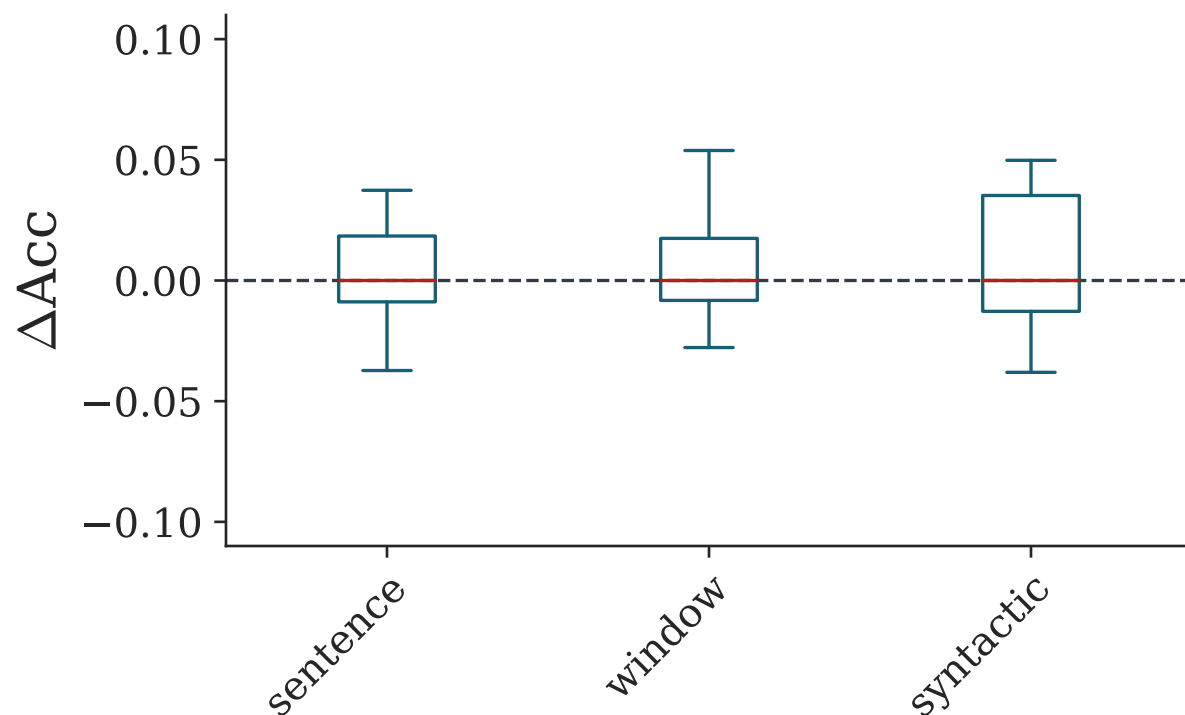
**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on trait type. For McRae dataset and Spanish corpora.

# McRAE (ENGLISH) — MULTI-CLASS



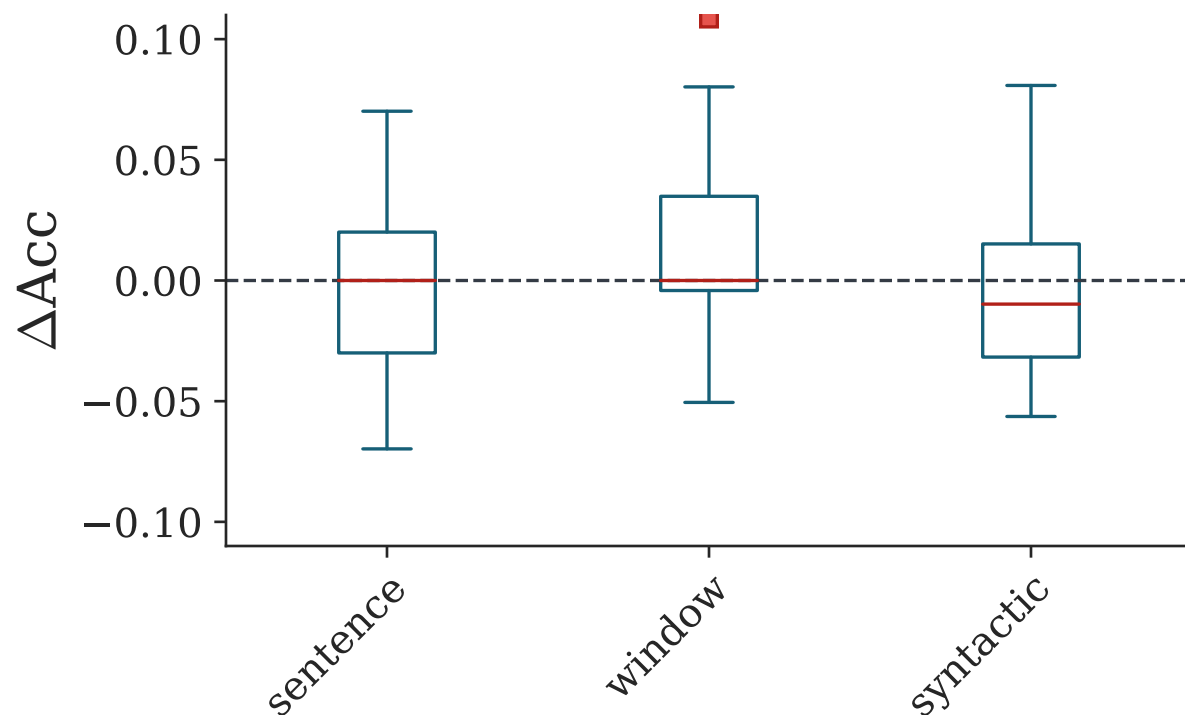
**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on co-occurrence removal method. For McRae dataset and English corpora.

# NORMS (ENGLISH) — MULTI-CLASS



**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on co-occurrence removal method. For Norms dataset and English corpora.

# McRAE (SPANISH) — MULTI-CLASS



**Figure 1:** Distribution of multi-class accuracy differences between models trained with and without co-occurrences separated based on co-occurrence removal method. For McRae dataset and Spanish corpora.

# BINARY RESULTS

---

- Typically quite high (80-90% accuracy)
- Require negative examples (clearly not trivial)
- Similar to multi-class (no real distinct patterns regarding removal of co-occurrences).

**END**

---

**THANKS. ANY QUESTIONS, REMARKS, OR CRITICISMS?**