





Be Different to Be Better! A Benchmark to Leverage the Complementarity of Language and Vision **Sandro Pezzelle**¹, Claudio Greco², Greta Gandolfi², Eleonora Gualdoni², Raffaella Bernardi^{2,3}

UNIVERSITÀ **DI TRENTO**

OVERVIEW

- Human communication, in real-life situations, is **multimodal** [1]
- Speakers do **not need to "repeat"** information provided by the environment
- In current computational approaches to language and vision the two modalities are *aligned* rather than *complementary*
- **bd2bb**: Novel task where either modality is necessary but not sufficient

DATASET

- Given an image depicting a real-life situation, 5 annotators provide 1) an **intention**, i.e., how they might feel/ behave if they were in that image (If I...); 2) a target action, i.e., what they would do based on that feeling/behavior (I will...)
- Intention: ungrounded; Action: grounded
- +10K valid <image, intention, action> datapoints + rich linguistic annotation





MODELS & SETTINGS



¹Institute for Logic, Language, and Computation, University of Amsterdam ²CIMeC, ³DISI, University of Trento sites.google.com/view/bd2bb

We introduce bd2bb, a novel benchmark that requires models combine complementary information from language and vision: +10K crowdsourced datapoints

While solving bd2bb is relatively easy for humans (~80% accuracy), SotA pre-trained multimodal encoders struggle to achieve similar results (~60% accuracy)

MULTIPLE-CHOICE TASK







L: attend a dinner like this man holding a gift L: buy him a cake and invite his friends to party **T:** act silly with this man and eat cake V: help my child cut their cake V: have cake with soldiers



I: If I am in the mood to act silly, I will.



sit next to the woman on the bench

get my face painted avert my eyes from the man who looks silly teach him how to tie a tie wear a costume and march in a parade

DISCUSSION

Multimodal integration is the key: L&V models outperform unimodal systems Best models are far from humans: gap of ~17% accuracy (20% in hard test set) **Pre-trained is better:** largely pre-trained systems outperform models trained from scratch on the bd2bb task

Importance of richly-annotated resources to test pre-trained systems

Importance of multimodal tasks where L&V are genuinely complementary

REFERENCES

[1] Gunther R Kress. 2010. *Multimodality: A social semiotic* approach to contemporary communication. Taylor & Francis. [2] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly optimized BERT pretraining approach. CoRR, abs/1907.11692. [3] Hao Tan and Mohit Bansal. 2019. LXMERT: Learning cross-modality encoder representations from transformers. In Proceedings of EMNLP-IJCNLP, pages 5103–5114.







