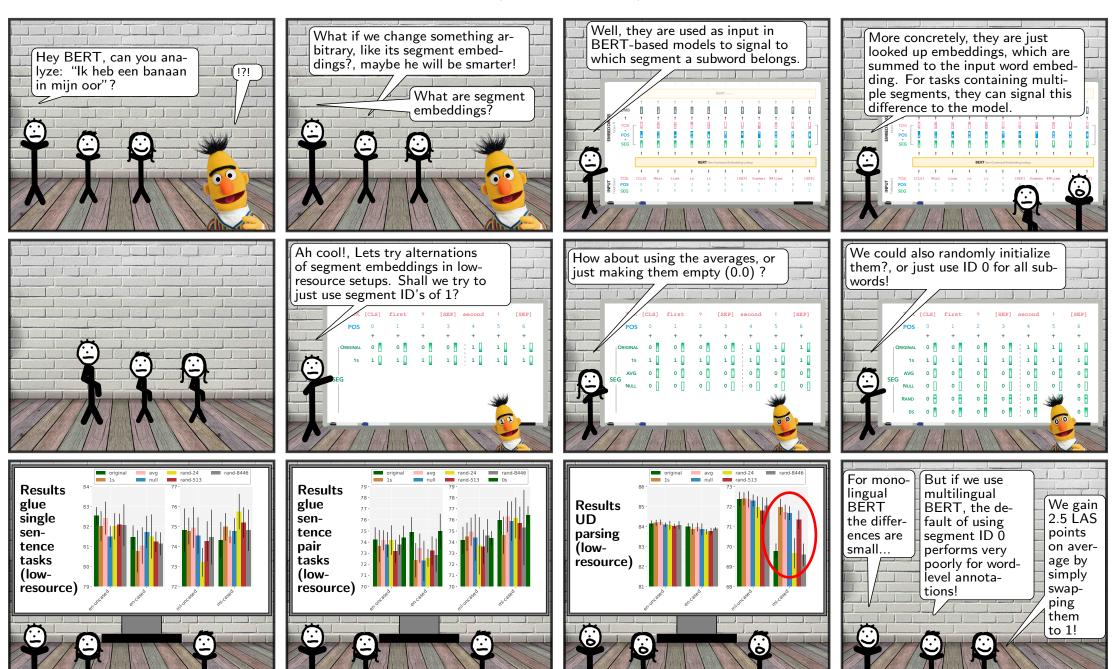# Frustratingly Easy Performance Improvements for Low-resource Setups:
## A Tale on BERT and Segment Embeddings

Rob van der Goot, Max Müller-Eberstein, Barbara Plank



*This is a work of fiction. Any similarities to persons or actual events is purely coincidental.