

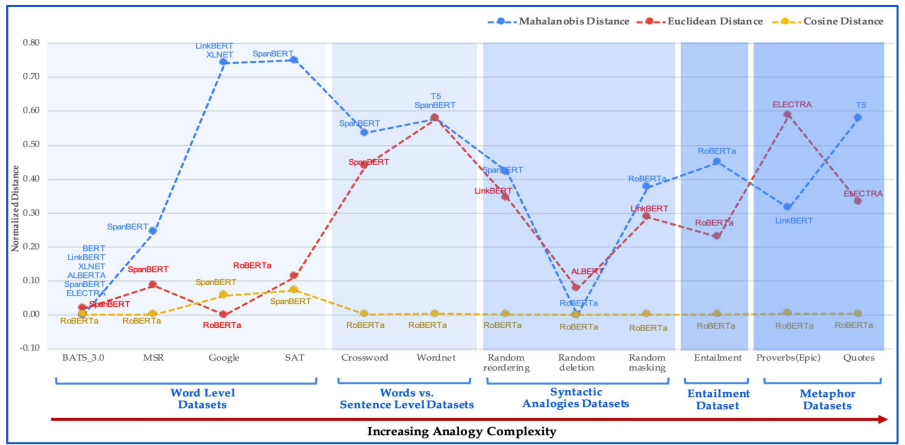
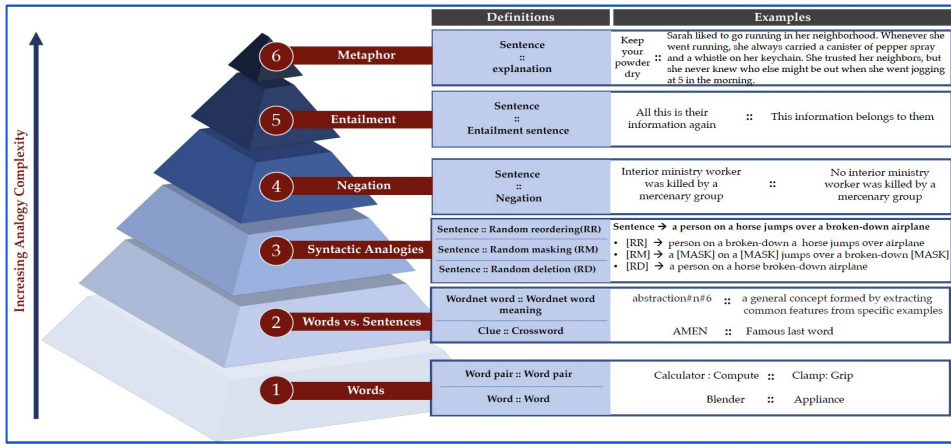


# ANALOGICAL

## A Novel Benchmark for Long Text Analogy Evaluation in Large Language Models

Thilini Wijesiriwardene<sup>1</sup>, Ruwan Wickramarachchi<sup>1</sup>, Bimal G. Gajera<sup>2</sup>, Shreyash Mukul Gowaikar<sup>3</sup>, Chandan Gupta<sup>4</sup>, Aman Chadha<sup>5,6,\*</sup>, Aishwarya Naresh Reganti<sup>7,\*</sup>, Amit Sheth<sup>1</sup>, Amitava Das<sup>1</sup>

<sup>1</sup>AI Institute, University of South Carolina, USA, <sup>2</sup>Nirma University, India, <sup>3</sup>BITS Pilani, Goa, India, <sup>4</sup>IIIT Delhi, India, <sup>5</sup>Amazon AI, USA, <sup>6</sup>Stanford, USA, <sup>7</sup>Amazon, USA



Word-level analogies intrinsically measure the quality of word embedding methods.

Large Language Models are primarily evaluated on extrinsic measures (e.g., GLUE, SuperGLUE) and there are only a few investigations on LLMs' ability to draw analogies between long texts.

ANALOGICAL, intrinsically evaluate LLMs across a six-level analogy taxonomy: (i) word, (ii) word vs. sentence, (iii) syntactic, (iv) negation, (v) entailment, and (vi) metaphor.

Abilities of 8 LLMs in identifying analogies in this taxonomy was evaluated (using 13 datasets and 3 different distance measures). We find that it is increasingly challenging for LLMs to identify analogies when going up the analogy taxonomy.