# Negative Text Similarity Analysis
## (SANW Motivation Evidence)
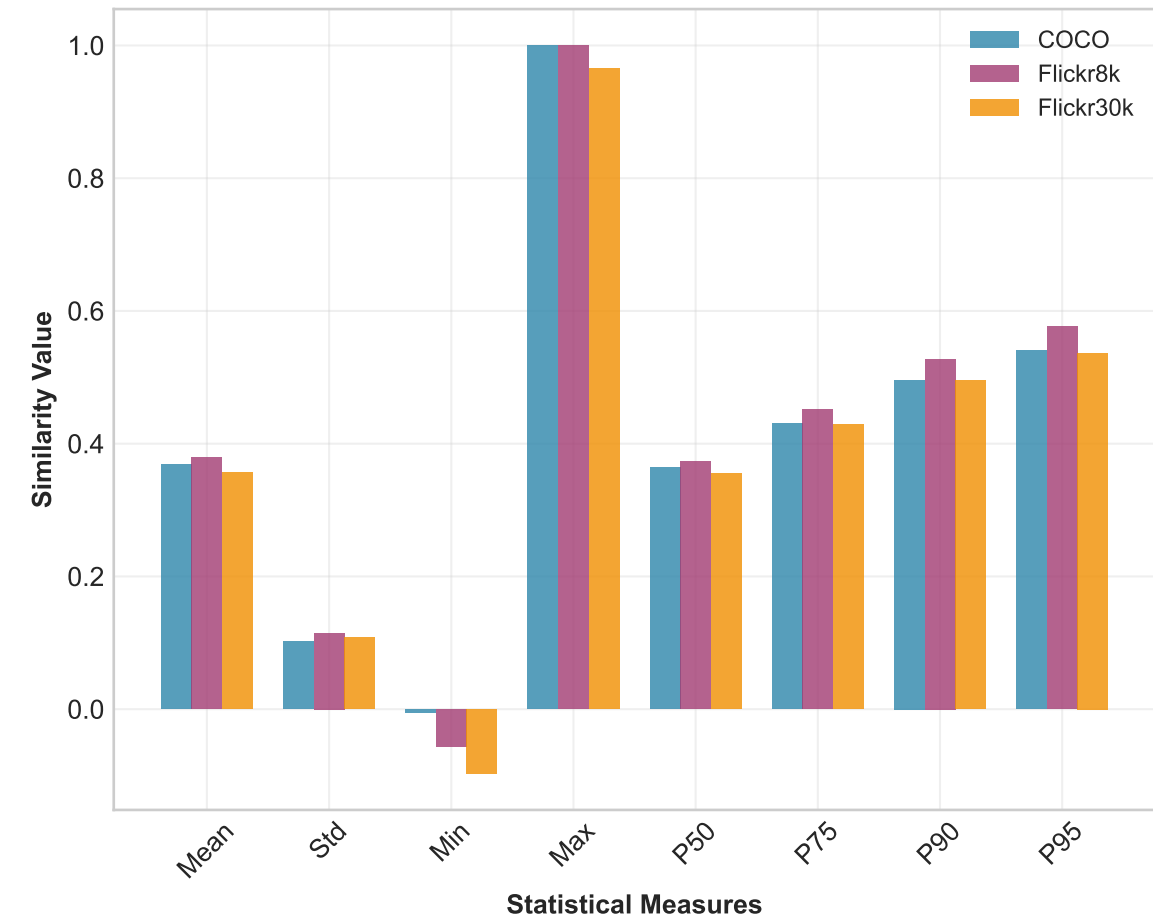


**Negative Text Similarity Analysis (SANW Motivation Evidence)** — grouped bar chart showing percentage of negative pairs by similarity threshold across datasets:
- COCO: 96.4% (≥0.2), 75.0% (≥0.3), 35.6% (≥0.4)
- FLICKR8K: 95.1% (≥0.2), 74.8% (≥0.3), 40.6% (≥0.4)
- FLICKR30K: 93.1% (≥0.2), 69.5% (≥0.3), 34.0% (≥0.4)

**Mean Similarity Across Datasets** — Mean Cosine Similarity:
- COCO: 0.369
- FLICKR8K: 0.379
- FLICKR30K: 0.357

**Similarity Threshold Analysis (%)** — heatmap:
- COCO: 96.4% (≥0.2), 75.0% (≥0.3), 35.6% (≥0.4)
- FLICKR8K: 95.1% (≥0.2), 74.8% (≥0.3), 40.6% (≥0.4)
- FLICKR30K: 93.1% (≥0.2), 69.5% (≥0.3), 34.0% (≥0.4)

**Statistical Distribution Analysis** — statistical measures (Mean, Std, Min, Max, P50, P75, P90, P95) for COCO, Flickr8k, Flickr30k.

**SANW Motivation: Current vs Proposed Approach** — Hard Negatives vs Semantic-aware Negatives:
- Current CLIP (Treats all negatives equally): 100% Hard Negatives, 35% Semantic-aware Negatives
- SANW Approach (Semantic-aware weighting): 65%

**Dataset Size Distribution** — pie chart:
- COCO: 50.4%
- Flickr30k: 39.5%
- Flickr8k: 10.1%

**KEY FINDINGS:**

- 93-96% of "negative" pairs have similarity ≥ 0.2 (very high!)

- 69-75% of "negative" pairs have similarity ≥ 0.3 (extremely high!)

- 34-41% of "negative" pairs have similarity ≥ 0.4 (significant!)

**CONCLUSION:**
Current CLIP training treats many semantically similar captions as "negatives" when they should be reweighted (SANW motivation).

**METHODOLOGY:**

1. Loaded captions from 3 datasets:
   - COCO (202K captions)
   - Flickr8k (40K captions)
   - Flickr30k (159K captions)

2. Used OpenCLIP ViT-B-32 (laion400m_e32) text encoder

3. Computed cosine similarities between all caption pairs in random batches (256 captions)

4. Analyzed off-diagonal similarities (excluding same-image pairs)

5. Calculated statistics across 50 batches per dataset

**Distribution of Negative Pair Similarities** — pie chart:
- 0.0-0.2: 5%
- 0.5+: 10%
- 0.4-0.5: 25%
- 0.3-0.4: 35%
- 0.2-0.3: 25%