

SCALABLE VISUAL TEXTUAL RETRIEVAL

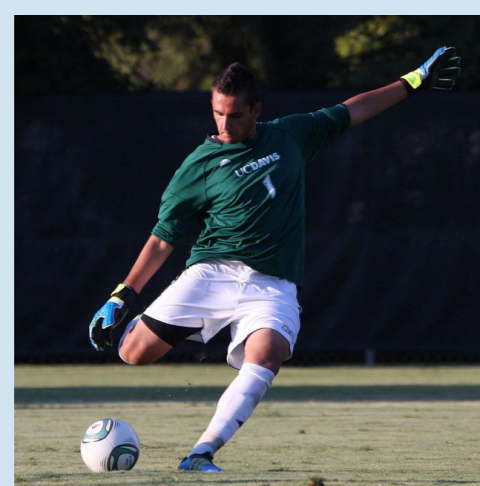
Query: "Player number 8 kicked the soccer ball with his foot."

$S(i_1, q) = 0.90$

$S(i_2, q) = 0.88$

$S(i_3, q) = 0.86$

$S(i_4, q) = 0.76$

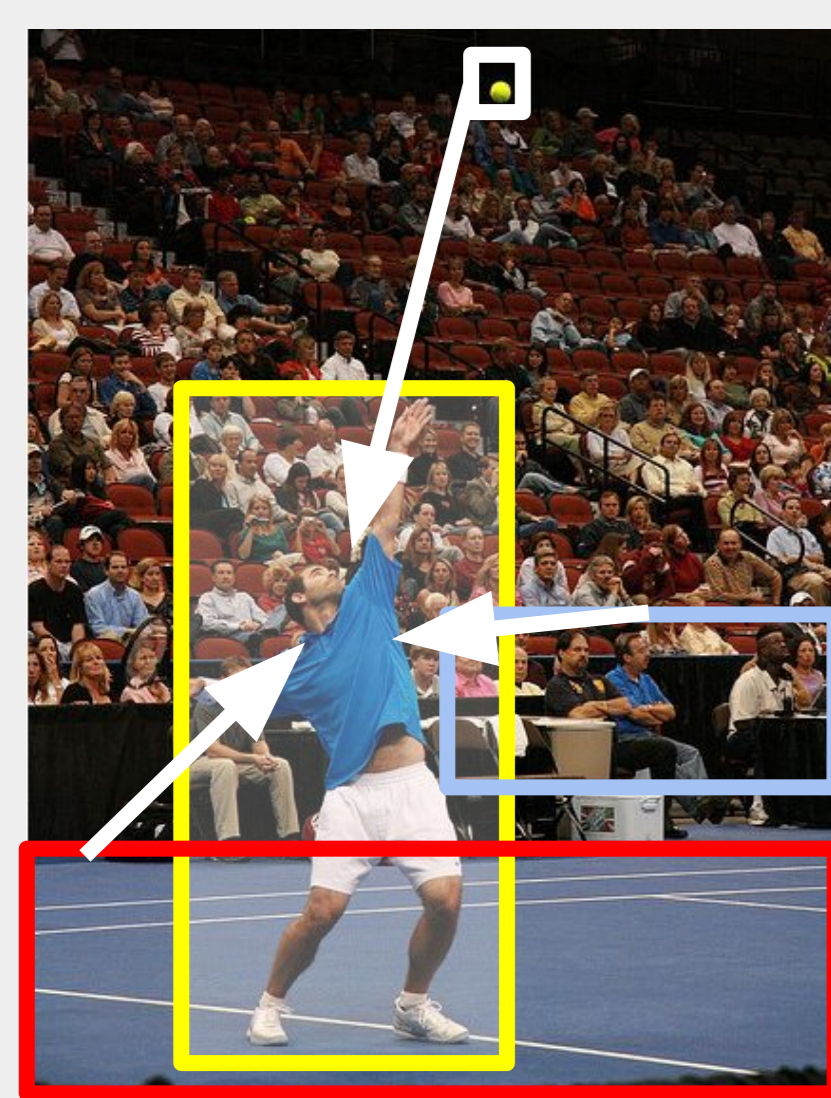


Decreasing Similarity

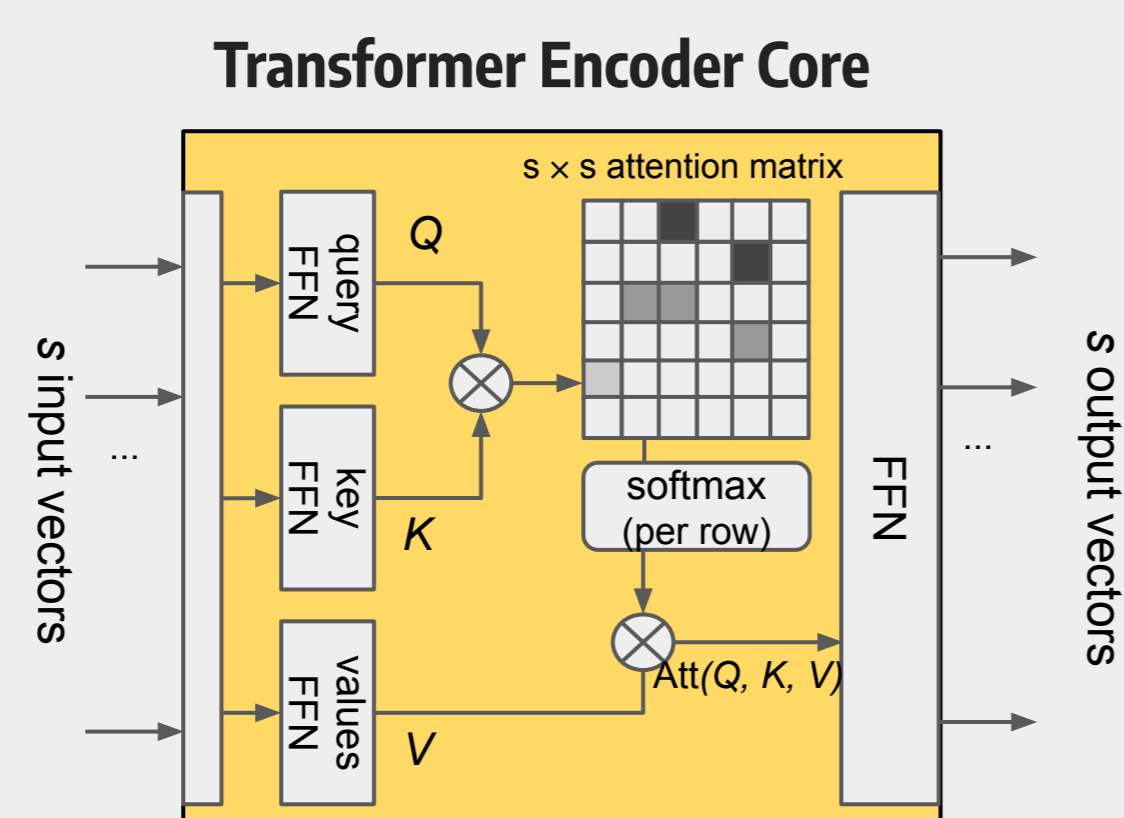
CHALLENGES

- Produce fixed-sized **visual** and **textual** features
- Easy to compare
- Can be indexed** using already existing text-based or metric space approaches
- They must carry **contextual information**
- It is **difficult** to represent **relationships between objects** within a single feature

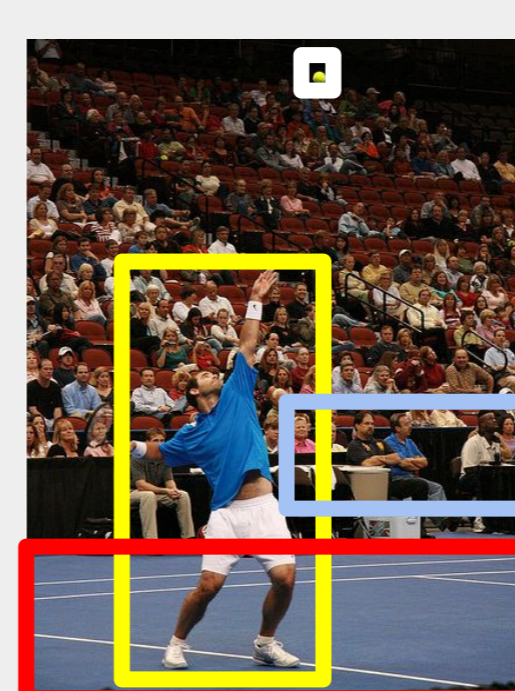
TRANSFORMER ENCODER REASONING NETWORK



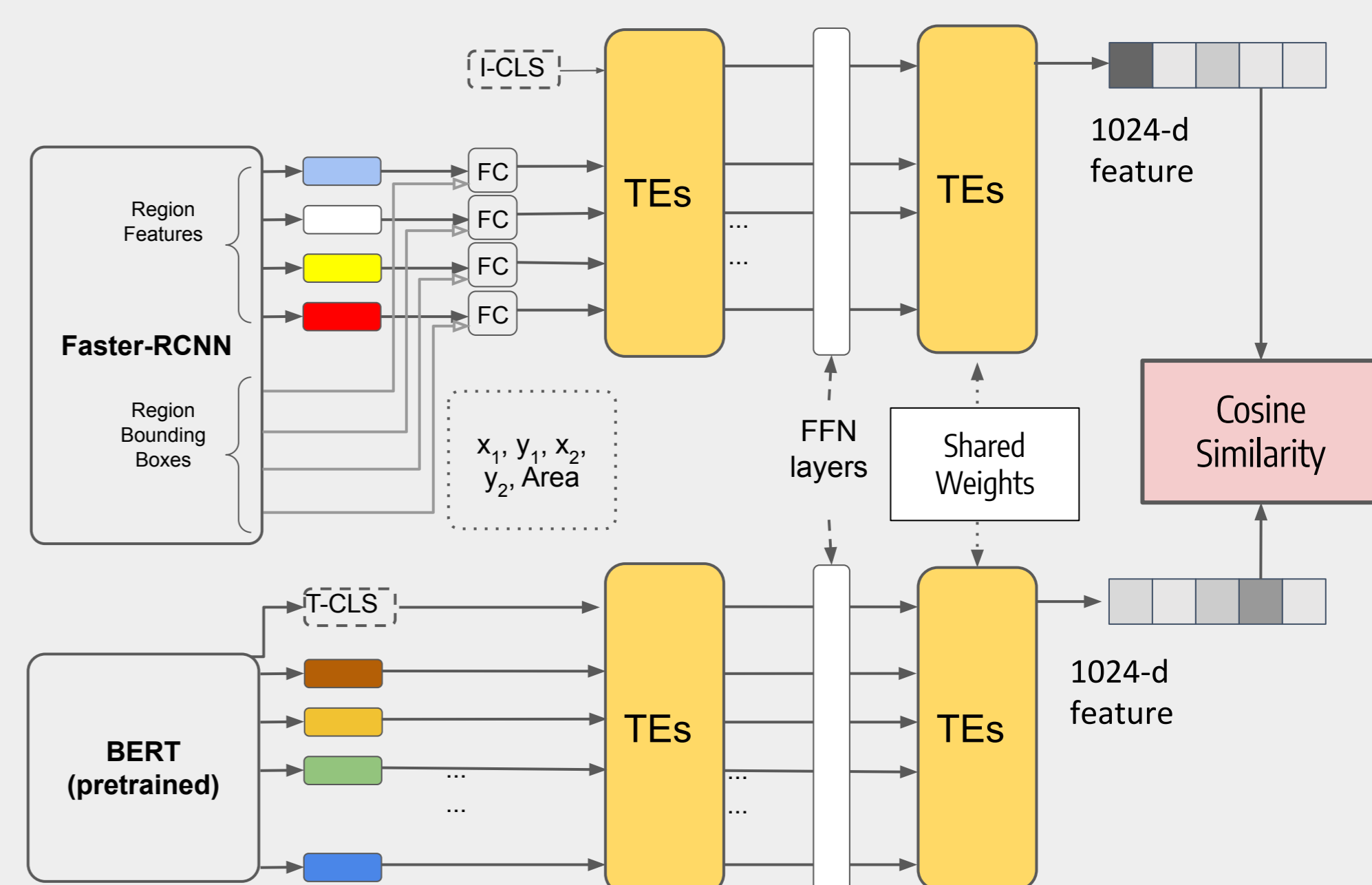
- Every region or word should look at its surroundings
- The information is accumulated through the **Transformer Encoder attention mechanism**



A tennis player serving a ball on the court



A tennis player serving a ball on the court



Training: Hinge-Based Triplet Ranking Loss

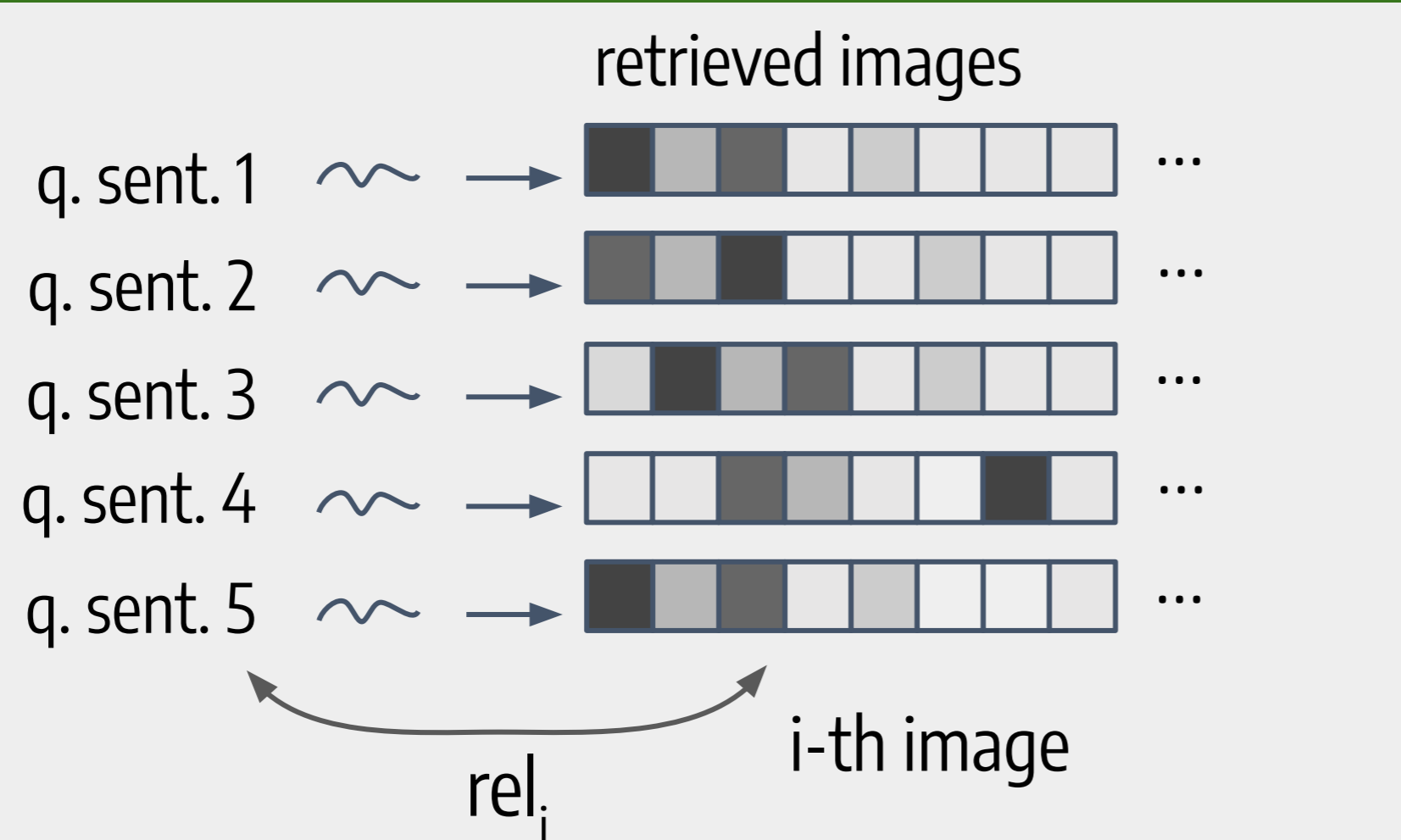
$$L_m(i, c) = \max_c [\alpha + S(i, c') - S(i, c)]_+ + \max_{i'} [\alpha + S(i', c) - S(i, c)]_+$$

Testing:

- Compute features
- Compute similarities
- Rank by decreasing sim.

EVALUATION: NDCG

- Non-exact matches
- High-level semantics



$$DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)}$$

where

$$rel_i = \text{ROUGE-L}(q, C_i)$$

$$rel_i = \text{SPICE}(q, C_i)$$

QUANTITATIVE RESULTS

- MS-COCO dataset (5 human-written sentences per image)
- 1K test set: 5-fold validation

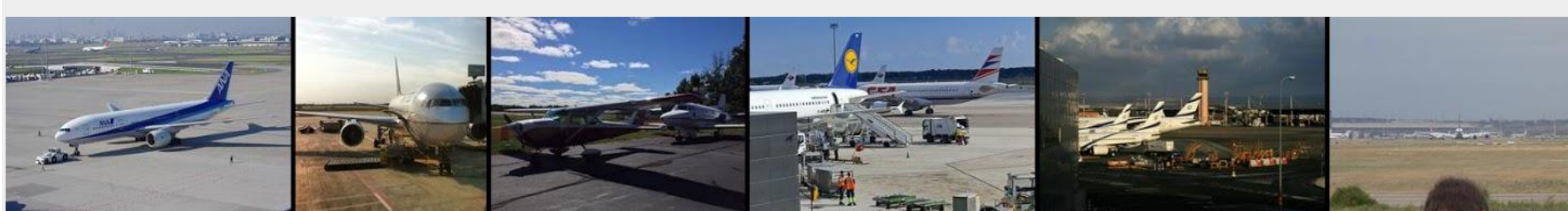
Model	ROUGE-L	SPICE
VSE-0	0.702	0.616
VSE++	0.712	0.617
VSRN	0.723	0.620
TERN (our)	0.725	0.653

1K test set

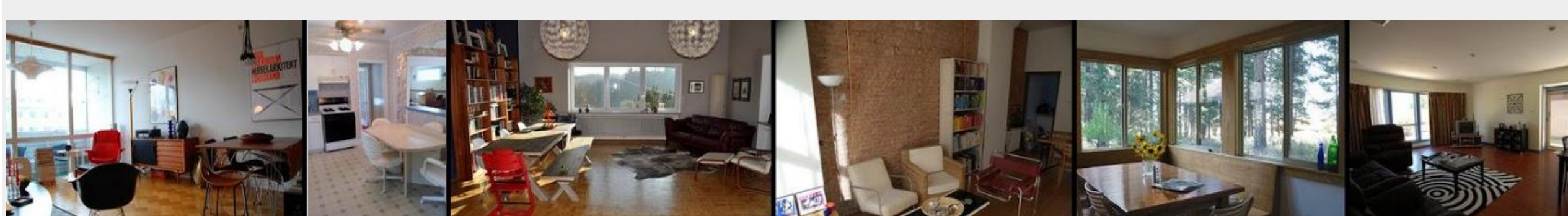
Model	ROUGE-L	SPICE
VSE-0	0.633	0.549
VSE++	0.656	0.577
VSRN	0.676	0.596
TERN (our)	0.665	0.600

5K test set

QUALITATIVE RESULTS

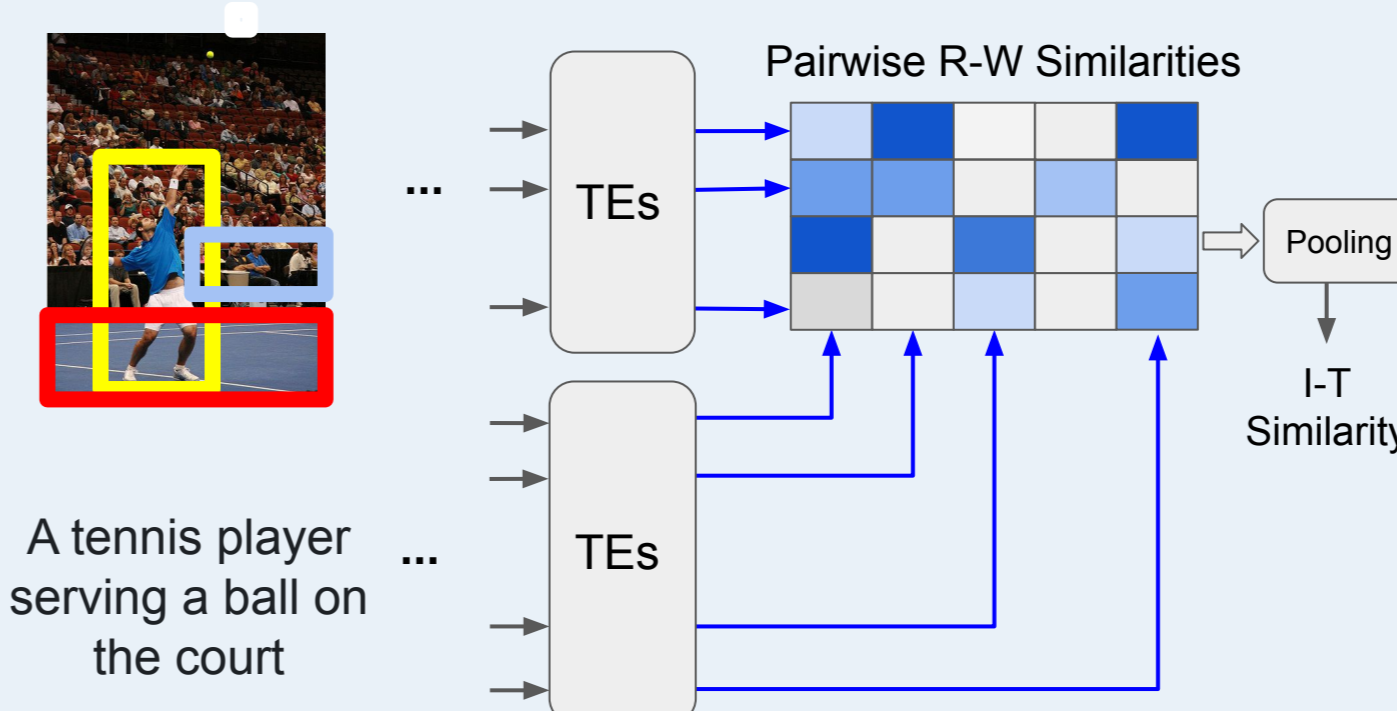


Query: A large jetliner sitting on top of an airport runway.



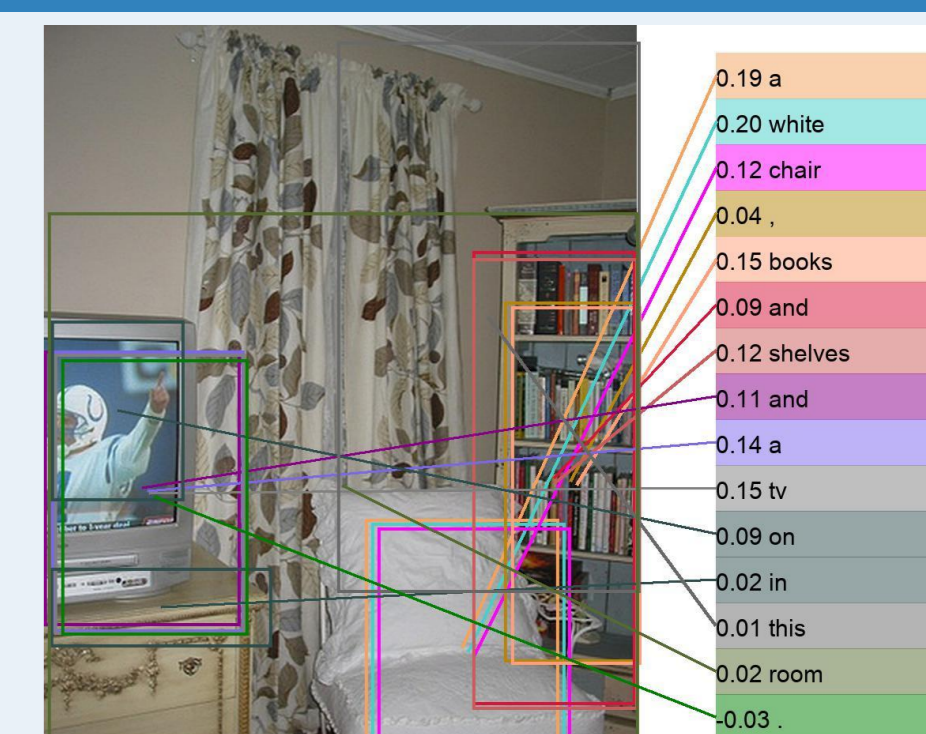
Query: An eating area with a table and a few chairs.

TOWARDS FINE-GRAINED ALIGNMENT: TERAN



Model	ROUGE-L	SPICE
TERN (our)	0.725	0.653
TERAN (our)	0.741	0.668

1K test set



Model	ROUGE-L	SPICE
TERN (our)	0.665	0.600
TERAN (our)	0.682	0.610

5K test set