

Multiple annotated text descriptions of each image

A man in a brown shirt and dark shorts plays on the beach with his two black dogs
A man in shorts with two black dogs holds a ball throwing toy at the beach
A man playing with two black dogs on the beach
A man with two dogs on a beach
Man at the beach with two dogs

TFIDF

CIDEr

Text
vector

T-T similarity

T_{11}	T_{12}	\dots	T_{1B}
T_{21}	T_{22}	\dots	T_{2B}
		\dots	
T_{B1}	T_{B1}	\dots	T_{BB}

Image

Preprocessing

Batch

Image
feature

SGRAF

Fully connected

Image
encoder

Image embeddings

v_1

v_2

v_3

\dots

v_B

I-I similarity

I_{11}	I_{12}	\dots	I_{1B}
I_{21}	I_{22}	\dots	I_{2B}
		\dots	
I_{B1}	I_{B1}	\dots	I_{BB}

Text

Batch

BiGRU

Text
encoder

Text embeddings

t_1

t_2

t_3

\dots

t_B

I-T similarity

C_{11}	C_{12}	\dots	C_{1B}
C_{21}	C_{22}	\dots	C_{2B}
		\dots	
C_{B1}	C_{B1}	\dots	C_{BB}

MSC loss

Triplet loss