Composed Query $q$ — Query Image $y$ — Query Text $t$: `dense`

Composed Encoder $h$

Visual Encoder $f$ — $d$ — Similarities — Distribution Normalization

Text Encoder $g$ — $d$ — Similarities — Distribution Normalization

$(1-\lambda)\cdot$ | 0.1 | 0.3 | ... | 0.2 | $\oplus$ | $\lambda\cdot$ | 0.2 | 0.1 | ... | 0.4 | $N$

| 0.2 | 0.1 | ... | 0.4 | $\rightarrow$ argmax $\rightarrow$

Retrieved Image $x$

Image Dataset $X$ — Visual Encoder $f$ — $N$, $d$