

FIRE-CIR: Fine-grained Reasoning for Composed Fashion Image Retrieval

Supplementary Material

1. Additional FIRE-CIR qualitative examples

To further illustrate the reasoning process, we present additional qualitative examples of FIRE-CIR on the Fashion IQ dataset (“dress” subset for Figure 1 and Figure 2, and “shirt” subset for Figure 3 and Figure 4).

For each example, the reference image and modification text are in the top-left corner. FIRE-CIR decomposes the modification text into a set of visual questions, listed on the left. These questions are then applied to the top-6 images retrieved by FashionBLIP-2 (from the first image on the top left to the sixth on the top right), and to the ground-truth target image (top-right corner). Each checkmark (resp. cross) corresponds to a predicted answer compatible (resp. incompatible) with the modification text. The probability of the answer being compatible with the text is displayed below the checkmark (resp. cross). Then, all answer probabilities are averaged for each candidate image to compute the VQA score, which is used to compute the final rank of each candidate image. The final rank is indicated at the bottom of each example. We also specify the rank difference with the FashionBLIP-2 results: a negative difference means that the candidate image gets a higher rank (so it is more relevant), while a positive difference means that the candidate image is further among the retrieved results (so it is less relevant).

Query	Image 1	Image 2	Image 3	Image 4	Image 5	Image 6	Image 7
 <p>“Has no sleeves with slim straps floral dress and is longer and has a circular pattern”</p>							
Is the dress sleeveless?	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00
Does the dress have slim straps?	✓ 0.56	✓ 0.53	✗ 0.47	✓ 0.85	✓ 0.95	✓ 1.00	✓ 0.85
Is the dress patterned with floral print?	✓ 0.65	✗ 0.15	✗ 0.20	✓ 0.50	✓ 0.96	✓ 0.50	✓ 0.50
Does the dress include a circular pattern?	✓ 0.62	✓ 1.00	✓ 0.73	✓ 1.00	✗ 0.38	✓ 0.59	✓ 1.00
Is the dress longer?	✓ 0.80	✓ 0.96	✗ 0.44	✓ 0.94	✓ 0.94	✓ 0.98	✓ 0.87
VQA score	0.73	0.73	0.57	0.86	0.84	0.81	0.84
New rank	3 (+2)	6 (+4)	23 (+20)	1 (-3)	2 (-3)	4 (-2)	5 (-4)

Figure 1. In this example, FIRE-CIR is able to accurately identify the pattern in the dresses and re-ranks the retrieved results according to their resemblance with what is described in the modification text.

Query	Image 1	Image 2	Image 3	Image 4	Image 5	Image 6	Image 7
 <p>“Plain black and sleeveless and is darker in color and solid colored”</p>							
Is the dress black?	✗ 0.00	✓ 1.00	✓ 0.50	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00
Is the dress sleeveless?	✓ 1.00	✓ 1.00	✓ 0.99	✓ 1.00	✓ 1.00	✓ 1.00	✓ 0.99
Is the dress plain?	✓ 0.56	✓ 1.00	✓ 0.98	✓ 0.96	✓ 0.99	✓ 0.50	✓ 0.91
Is the dress darker?	✓ 0.62	✓ 1.00	✓ 0.97	✓ 0.98	✓ 0.97	✓ 0.98	✓ 0.96
Is the dress a solid color?	✗ 0.29	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00	✓ 1.00
VQA score	0.50	1.00	0.89	0.99	0.99	0.90	0.97
New rank	73 (+72)	1 (-1)	3 (+0)	2 (-1)	4 (-1)	5 (-2)	8 (-4)

Figure 2. While FashionBLIP-2 ranks highly the reference image as some of its visual aspects remain compatible with the CIR query, FIRE-CIR correctly identifies that this dress does not have a solid black color, and thus re-ranks it outside of the top-50 results. The ground-truth target image is at the 8th rank, but the 7 other top-ranked images are also compatible with the given query.



Figure 3. Similarly, visual similarity and the green design contribute to having candidate images similar to the reference one in the top-retrieved results. However, FIRE-CIR accurately promotes the ground-truth image, as its color is more coherent with what is written in the modification text.

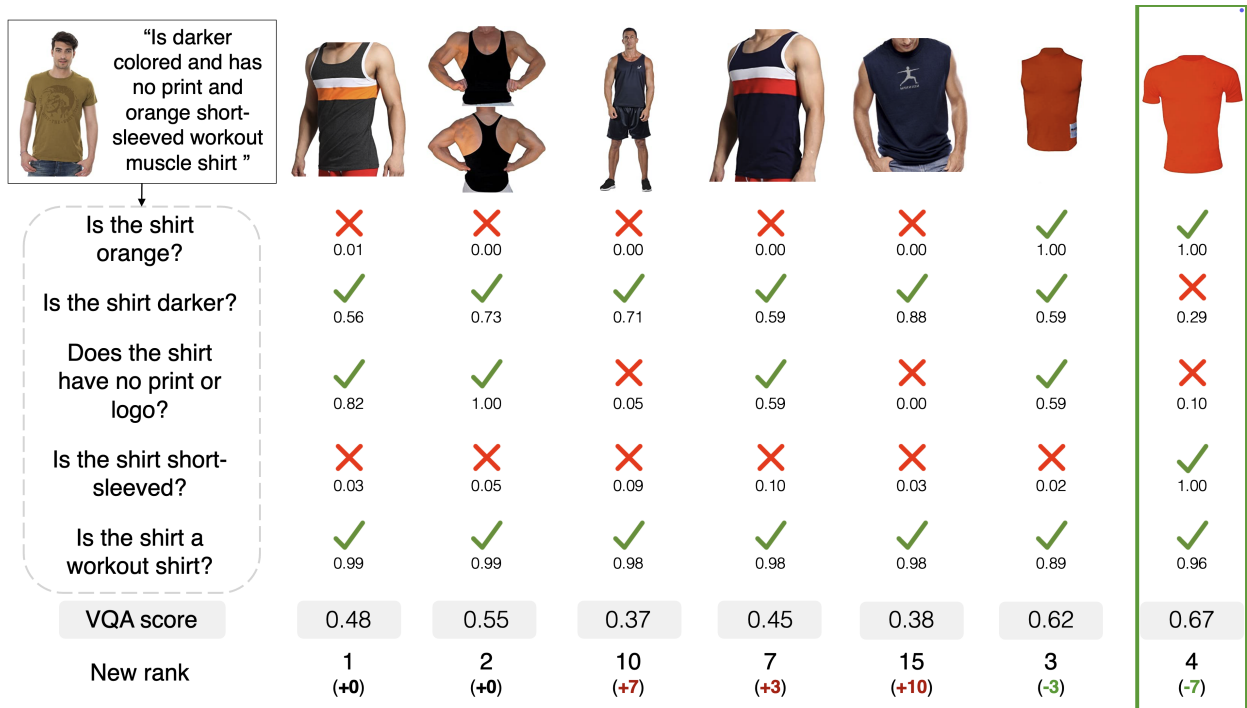


Figure 4. Contrary to FashionBLIP-2 which focuses specifically on the “workout muscle” characteristic, FIRE-CIR gives more importance to all the features, including the orange color.