**Supplementary Material**



the remains of the small creature called tiny were found in the scottish borders . ###-### million years ago, scotland lay close to the equator and its land was hot . researcher say tiny was one of the first four-legged creatures to move onto land . the findings fills in a ##-million year fossil gap when fish transitioned to land life .

**Pictorial Reference**

## Model Output

the remains of the creature, dubbed 'tiny,' were found in the scottish borders in south eastern scotland in a piece of rock smaller than a clenched fist, . 'tiny' was one of the first four-legged creatures to move onto land - making it our ancestor and filling in a -million year fossil gap when fish transitioned to becoming land [UNK] .
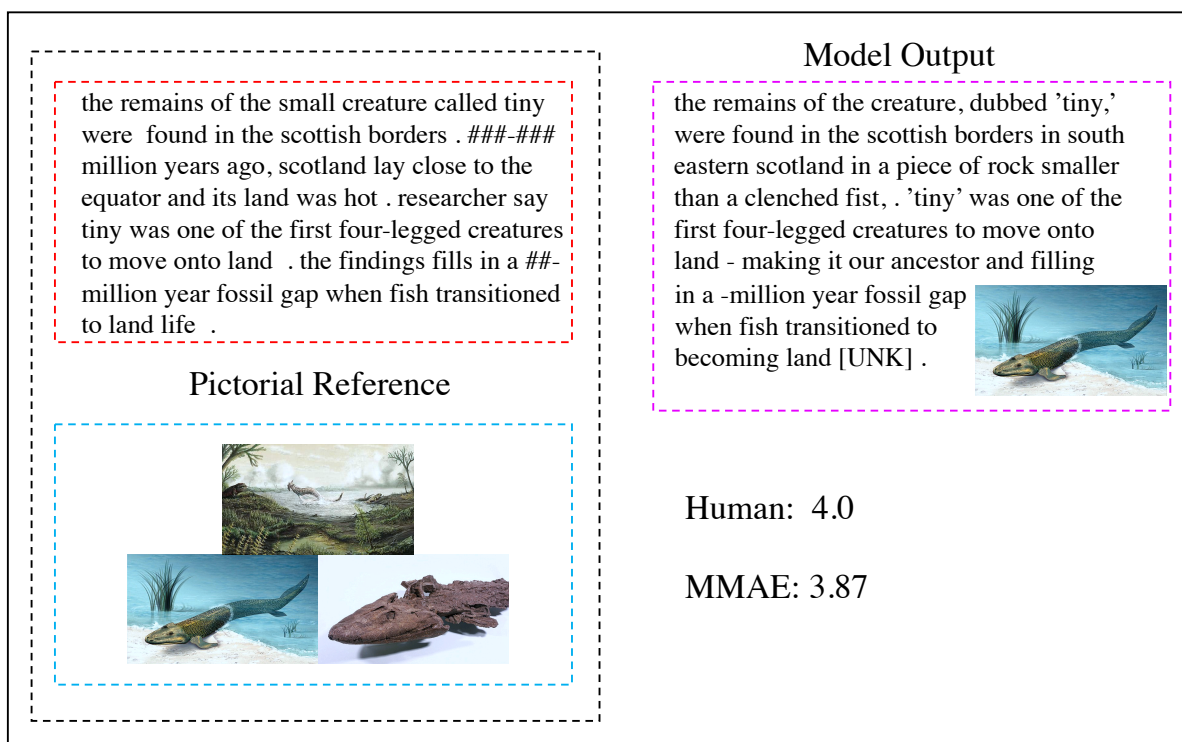
Human: 4.0

MMAE: 3.87

Figure 1: An example for Our MMAE. The left half is the pictorial reference and the right half gives the model output (from ATG), the human judgment score, and MMAE score. "#" denotes the digit number in the articles. In this example, the pictorial reference has three images. The model selects an image related to the topic and outputs a good quality text summary. Meanwhile, the image is also related to the text in the summary. Thus the output get a high human judgment score of 4, and the score by MMAE is 3.87, which is very close to the human judgment score.
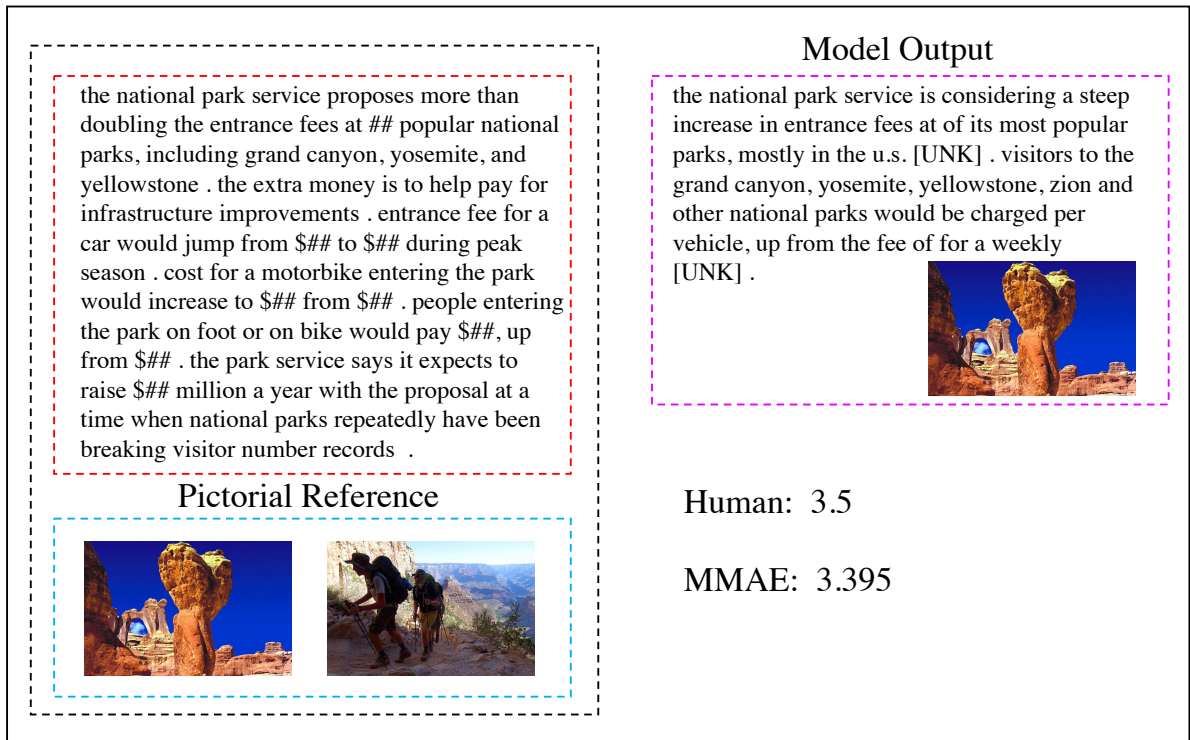
Figure 2: An example for Our MMAE. The left half is the pictorial reference and the right half gives the model pictorial output (from ATG), the human judgment score, and MMAE score. "#" denotes the digit number in the articles. The model selects an image that is very relevant to the topic and most of the contents of the summary are suitable. But the "[UNK]" in the output affects the readability and the quality of the summary. Thus this output gets a moderate human judgment score (3.5), and the score by MMAE is 3.395, which correlates well with the human judgment score.

hundreds of photographers show the lives of military personnel life in british army photographic competition . shots show exercises, marches and combat training as gunners and sergeants document daily service life . prestigious annual contest is open to all regular and reserve personnel, staff cadets and military contractors . photography is a recognised trade in the royal logistic corps and ## professional snappers are enlisted .

Pictorial Reference

Model Output

intimate shots of british army life have showcased the finest photographers working with our military as part of an annual competition . bombardier murray kerr's stunning still of servicemen battling with batons and shields . the -year-old from glasgow served as a gunner in afghanistan has been a reservist photographer for years .
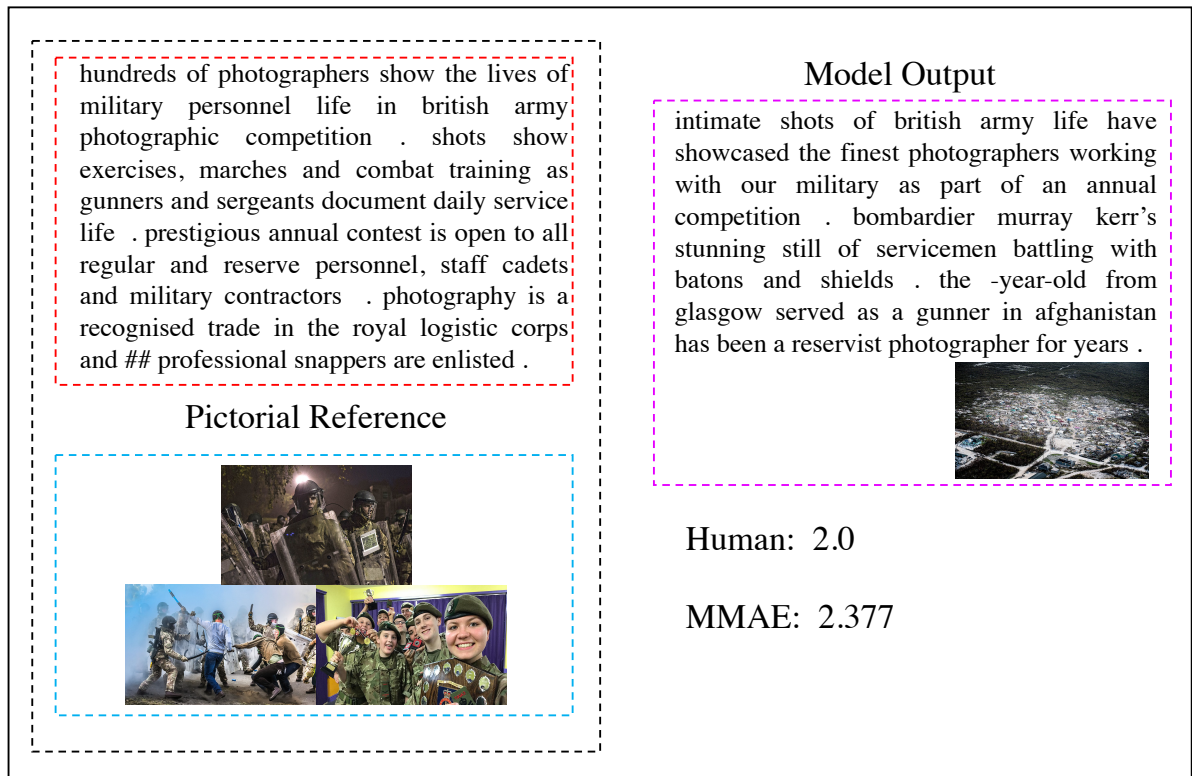
Human: 2.0

MMAE: 2.377

Figure 3: An example for Our MMAE. The left half is the pictorial reference and the right half gives the model pictorial output (from ATG), the human judgment score, and MMAE score. "#" denotes the digit number in the articles. There is no suitable image in the output of the model, and there is no correspondence between the image and the text. Thus our MMAE gives a very low score (2.377) for the output. Meanwhile, the human judgment score of this output is 2.0.