# Supplementary Material

# Speaking the Same Language: Matching Machine to Human Captions by Adversarial Training

Rakshith Shetty[1]    Marcus Rohrbach[2,3]    Lisa Anne Hendricks[2]

Mario Fritz[1]    Bernt Schiele[1]

[1]Max Planck Institute for Informatics, Saarland Informatics Campus, Saarbrücken, Germany
[2]UC Berkeley EECS, CA, United States    [3]Facebook AI Research

We present several qualitative examples to illustrate the strengths of our adversarially trained caption generator. All the examples are from the sampled versions of the adversarial (*adv-samp*) and the baseline (*base-samp*) models presented in the main paper. We show qualitative examples to highlight two main merits of the adversarial caption generator. First, we illustrate diversity across images in Section 1. Next, we demonstrate diversity when sampling multiple captions for each image in Section 2.

## 1. Illustrating diversity across images

Our adversarial model produces diverse captions across different images containing similar content, whereas the baseline model tends to use common generic captions which is mostly correct but not descriptive. We quantify this by visualizing the most frequently generated captions by the baseline model on the test set in Table 1a. Note that only the most likely caption according to the model is considered. Table 1a shows that the baseline model repeatedly generates identical captions for different images. However, the adversarial model is less prone to repeating generic captions, as seen in Table 1b. This is visualized in Figures 1, 2, and 3. Here we show sets of five images for which the baseline model generates identical generic caption. The five images are picked from among the images corresponding to captions in Table 1a, starting from the most frequent caption. Some entries are skipped, for example the caption in the last row, to avoid repeated concepts. While the baseline model produces a fixed generic caption for these images, we see that the adversarial model produces diverse captions which are often more specific to the contents of the image. For example, in the last row of Figure 2 we can see that the baseline model produces the generic caption "*a man rid-*

*ing skis down a snow covered slope*", whereas captions produced by the adversarial model include more image specific terms like "*jumping*", "*turn*", "*steep*", and "*corss country skiing*".

## 2. Illustrating diversity in captions for a single image

To qualitatively demonstrate the diverse captions produced by our adversarial model for each image, we visualize three captions produced by the adversarial and the baseline model for each input image. This is shown in figures 4 and 5. The captions are obtained by retaining the top three caption samples out of five (ranked by models' probability) from each model. Here bi-grams which are top-100 frequent bi-grams in the training set are highlighted in red (e.g., "a group" and "group of"). Additionally captions which are replicas from training set are marked with a '•' symbol. We can see that adversarial generator produces more diverse sets of captions for each image without over-using more frequent bi-grams and producing more novel sentences. For example, in the two figures, we see that the baseline model produces 22 captions (out of 45) which are copies from the training set, whereas the adversarial model does so only six times.

| Sentence | # baseline | # adversarial | Sentence | # adversarial | # baseline |
|---|---|---|---|---|---|
| a man riding a wave on top of a surfboard | 54 | 20 | a man riding a wave on top of a surfboard | 20 | 54 |
| a bathroom with a toilet and a sink | 44 | 1 | a skateboarder is attempting to do a trick | 16 | 0 |
| a baseball player swinging a bat at a ball | 37 | 2 | a female tennis player in action on the court | 16 | 0 |
| a man riding skis down a snow covered slope | 29 | 5 | a living room filled with furniture and a flat screen tv | 15 | 0 |
| a man holding a tennis racquet on a tennis court | 26 | 3 | a bus that is sitting in the street | 15 | 3 |
| a bathroom with a sink and a mirror | 25 | 3 | a long long train on a steel track | 11 | 0 |
| a man riding a snowboard down a snow covered slope | 24 | 4 | a close up of a sandwich on a plate | 10 | 0 |
| a baseball player holding a bat on a field | 22 | 1 | a baseball player swinging at a pitched ball | 10 | 0 |
| a man riding a skateboard down a street | 21 | 2 | a bus that is driving down the street | 9 | 1 |
| a bus is parked on the side of the road | 20 | 0 | a boat that is floating in the water | 8 | 3 |
| . . . | . . . | . . . | . . . | . . . | . . . |

(a)  (b)

Table 1: Most frequently repeated captions generated by (a) the baseline model and (b) the adversarial model on the test set of 5000 images. Columns show the number of times the caption was produced by the two models. The caption "a man riding a wave on top of a surfboard" is also the most frequently generated caption by the adversarial model, albeit less than half the times of the baseline model.



| | | | | |
|---|---|---|---|---|
| Adv-samp: a group of men sitting around a meeting room | a group of people sitting at a bar drinking wine | a group of friends enjoying lunch outdoors at a outdoor event | a group of people sitting at tables outside | a couple of men that are working on laptops |

Base-samp: . . . . . . . . . . . . . . . a group of people sitting around a table . . . . . . . . . . . . . . .

| | | | | |
|---|---|---|---|---|
| Adv-samp: a laptop and a desktop computer sit on a desk | a person is working on a computer screen | a cup of coffee sitting next to a laptop | a laptop computer sitting on top of a desk next to a desktop computer | a picture of a computer on a desk |

Base-samp: . . . . . . . . . . . . . . . a laptop computer sitting on top of a desk . . . . . . . . . . . . . . .

Figure 1: Illustrating diversity across images

Adv-samp: a surfer rides a large wave in the ocean

a surfer is falling off his board as he rides a wave

a person on a surfboard riding a wave

a man surfing on a surfboard in rough waters

a surfer rides a small wave in the ocean

Base-samp: · · · · · · · · · · · · · · a man riding a wave on top of a surfboard · · · · · · · · · · · · · ·

Adv-samp: a bathroom with a walk in shower and a sink

a dirty bathroom with a broken toilet and sink

a view of a very nice looking rest room

a white toilet in a public restroom stall

a small bathroom has a broken toilet and a broken sink

Base-samp: · · · · · · · · · · · · · · a bathroom with a toilet and a sink · · · · · · · · · · · · · ·

Adv-samp: a baseball player getting ready to swing at a pitch

a boy in a baseball uniform swinging a bat

a group of kids playing baseball on a field

a baseball game in progress with the batter upt to swing

a crowd watches a baseball game being played

Base-samp: · · · · · · · · · · · · · · a baseball player swinging a bat at a ball · · · · · · · · · · · · · ·
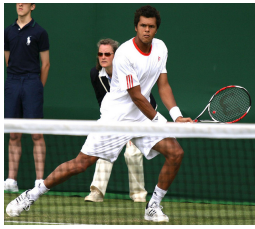
Adv-samp: a person on skis jumping over a ramp

a skier is making a turn on a course

a cross country skier makes his way through the snow

a skier is headed down a steep slope

a person cross country skiing on a trail

Base-samp: · · · · · · · · · · · · · · a man riding skis down a snow covered slope · · · · · · · · · · · · · ·

Figure 2: Illustrating diversity across images

| Adv-samp | a tennis player gets ready to return a serve | two men dressed in costumes and holding tennis rackets | a tennis player hits the ball during a match | a male tennis player in action on the court | a man in white is about to serve a tennis ball |

Base-samp · · · · · · · · · · · · · a man holding a tennis racquet on a tennis court · · · · · · · · · · · · · ·

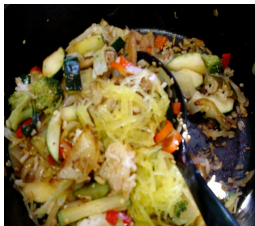| Adv-samp | a young boy riding a skateboard down a street | a skateboarder is attempting to do a trick | a boy wearing a helmet and knee pads riding a skateboard | a boy in white shirt doing a trick on a skateboard | a boy is skateboarding down a street in a neighborhood |

Base-samp · · · · · · · · · · · · · a man riding a skateboard down a steet · · · · · · · · · · · · · ·

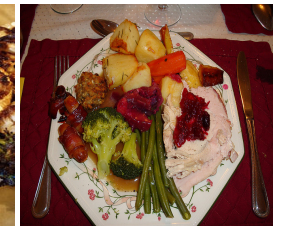| Adv-samp | a dish with noodles and vegetables in it | a plate of food that has some fried eggs on it | a meal consisting of rice meat and vegetables | a close up of some meat and potatoes | a plate with some meat and vegetables on it |

Base-samp · · · · · · · · · · · · · a plate of food with meat and vegetables · · · · · · · · · · · · · ·

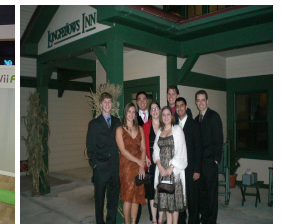| Adv-samp | a group of people standing around a shop | a group of young people standing around talking on cell phones | a group of soldiers stand in front of microphones | a couple of women standing next to a man in front of a store | a group of people posing for a photo in formal wear |

Base-samp · · · · · · · · · · · · · a group of people standing around a table · · · · · · · · · · · · · ·

Figure 3: Illustrating diversity across images

**Adv-samp**
- a red motorcycle parked on the side of the road
- a motorcycle is parked on a city street

a red motorcycle parked on a street in a city

**Base-samp**
- a man riding a motorcycle down a street
- a person riding a motorcycle on a city street
- a motorcycle is parked on the side of the road

**Adv-samp**
a motor cycle parked outside a building with people nearby
- a motor bike parked in front of a building

a row of bicycles parked outside of a building

**Base-samp**
- a group of people riding bikes down a street
a group of people on a street with motorcycles
a group of people on a street with motorcycles

**Adv-samp**
a motorcycle parked in front of a group of people

a police officer on a motorcycle in front of a crowd of people

a police officer on his motorcycle in front of a crowd

**Base-samp**
- a motorcycle parked on the side of a road
a motorcycle is parked on the side of a road
a motorcycle parked on the side of a road with a person walking by

**Adv-samp**
a skier is jumping over a snow covered hill

a person on skis jumping over a hill

a skier is in mid air after completing a jump

**Base-samp**
- a man riding skis down a snow covered slope
- a man riding skis down a snow covered slope
a person on a snowboard jumping over a ramp

**Adv-samp**
a group of people watching a skateboarder do stunts

a group of skateboarders performing tricks at a skate park

a group of skateboarders watch as others watch

**Base-samp**
- a man doing a trick on a skateboard at a skate park
- a man riding a skateboard down a rail
- a man doing a trick on a skateboard in a park

**Adv-samp**
- a stop sign with graffiti written on it

a stop sign with a few stickers on it

a stop sign has graffiti written all over it

**Base-samp**
a stop sign with a street sign on top of it

a stop sign with a street sign on top of it
- a stop sign with a street sign above it

**Adv-samp**
a man standing next to a mail truck
a picture of people standing outside a business

a couple of people standing by a bus

**Base-samp**
a group of people standing around a white truck
a group of people standing around a white truck
a group of people standing on a street next to a white truck

**Adv-samp**
a bouquet of flowers in a vase on a table
a bouquet of flowers in a vase on a table

a vase full of purple flowers sitting on a table

**Base-samp**
- a vase filled with flowers sitting on top of a table
- a vase filled with flowers sitting on top of a table
- a vase of flowers sitting on a table

**Adv-samp**
a plant in a vase on a wooden porch
a small blue flower vase sitting on a wooden porch

a large flower arrangement in a vase on a corner

**Base-samp**
- a vase with flowers in it sitting on a table
- a vase of flowers sitting on a table
- a vase with flowers in it on a table

Figure 4: Comparing 3 captions sampled from adversarial model to the baseline model. Bi-grams which are top-100 frequent bi-grams in the training set are highlighted in red. Captions which are replicas from training set are marked with a ●.

|  | | | |
|---|---|---|---|
| Adv-samp | a long line of stairs leading to a church<br>• a large cathedral filled with lots of pews<br>a cathedral with stained glass windows and a few people | a large church with a very tall tower<br>a large tall brick building with a clock on it<br>a church steeple with a clock on its side | several stop signs in front of some buildings<br>a stop sign in front of some graffiti writing<br>several stop signs lined up in a row |
| Base-samp | a large building with a large window and a building<br>a row of benches in front of a building<br>a church with a large window and a large building | a large clock tower in a city<br><br>a large clock tower in a city with a sky background<br>• a large clock tower in the middle of a park | • a stop sign with graffiti on it<br><br>a stop sign is shown with a lot of graffiti on it<br>a stop sign and a stop sign in front of a building |



|  | | | |
|---|---|---|---|
| Adv-samp | a family enjoying pizza at a restaurant party<br><br>a group of friends enjoying pizza and drinking beer<br>a group of kids enjoying pizza and drinking beer | a view of a city street at dusk<br><br>a city street with many buildings and buildings<br><br>a view of a city intersection in the evening | a white toilet sitting underneath a shower curtain<br>• a white toilet sitting underneath a bathroom window<br>a bathroom with a shower curtain open and a toilet in it |
| Base-samp | a group of people sitting at a table with pizza<br>a group of people sitting at a table with pizza and drinks<br>a group of people sitting at a table with pizza | a traffic light in a city with tall buildings<br>a traffic light in a city with tall buildings<br><br>• a traffic light and a street sign on a city street | • a bathroom with a toilet and a shower<br>• a bathroom with a toilet and a shower<br><br>a white toilet sitting next to a shower in a bathroom |

Figure 5: Comparing 3 captions sampled from adversarial model to the baseline model. Bi-grams which are top-100 frequent bi-grams in the training set are highlighted in red. Captions which are replicas from training set are marked with a •.