

AI504: Programming for Artificial Intelligence

Week 11: Image-to-Text

Edward Choi

Grad School of AI

edwardchoi@kaist.ac.kr

Today's Topic

- Image-to-text (a.k.a Image captioning)
- Show and Tell
- Show, Attend and Tell
- Text-to-Image

Image Captioning

Image-to-Text

- Sequence to sequence
 - Text in, text out
 - e.g. Translate French to English
- Image to sequence
 - Image in, text out
 - e.g. Describe a given image in text (i.e. Image Captioning)

Image Captioning

A person riding a motorcycle on a dirt road.



A group of young people playing a game of frisbee.



A herd of elephants walking across a dry grass field.



Encoder-Decoder Architecture

- Sequence to sequence
 - Encoder: RNN
 - Decoder: RNN
- Image to sequence
 - Encoder: ???
 - Decoder: ???

Encoder-Decoder Architecture

- Sequence to sequence
 - Encoder: RNN
 - Decoder: RNN
- Image to sequence
 - Encoder: CNN
 - Decoder: RNN

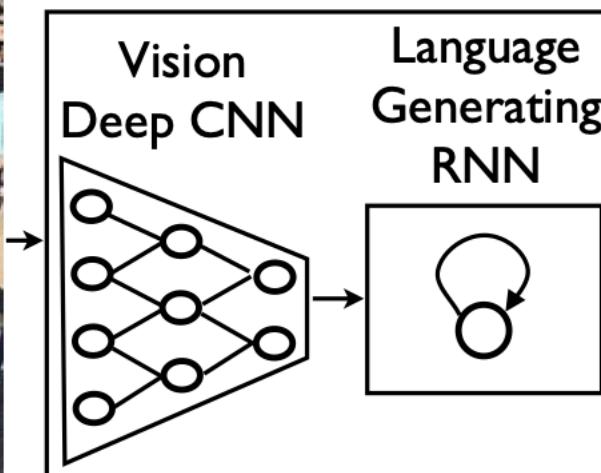
Show and Tell

Show and Tell

- Show and Tell: A Neural Image Caption Generator
 - Vinyals et al. CVPR 2015
- First paper to perform neural image captioning without any domain knowledge
 - No object detection, language modeling, description templates
 - Not text ranking, but pure generation
 - End-to-end training

Show and Tell

- Very simple architecture

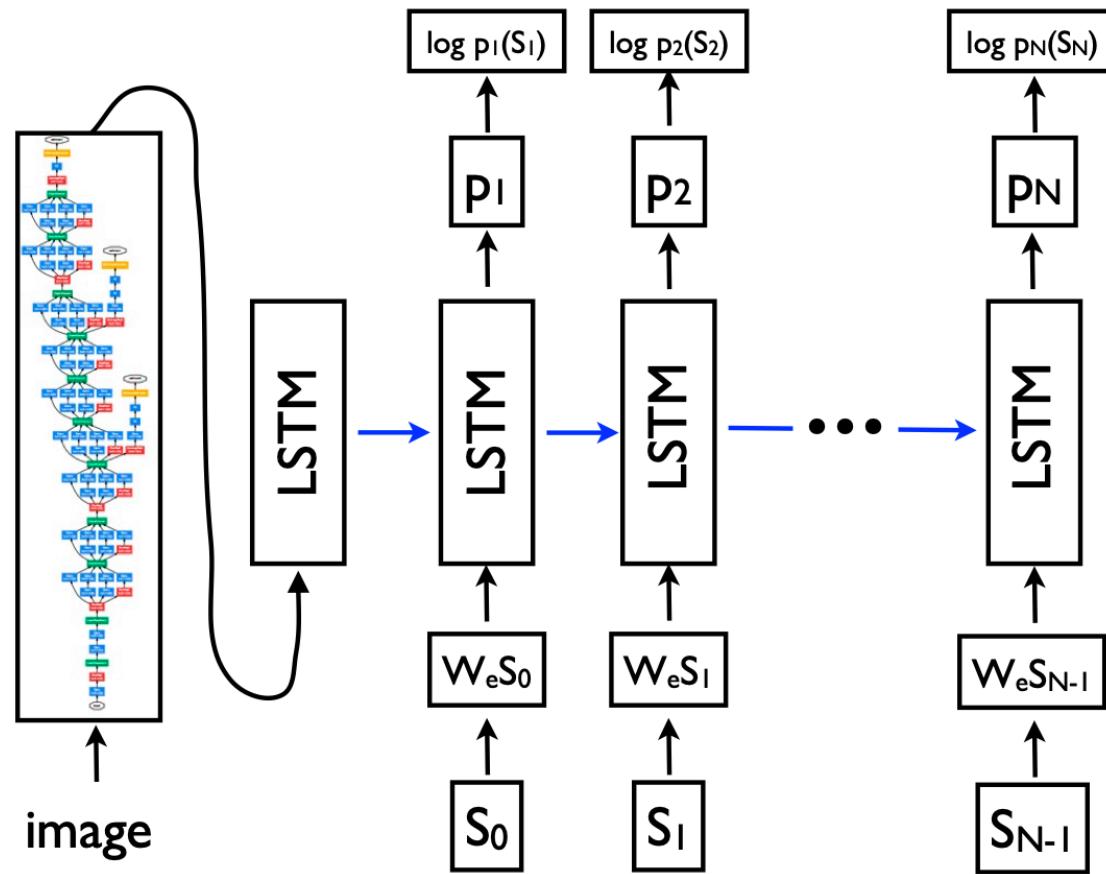


A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.

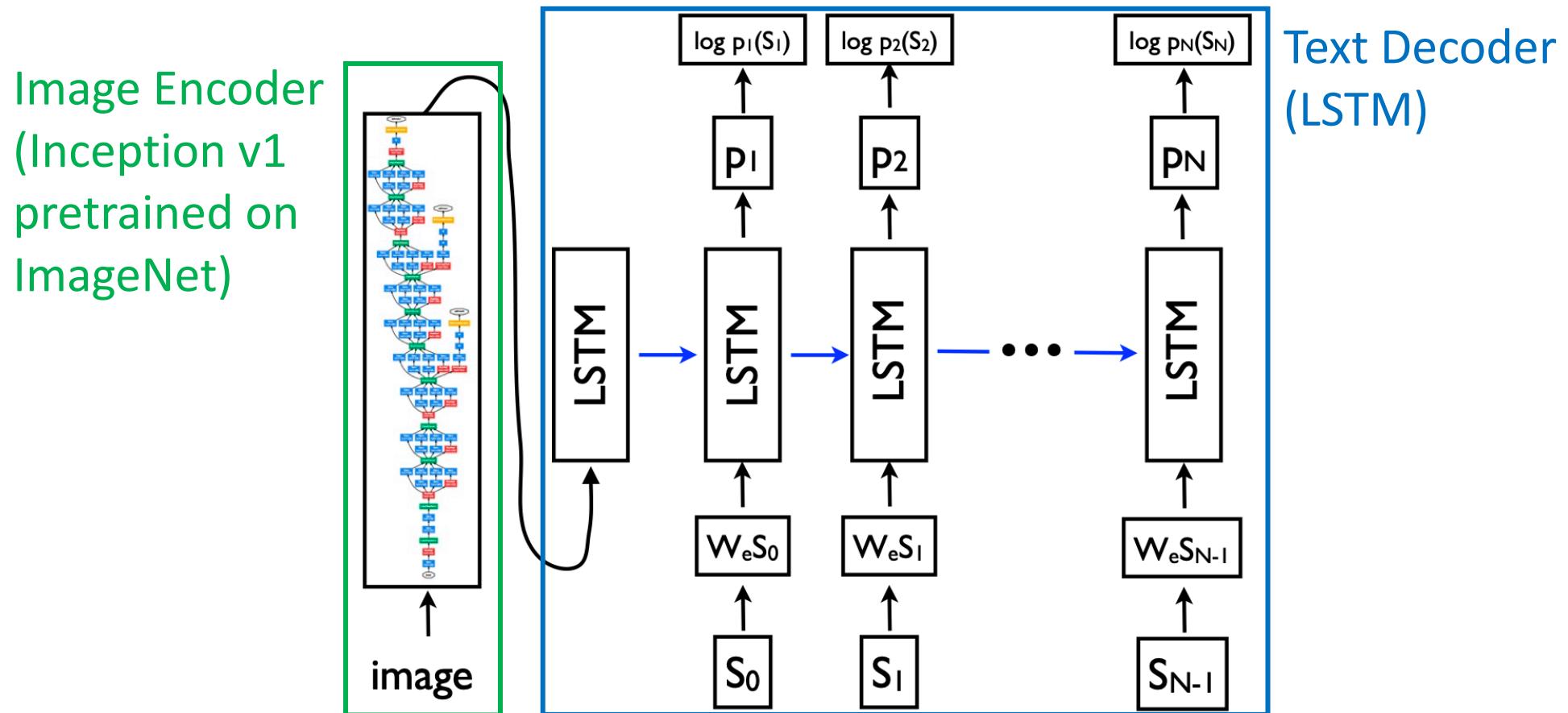
Show and Tell

- A bit more detailed architecture depiction



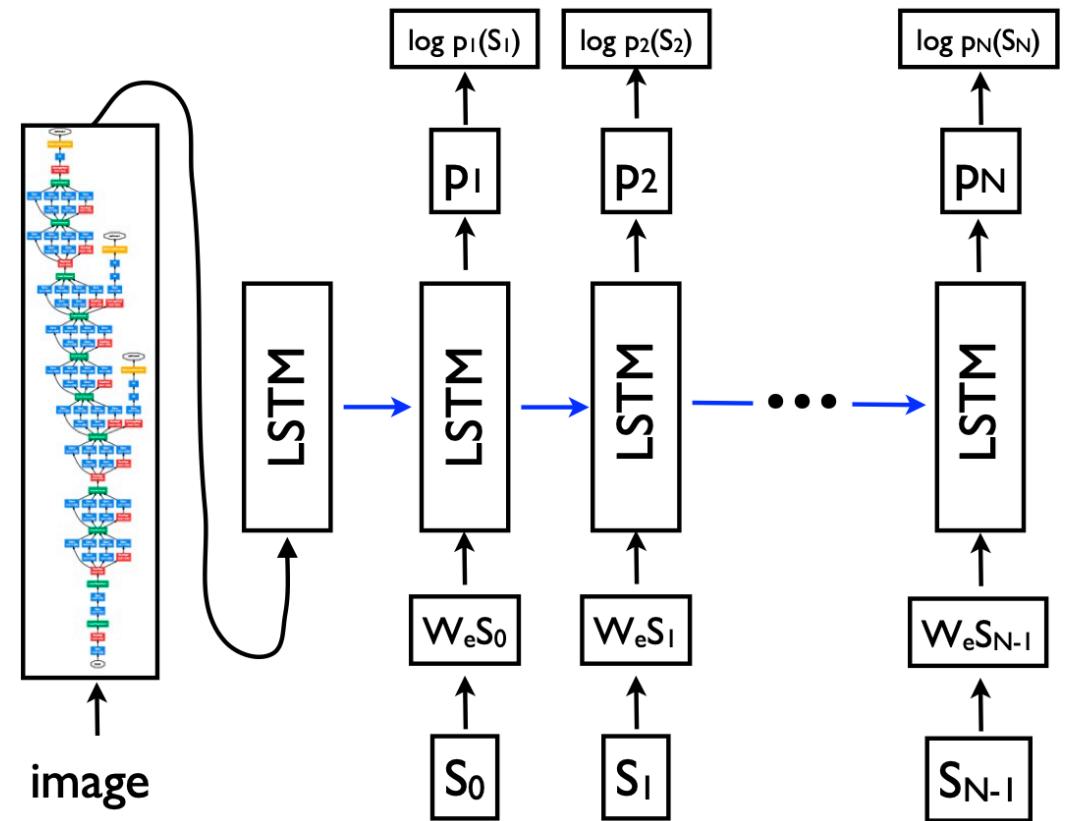
Show and Tell

- A bit more detailed architecture depiction



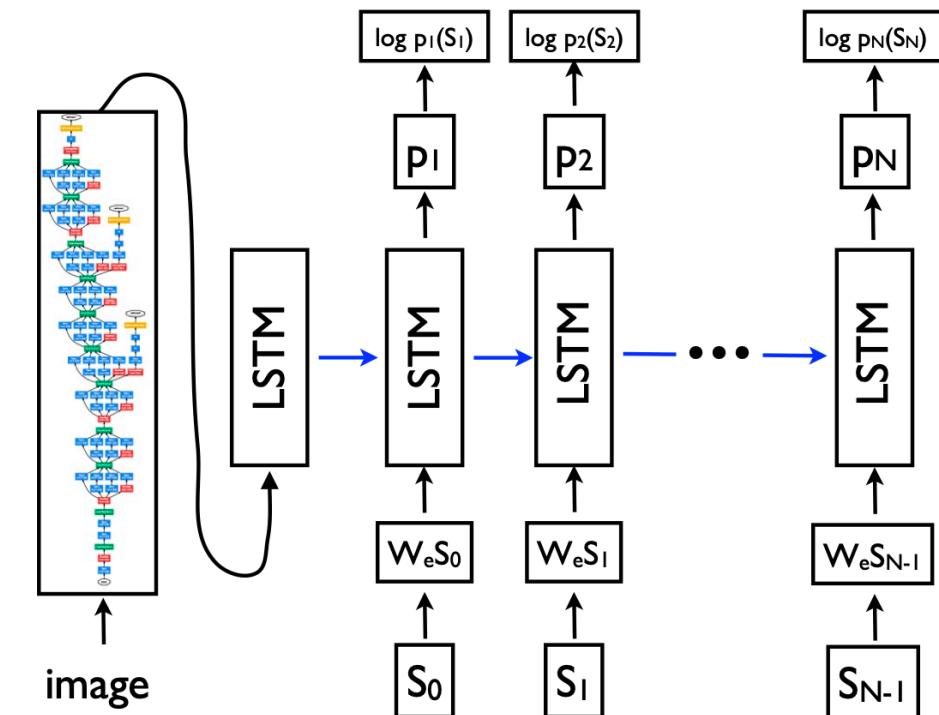
Show and Tell

- Each S_i is predicted based on p_i
 - $S_i = \text{Softmax}(W_s p_i + b)$
- Each p_i is derived based on p_{i-1}, S_{i-1}
 - $p_i = \text{RNN}(p_{i-1}, W_e S_{i-1})$
- W_e = Word embedding
- $S_{-1} = \text{CNN}(\text{Image})$
- $S_0: \text{<START>}, S_N: \text{<END>}$



Show and Tell

- **Some technical details**
- 512 embedding size & RNN size
 - Output of CNN is also 512-dimensional
- Image embedding is “fed” into LSTM at time -1
 - Not used to initialized the LSTM hidden vector.
 - Hidden layers are probably initialized to 0
- Pretrained word embeddings didn’t help much
 - Specifically, Word2Vec
- Beam search is used with beam size 20
- Trained with negative log-likelihood



Popular Datasets

Dataset name	size		
	train	valid.	test
Pascal VOC 2008 [6]	-	-	1000
Flickr8k [26]	6000	1000	1000
Flickr30k [33]	28000	1000	1000
MSCOCO [20]	82783	40504	40775
SBU [24]	1M	-	-

Model Performance

Metric	BLEU-4	METEOR	CIDER
NIC	27.7	23.7	85.5
Random	4.6	9.0	5.1
Nearest Neighbor	9.9	15.7	36.5
Human	21.7	25.2	85.4

Table 1. Scores on the MSCOCO development set.

Approach	PASCAL (xfer)	Flickr 30k	Flickr 8k	SBU
Im2Text [24]				11
TreeTalk [18]				19
BabyTalk [16]	25			
Tri5Sem [11]			48	
m-RNN [21]		55	58	
MNLM [14] ⁵		56	51	
SOTA	25	56	58	19
NIC	59	66	63	28
Human	69	68	70	

Table 2. BLEU-1 scores. We only report previous work results when available. SOTA stands for the current state-of-the-art.

Evaluation Results (grouped by human rating)

A person riding a motorcycle on a dirt road.



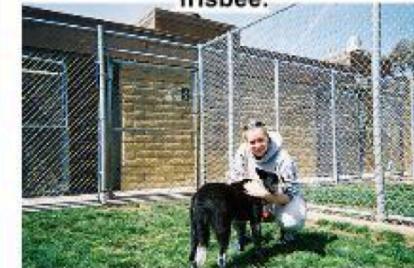
Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



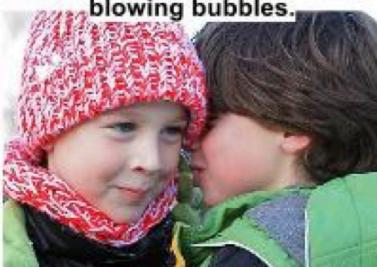
A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image

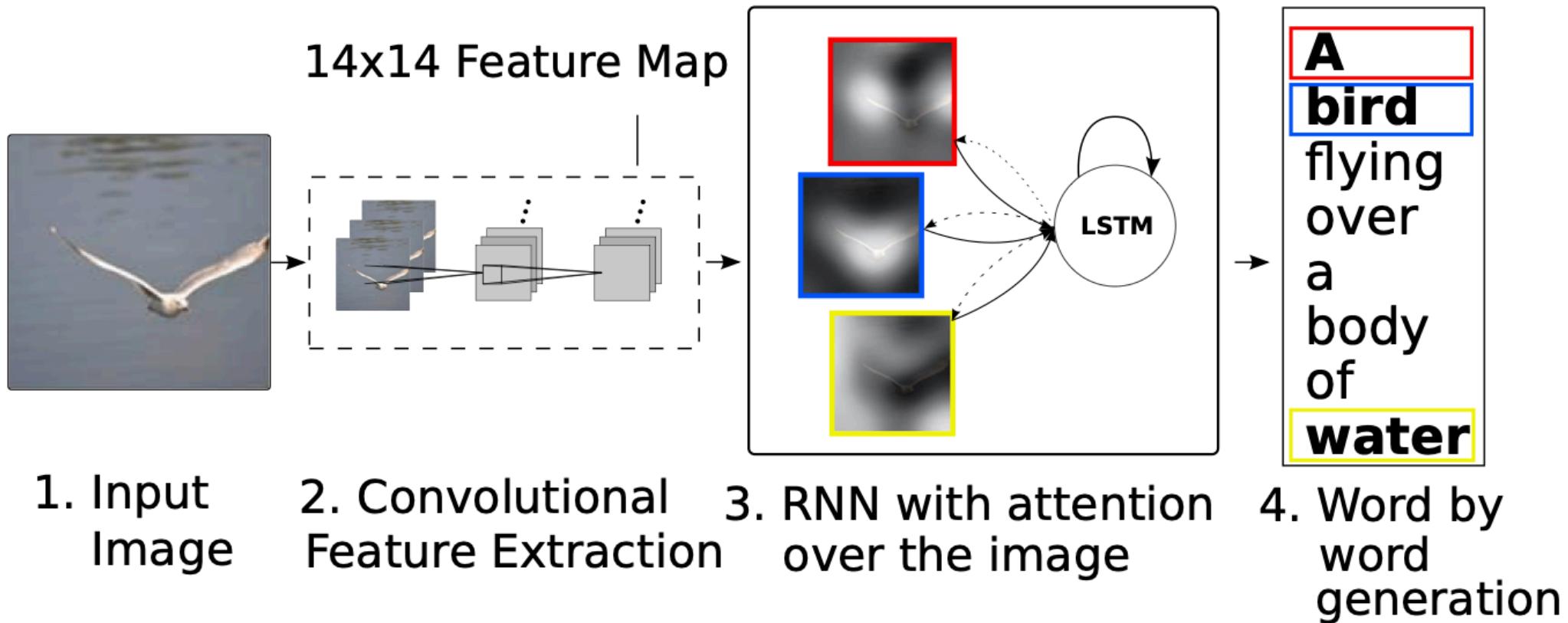
Show, Attend and Tell

Show, Attend and Tell

- Show, Attend and Tell: Neural Image Caption Generation with Visual Attention
 - Xu et al. ICML 2015
- Mixing attention mechanism with image captioning

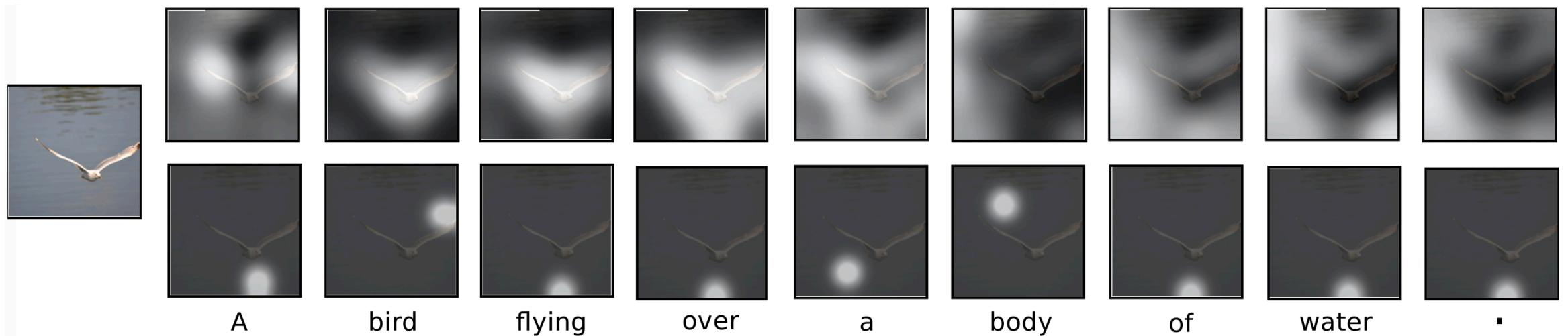
Show, Attend and Tell

- High-level architecture



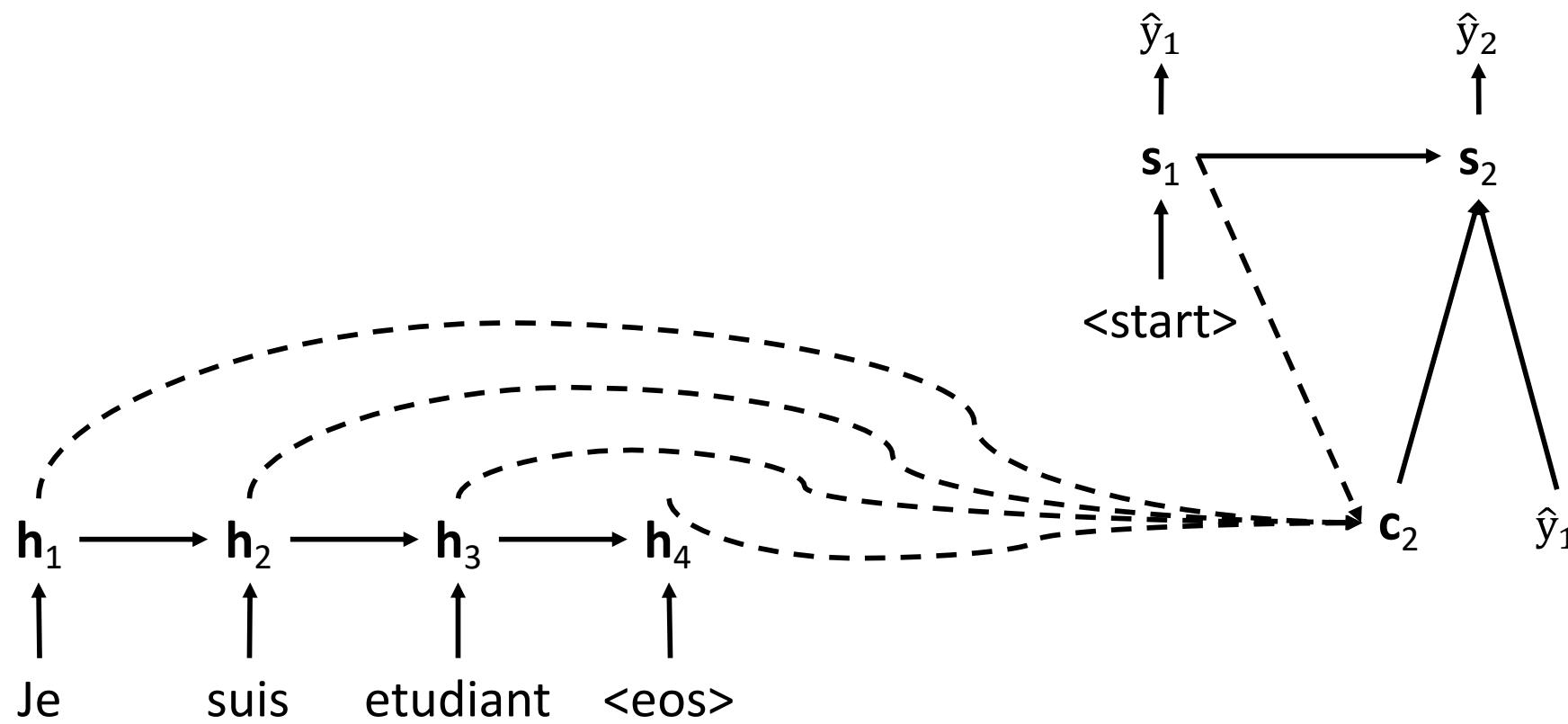
Show, Attend and Tell

- Example: “A bird flying over a body of water .”
 - Top row is “soft” attention, bottom row is “hard” attention.
- Model is “attending” to relevant part of image when generating word



Encoder-Decoder Architecture

- Seq2seq with attention



Encoder-Decoder Architecture

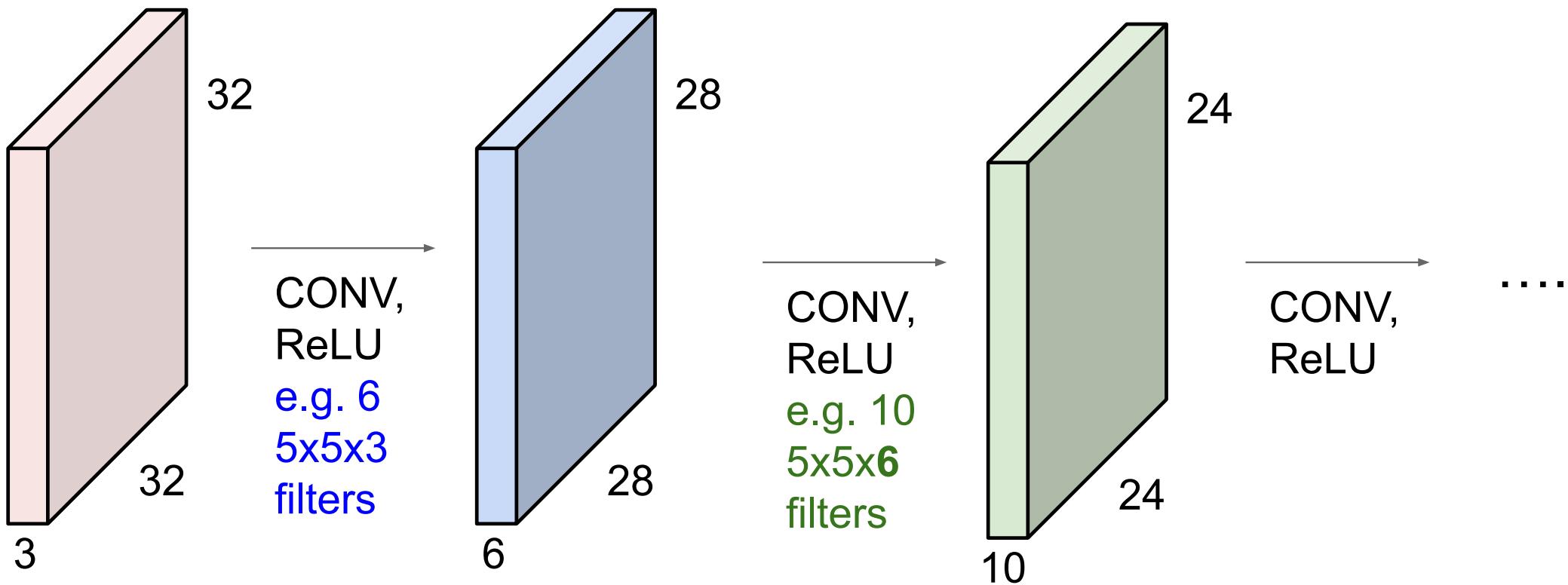
- **What we need:**
- Encoder to obtain image representation
- Decoder to generate caption
- Attention module to calculate attention weights

Encoder-Decoder Architecture

- **What we need:**
- Encoder to obtain image representation
 - Oxford VGGnet
- Decoder to generate caption
 - LSTM
- Attention module to calculate attention weights
 - MLP

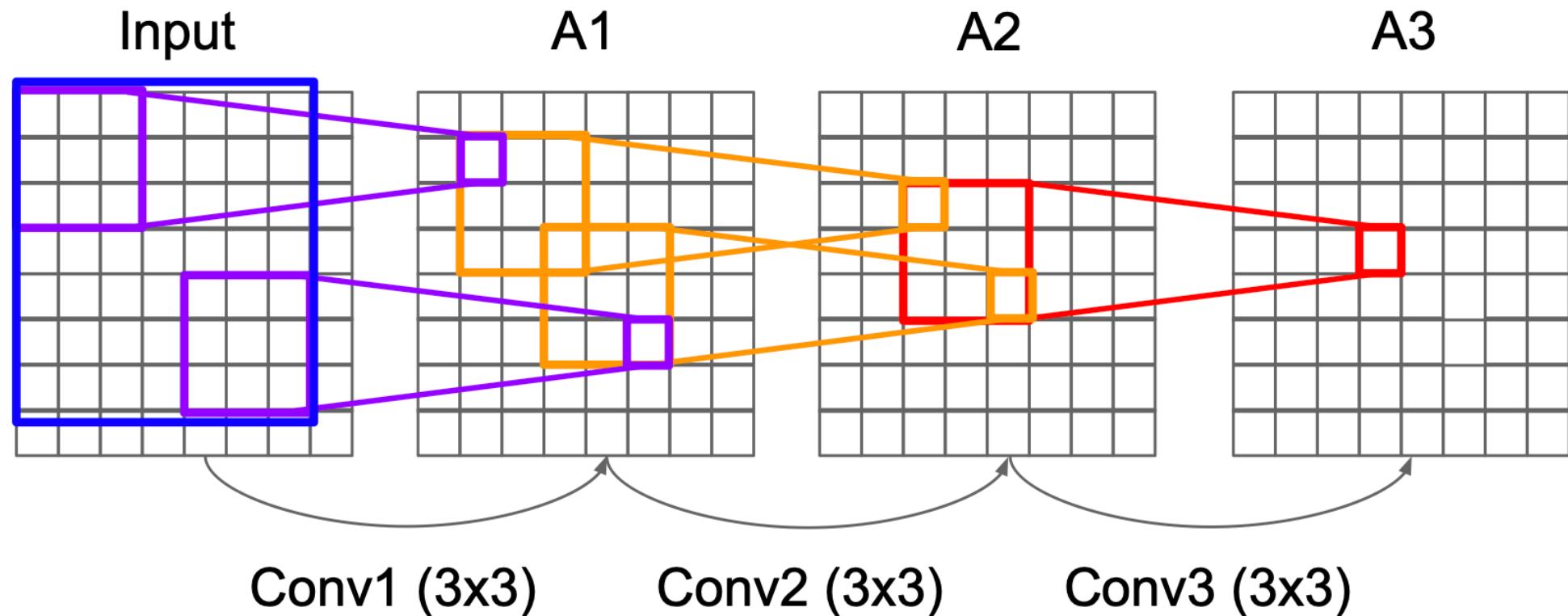
How to Attend to Part of Image

- Remember Convolution?



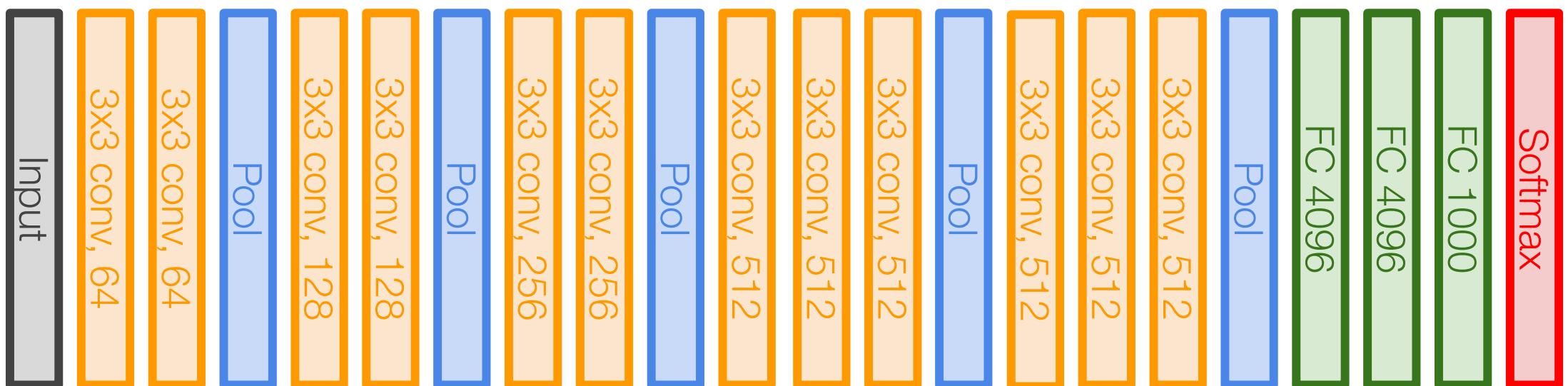
How to Attend to Part of Image

- Remember receptive field?



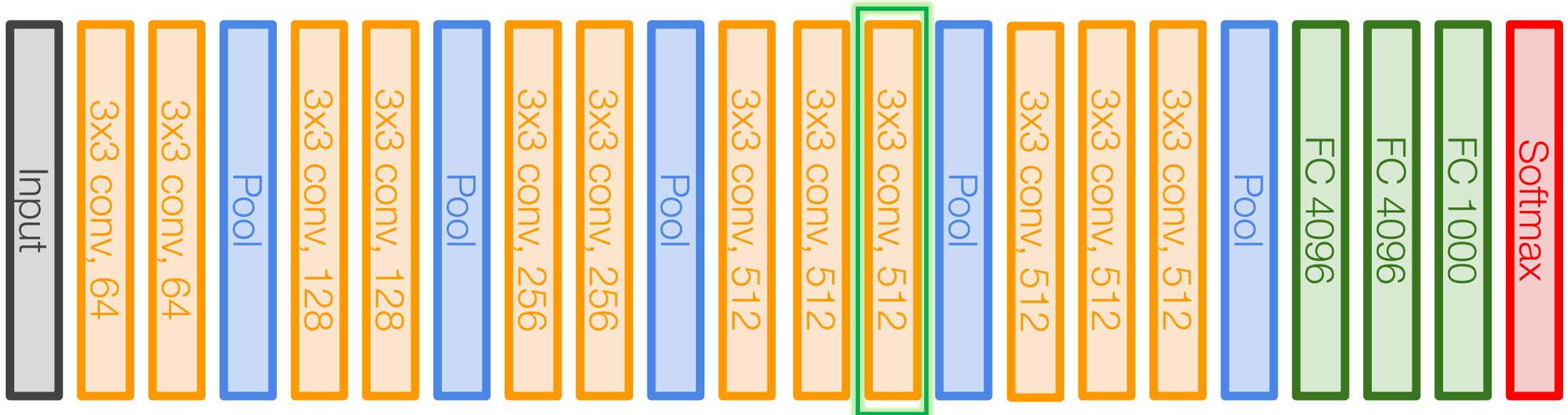
How to Attend to Part of Image

- Remember VGG 16?



How to Attend to Part of Image

- Remember VGG 16?



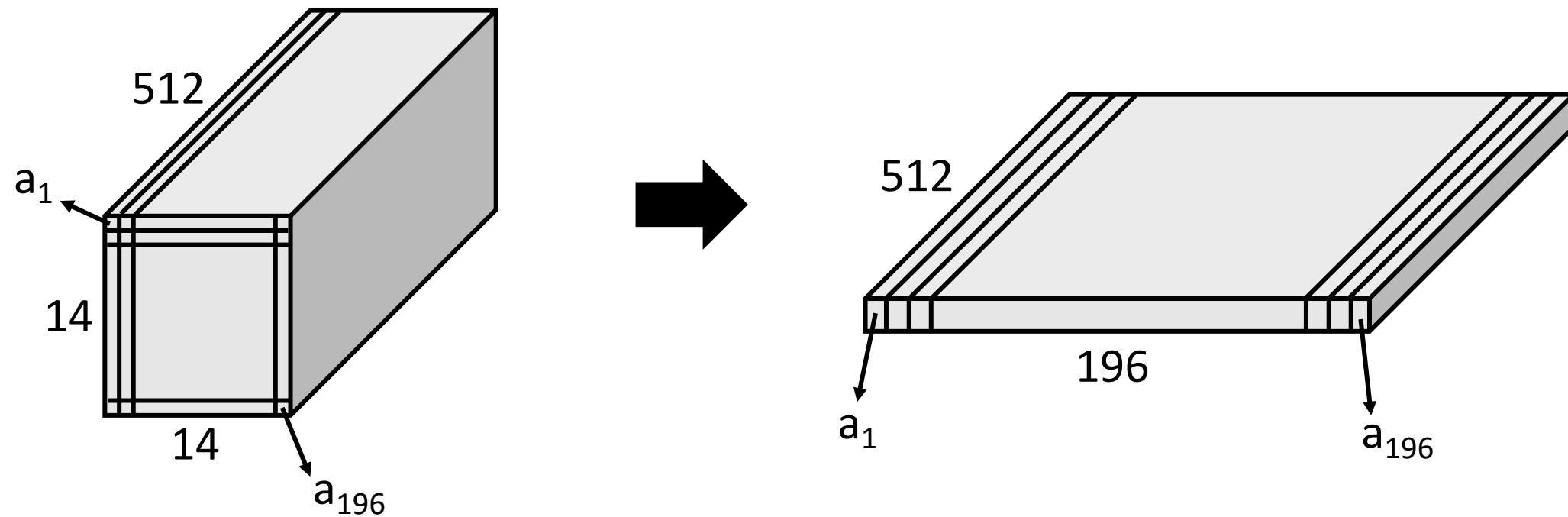
Output of this convolution layer:

$14 \times 14 \times 512$ feature map

→ 196×512 image representation vector

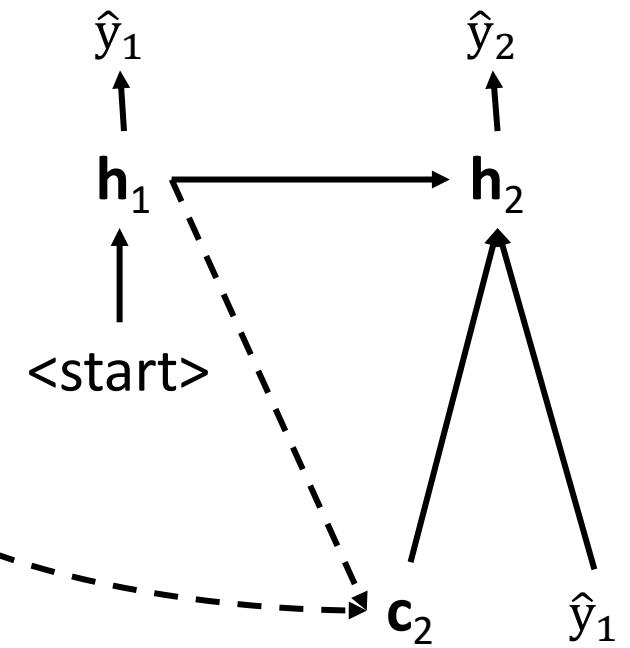
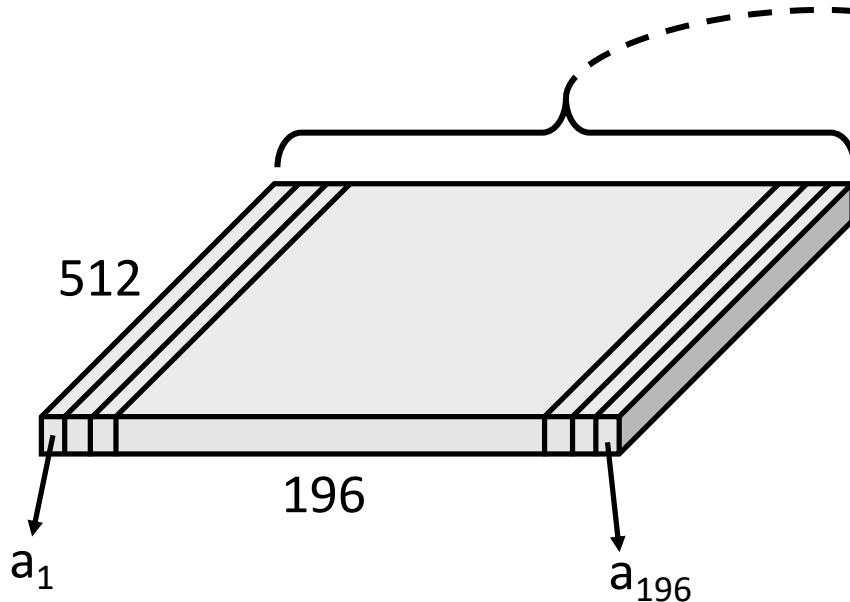
Model Architecture

- Flattening the image feature maps



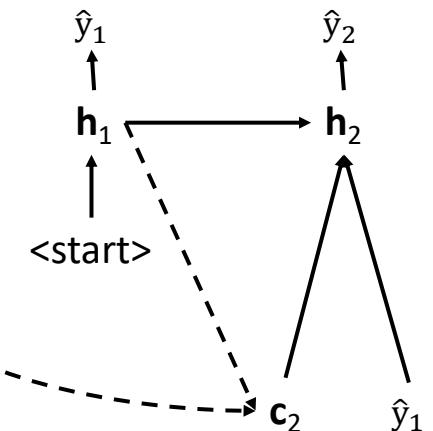
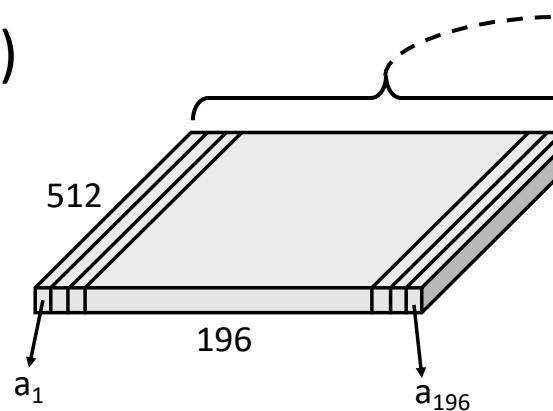
Show, Attend and Tell

- Each y_i is predicted based on h_i
- Each h_i is derived based on h_{i-1} , y_{i-1} , c_i
- c_i is derived from h_{i-1} and $a_{1:196}$



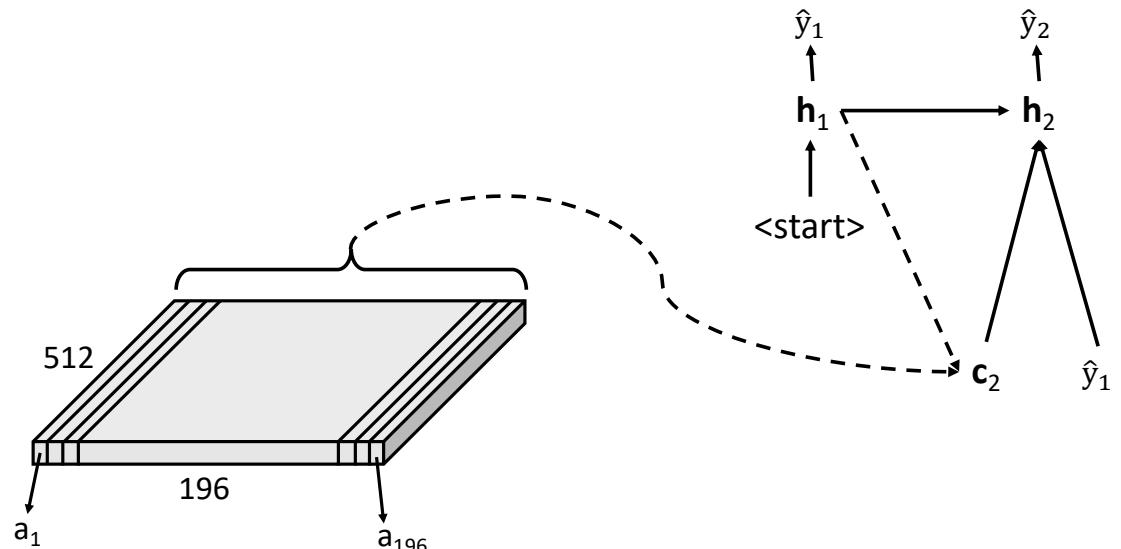
Show, Attend and Tell

- Each y_i is predicted based on h_i
 - $\hat{y}_1 = \text{Softmax}(W_w h_i + b)$
 - Each h_i is derived based on h_{i-1} , y_{i-1} , c_i
 - $h_i = \text{RNN}(h_{i-1}, [y_{i-1}; c_i]_{\text{concat}})$
 - c_i is derived from h_{i-1} and $a_{1:196}$
 - $c_i = \text{sum}(\alpha_i * a_i)$
 - $\alpha_i = \text{Softmax}(f(h_{i-1}, a_1), \dots, f(h_{i-1}, a_{196}))$
 - $f(h_{i-1}, a_j) = h_{i-1}^T W_f a_j$



Show, Attend and Tell

- **Some technical details**
- RNN's initial hidden state is learned
 - $h_0 = \text{MLP}\left(\frac{1}{L} \sum_{i=1}^L a_{1:L}\right)$
- Authors also tried “hard” attention.
 - Stochastically select only one a_i at each step.
 - Use reinforcement learning to train.
- Encourage $\sum_t \alpha_{ti} \approx 1$
 - Make the model pay equal attention to every part of image during text generation.



Model Performance

Dataset	Model	BLEU				METEOR
		BLEU-1	BLEU-2	BLEU-3	BLEU-4	
Flickr8k	Google NIC(Vinyals et al., 2014) $^{\dagger\Sigma}$	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) $^{\circ}$	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	67	44.8	29.9	19.5	18.93
	Hard-Attention	67	45.7	31.4	21.3	20.30
Flickr30k	Google NIC $^{\dagger\circ\Sigma}$	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	18.49
	Hard-Attention	66.9	43.9	29.6	19.9	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) a	—	—	—	—	20.41
	MS Research (Fang et al., 2014) $^{\dagger a}$	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) $^{\circ}$	64.2	45.1	30.4	20.3	—
	Google NIC $^{\dagger\circ\Sigma}$	66.6	46.1	32.9	24.6	—
	Log Bilinear $^{\circ}$	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	23.90
	Hard-Attention	71.8	50.4	35.7	25.0	23.04

Correction Attention Examples

Figure 3. Examples of attending to the correct object (white indicates the attended regions, *underlines* indicated the corresponding word)



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



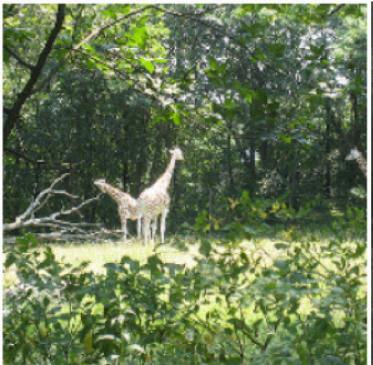
A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Incorrect Attention Examples

Figure 5. Examples of mistakes where we can use attention to gain intuition into what the model saw.



A large white bird standing in a forest.



A woman holding a clock in her hand.



A man wearing a hat and a hat on a skateboard.



A person is standing on a beach with a surfboard.



A woman is sitting at a table with a large pizza.



A man is talking on his cell phone while another man watches.

A woman is throwing a frisbee in a park.



Text-to-Image

Text-to-Image

- Generative Adversarial Text to Image Synthesis
 - Reed et al. ICML 2016
- Text-conditioned image generation with GAN

this small bird has a pink breast and crown, and black primaries and secondaries.



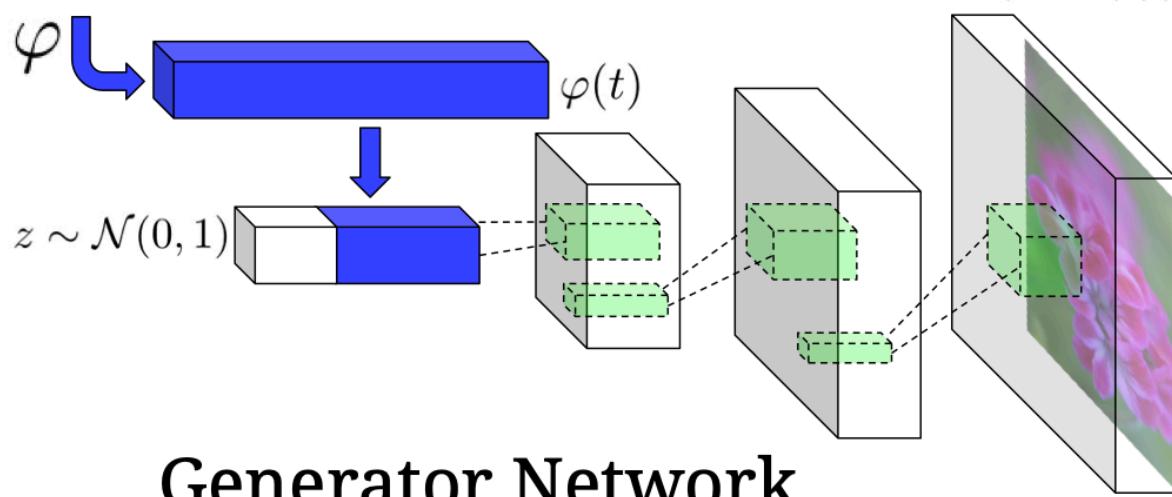
the flower has petals that are bright pinkish purple with white stigma



Model Architecture

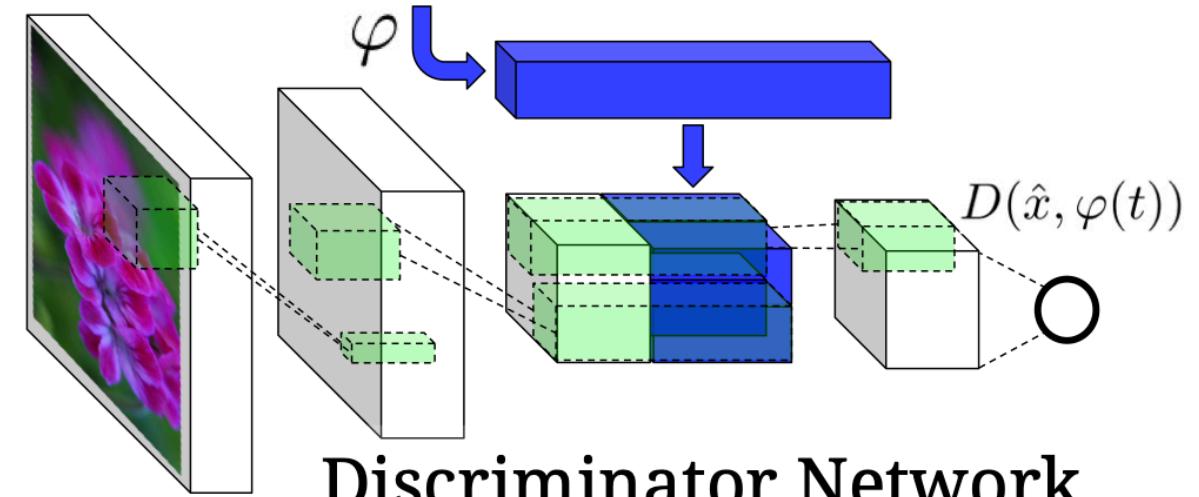
- Encode text with RNN
- Decode (i.e. generate) image with GAN
 - Use deconvolution (like DC-GAN) to upsample.

This flower has small, round violet petals with a dark purple center



Generator Network

This flower has small, round violet petals with a dark purple center



Discriminator Network

Training Strategy

- Discriminator's job is complicated
 - Real image with right text? \rightarrow Real!
 - Fake image with right text? \rightarrow Fake!
 - Real image with wrong text? \rightarrow Fake!
 - Fake image with wrong text? \rightarrow Fake!
- Discriminator is fed three cases
 - Real image, right text
 - Real image, wrong text
 - Fake image, right text

Algorithm 1 GAN-CLS training algorithm with step size α , using minibatch SGD for simplicity.

- 1: **Input:** minibatch images x , matching text t , mis-matching \hat{t} , number of training batch steps S
- 2: **for** $n = 1$ **to** S **do**
- 3: $h \leftarrow \varphi(t)$ {Encode matching text description}
- 4: $\hat{h} \leftarrow \varphi(\hat{t})$ {Encode mis-matching text description}
- 5: $z \sim \mathcal{N}(0, 1)^Z$ {Draw sample of random noise}
- 6: $\hat{x} \leftarrow G(z, h)$ {Forward through generator}
- 7: $s_r \leftarrow D(x, h)$ {real image, right text}
- 8: $s_w \leftarrow D(x, \hat{h})$ {real image, wrong text}
- 9: $s_f \leftarrow D(\hat{x}, h)$ {fake image, right text}
- 10: $\mathcal{L}_D \leftarrow \log(s_r) + (\log(1 - s_w) + \log(1 - s_f))/2$
- 11: $D \leftarrow D - \alpha \partial \mathcal{L}_D / \partial D$ {Update discriminator}
- 12: $\mathcal{L}_G \leftarrow \log(s_f)$
- 13: $G \leftarrow G - \alpha \partial \mathcal{L}_G / \partial G$ {Update generator}
- 14: **end for**

Examples

GT an all black bird with a distinct thick, rounded bill.



this small bird has a yellow breast, brown crown, and black superciliary



a tiny bird, with a tiny beak, tarsus and feet, a blue crown, blue coverts, and black cheek patch



this bird is different shades of brown all over with white and black spots on its head and back



the gray bird has a light grey head and grey webbed feet



GAN



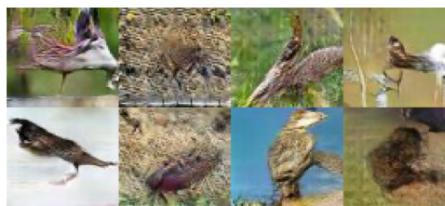
GAN - CLS



GAN - INT



GAN - INT - CLS



Examples

GT
this flower is white and pink in color, with petals that have veins.



these flowers have petals that start off white in color and end in a dark purple towards the tips.



bright droopy yellow petals with burgundy streaks, and a yellow stigma.



a flower with long pink petals and raised orange stamen.



the flower shown has a blue petals with a white pistil in the center



GT



GAN



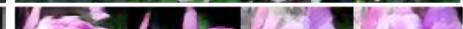
GAN - CLS



GAN - INT



GAN - INT - CLS



Examples

	GT	Ours		GT	Ours		GT	Ours
a group of people on skis stand on the snow.			a man in a wet suit riding a surfboard on a wave.			a pitcher is about to throw the ball to the batter.		
a table with many plates of food and drinks			two plates of food that include beans, guacamole and rice.			a picture of a very clean living room.		
two giraffe standing next to each other in a forest.			a green plant that is growing out of the ground.			a sheep standing in a open grass field.		
a large blue octopus kite flies above the people having fun at the beach.			there is only one horse in the grassy field.			a toilet in a small room with a window and unfinished walls.		

AI504: Programming for Artificial Intelligence

Week 11: Image-to-Text

Edward Choi

Grad School of AI

edwardchoi@kaist.ac.kr