

## تمرین سوم : Image Captioning with CNN-RNN Architecture

رضا حاجی علیزاده

بخش مهندسی کامپیوتر دانشگاه شهید باهنر کرمان

rezahajializadeh@eng.uk.ac.ir

### ۱. معرفی و هدف پروژه

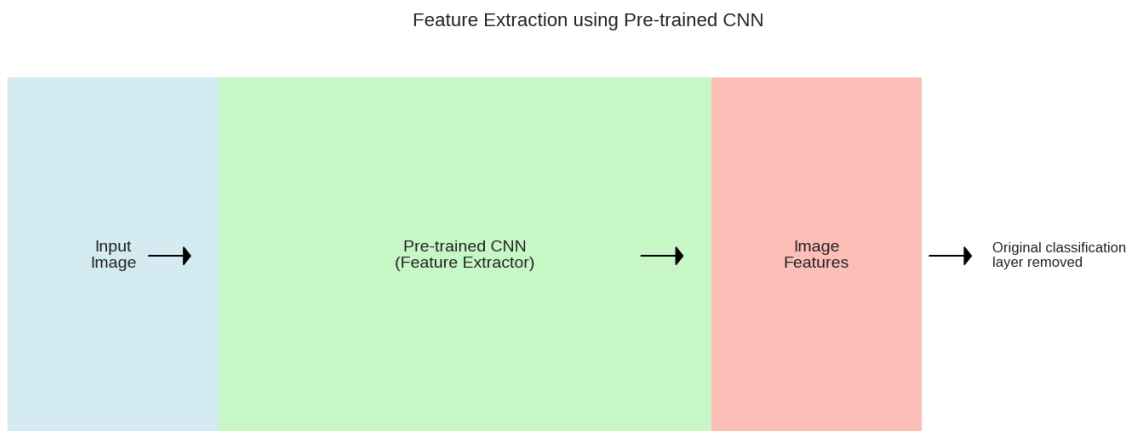
مسئله ای که در این پروژه بررسی می شود توصیف تصاویر (Image Captioning) است. هدف این است که مدلی طراحی و آموزش داده شود که با دریافت یک تصویر به عنوان ورودی بتواند جمله ای معنی دار و توصیفی از محتوای آن تصویر تولید کند. با توجه به پیچیدگی بالای این مسئله به دلیل نیاز به درک بصری دقیق از تصویر و همچنین تولید عبارت های زبان طبیعی، مدل هایی که برای حل آن استفاده می شوند نیاز به ترکیب ویژگی های شبکه های عصبی کانولوشنی (CNNs) برای استخراج ویژگی های بصری و شبکه های عصبی بازگشتی (RNNs) برای مدل سازی توالی کلمات دارند. در این پروژه از یک معماری ترکیبی مبتنی بر CNN-RNN استفاده می شود که ابتدا با استفاده از یک شبکه عصبی کانولوشنی از پیش آموزش دیده مانند ResNet یا MobileNet ویژگی ها سطح بالا از تصویر استخراج می شود و سپس این ویژگی ها به عنوان ورودی به یک شبکه عصبی بازگشتی مثل LSTM یا GRU داده می شود که وظیفه تولید گام به گام کلمات توصیف را دارد. در نهایت مدل باید قادر باشد تا با دیدن تصویر جدیدی که قبلاً مشاهده نکرده است یک جمله توصیفی جدید و معنادار تولید کند.

### ۲. جزئیات معماری مدل

مدلی که در این پروژه برای تولید توصیف تصاویر طراحی شده از یک معماری ترکیبی Encoder-Decoder استفاده می کند. این معماری شامل دو مؤلفه اصلی است:

- رمزگذار (Encoder): برای درک محتوای تصویر ابتدا باید ویژگی های بصری سطح بالای آن استخراج شوند، برای این منظور از یک شبکه عصبی کانولوشنی از پیش آموزش دیده به عنوان رمزگذار استفاده می شود. ورودی این یک تصویر و خروجی آن یک بردار ویژگی است که نشان دهنده ویژگی های بصری سطح بالا تصویر است.
- رمزگشا (Decoder): پس از استخراج ویژگی های تصویر این ویژگی ها به یک شبکه عصبی بازگشتی داده می شوند تا توصیف متنی گام به گام تولید شود.

شکل ۱ نشان دهنده معماری Encoder این مدل است، قسمت Encoder این مدل تصاویری با ابعاد  $224 \times 224$  دریافت می کند و یک بردار ویژگی 256 بعدی تولید می کند. ویژگی های استخراج شده توسط Encoder به Decoder داده می شوند تا متن توصیفی تصویر را تولید کند.



شکل ۱

### ۳. Experimental Setup

برای آموزش و ارزیابی مدل از مجموعه داده Flickr8k استفاده شده است. این مجموعه داده شامل 8091 تصویر واقعی است که برای هر تصویر 5 Caption وجود دارد و تعداد کل Caption ها 40455 است. شکل ۲ نشان دهنده 5 تصویر و Caption های آن تصاویر است. شکل ۳ نشان دهنده نمودار توضیح طول Caption ها است، کوتاه ترین Caption فقط شامل 1 کلمه است و طولانی ترین Caption شامل 36 کلمه است. اندازه Vocabulary برابر 8827 است و تعداد کلمه هایی که فقط یک بار استفاده شده اند 3608 است. شکل ۴ نشان دهنده نمودار 30 کلمه رایج در Caption ها است، کلمه ای که بیشترین تکرار را در Caption ها دارد کلمه a است. شکل ۵ نشان دهنده نمودار Vocabulary Size و Vocabulary Coverage برحسب Frequency Threshold است. Frequency Threshold نشان دهنده تعداد باری که یک کلمه باید در Caption ها ظاهر شود تا در واژگان جدید باقی بماند است، این نمودار نشان می دهد که با افزایش Frequency Threshold اندازه Vocabulary کوچک می شود و Vocabulary Coverage کاهش می یابد چون کلمات نادر تر حذف می شوند. اگر دقت اهمیت داشته باشد باید Frequency Threshold کمتری انتخاب کنیم تا اندازه Vocabulary بزرگ تر شود و Vocabulary Coverage بالاتری داشته باشیم و اگر سرعت و کارایی اهمیت داشته

باشد باید Frequency Threshold بزرگ تری انتخاب کنیم تا اندازه Vocabulary کوچک تر شود. برای آموزش مدل دیتاست را به سه دیتاست Train، Validation و Test تقسیم می کنیم که Train Set شامل 6000 تصویر و 30000 Caption است، Validation Set شامل 1000 تصویر و 5000 Caption است و Test Set شامل 1000 تصویر و 5000 Caption است. شکل ۶ نشان دهنده نمودار ابعاد تصاویر است.

A black-and-white dog bounds off the ground , all feet in the air , of a yellow field .  
 A black and white dog is jumping over high yellow grass .  
 A black and white dog is jumping through a field of brown grass .  
 a dog runs through the dry grass .  
 The black and white dog runs through the field .



A child staring at Santa .  
 A father and son looking at a funny looking Santa .  
 A little boy and his father talking to a man dressed as Santa Claws .  
 The child is looking at Santa Claus .  
 The little boy has a yellow crown and the man is wearing red velvet .



Three woman walk down a city street and one has a pink purse .  
 Three women , two with tattoos , walking down the street  
 Women are walking through the street drinking iced coffee .  
 Women walk down a buzy sidewalk .  
 Women walking down the street .

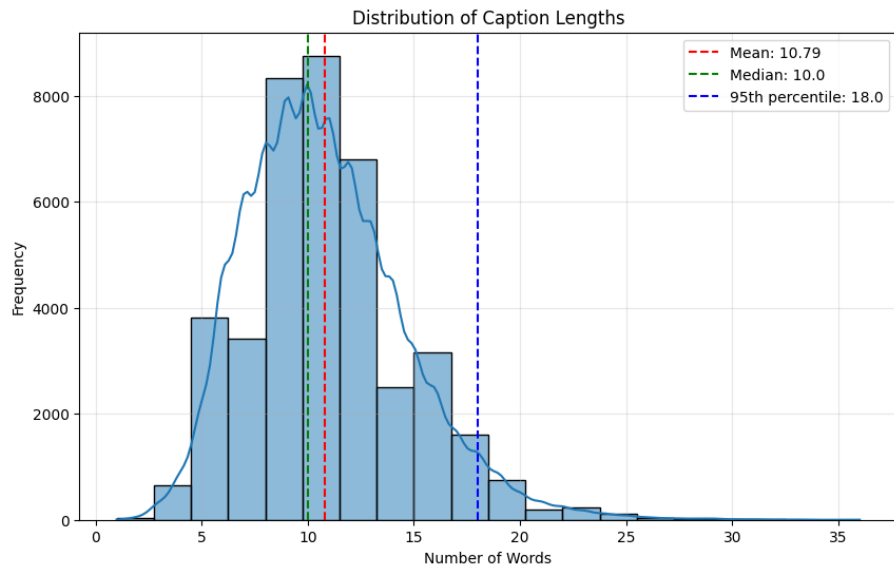


The two people are standing on a rock , holding themselves up against another rock , while looking down  
 Three people make their way through rocky terrain .  
 Three people participate in rock climbing .  
 Two female hikers hang onto a rock in front of them while looking down into a deep crevice .  
 Two people wearing backpacks and a woman in red shorts are standing on some large rock formations .

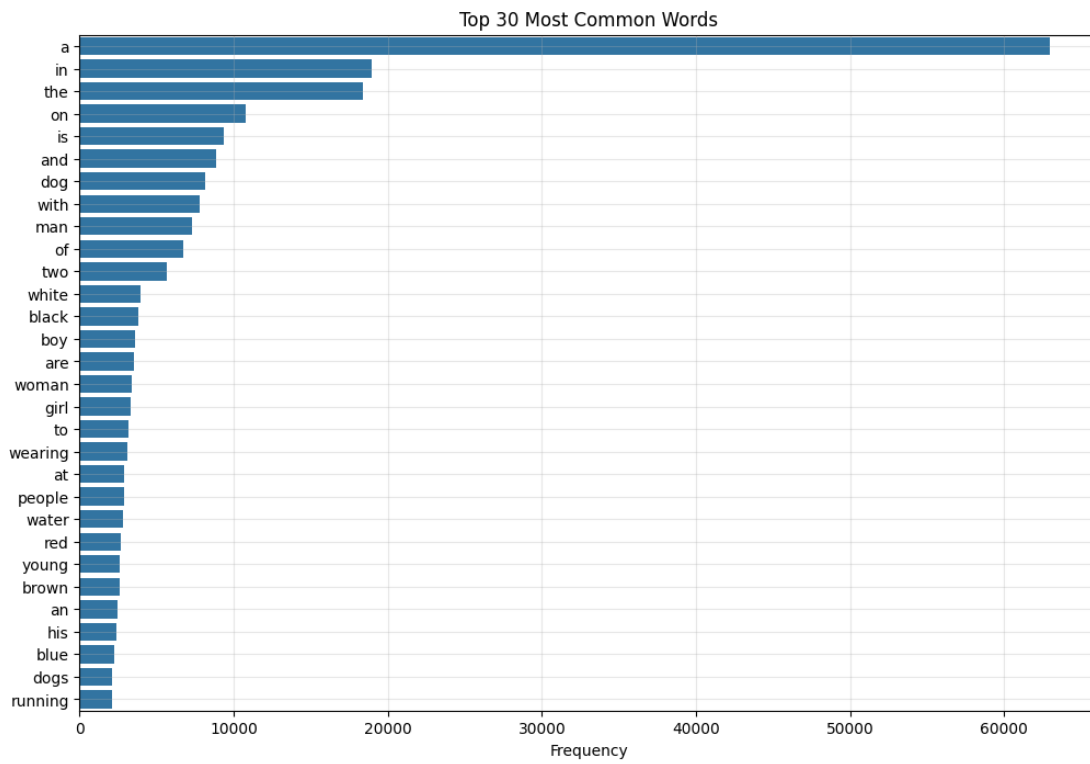


A closeup of a little girl on a swing .  
 A little girl in a dress with pink flowers swings on a red-seated swing .  
 A little girl in a pink and white flowered dress and blue sweater swinging .  
 A young girl is swinging in a backyard .  
 The little girl swings in the backyard .

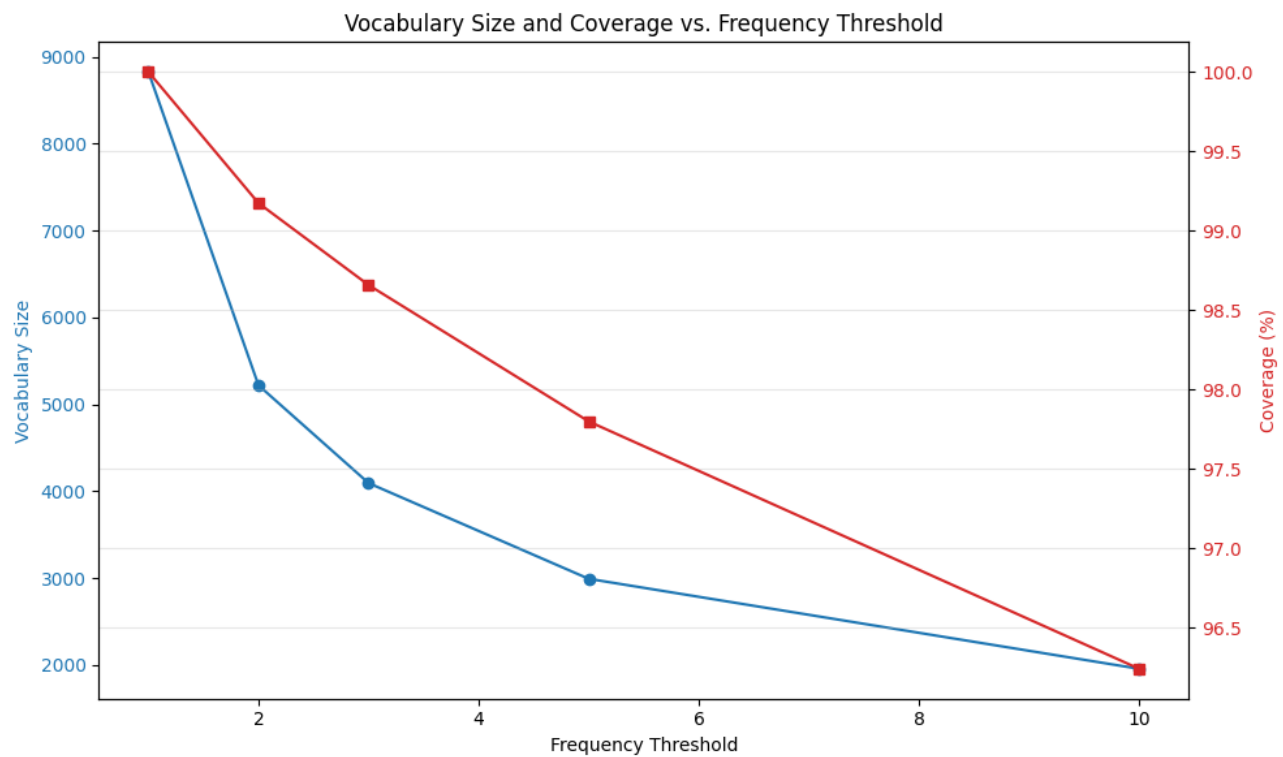




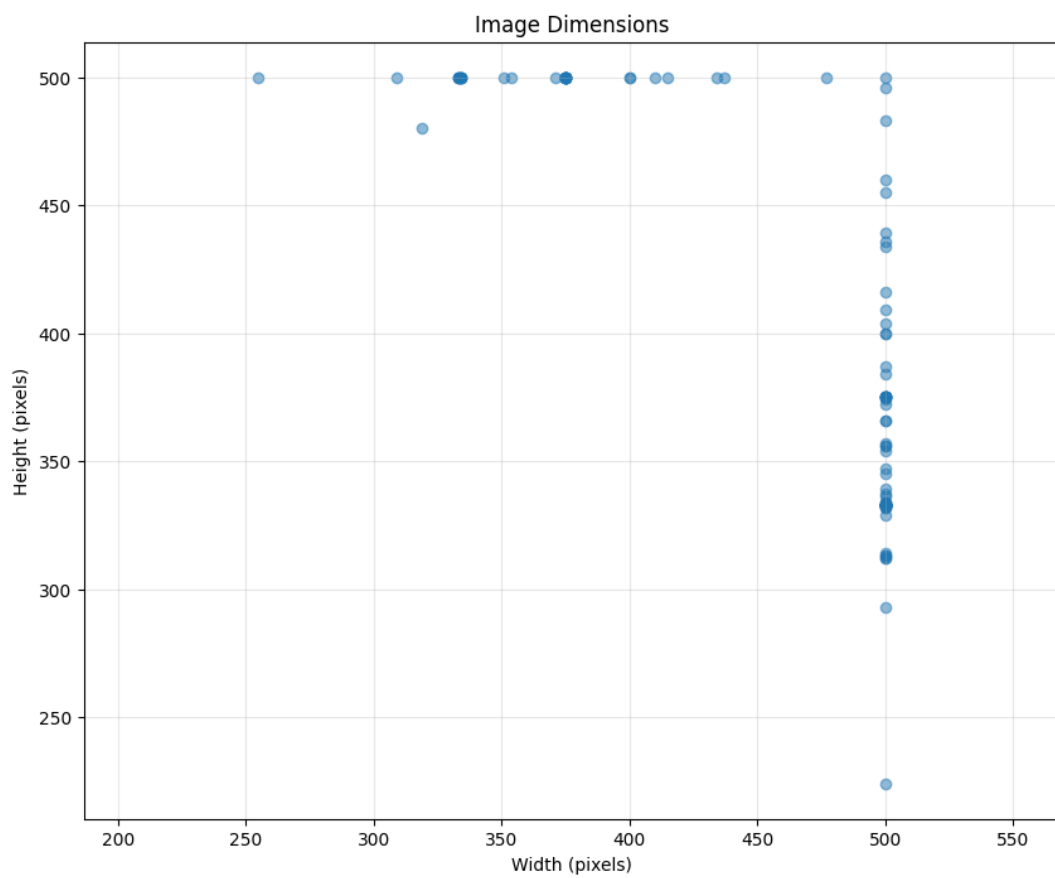
شکل ۳



شکل ۴



شکل ۵

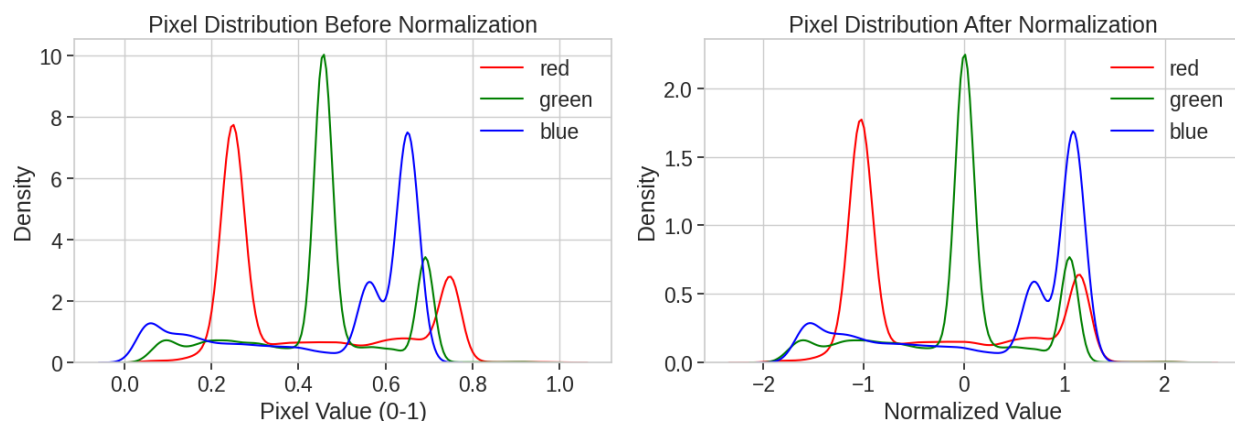


شکل ۶

شکل ۷ نشان دهنده نحوه اعمال مراحل پیش پردازش روی تصاویر است، هدف پیش پردازش تبدیل تصاویر ورودی به شکل مناسب برای شبکه های عصبی کانولوشنی است. در پیش پردازش تصویر ابتدا تصاویر Resize می شوند سپس Crop می شوند و در نهایت Normalize می شوند. شکل ۸ نشان دهنده نمودار توزیع پیکسل ها قبل و بعد از نرمال سازی است.

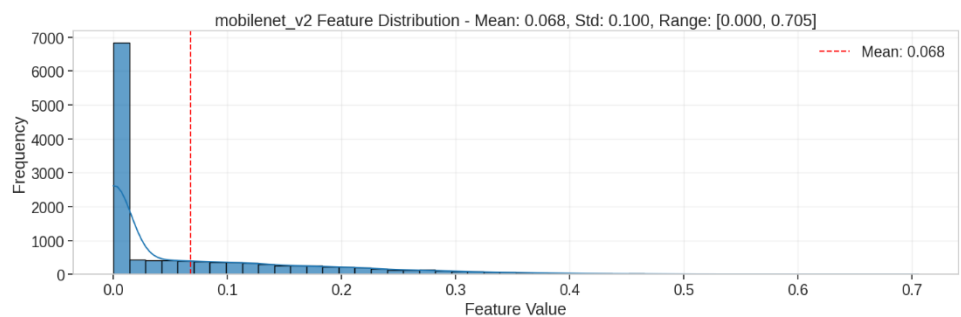
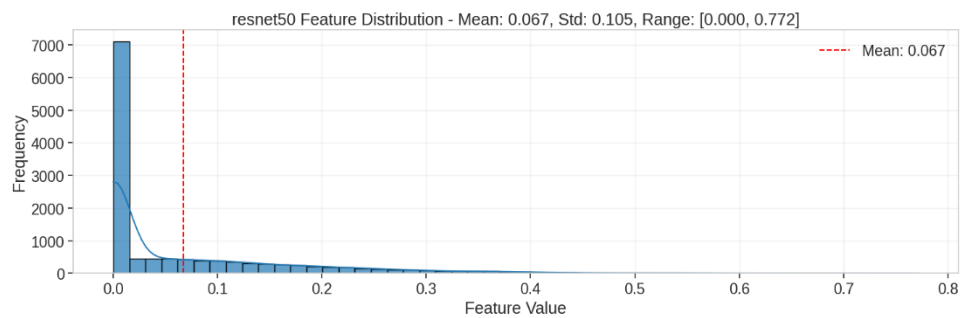
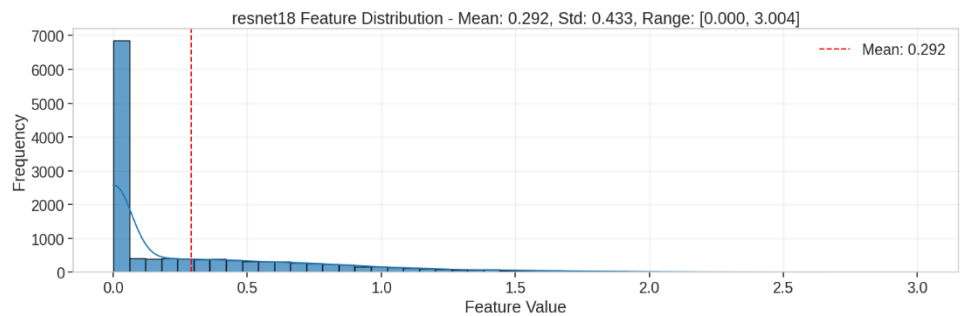


شکل ۷

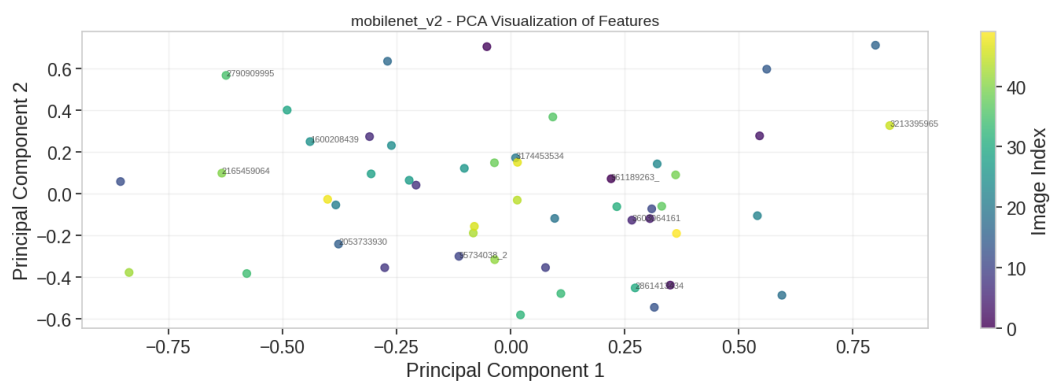
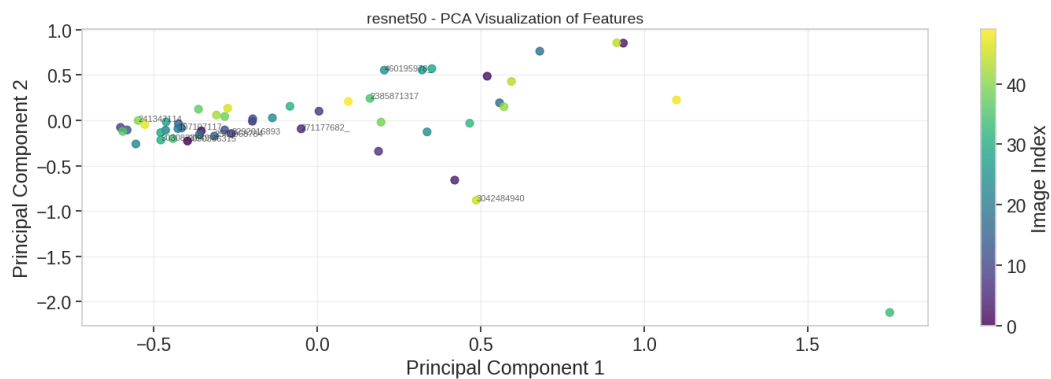
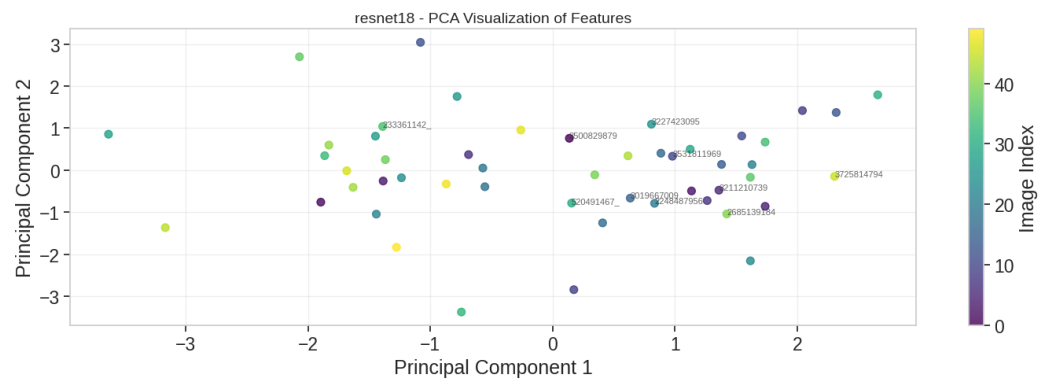


شکل ۸

برای استخراج ویژگی ها از سه مدل از پیش آموزش دیده ResNet18، ResNet50 و MobilenetV2 استفاده می کنیم. شکل ۹ نشان دهنده نمودار توزیع ویژگی های برای این سه مدل است، مقادیر ویژگی های در ResNet18 محدوده گسترده تری قرار دارند و انحراف معیار بالایی دارند اما در ResNet50 و MobilenetV2 مقادیر ویژگی های در محدوده کوچک تری قرار دارند و انحراف معیار کمتری دارند. شکل ۱۰ نشان دهنده نمودار PCA برای تحلیل ویژگی های استخراج شده توسط این مدل ها را نشان می دهد، هدف اصلی این نمودار ها مقایسه نحوه استخراج ویژگی ها توسط این مدل ها و نمایش توزیع ویژگی ها در فضای Principal Component است. شکل ۱۱ نشان دهنده نمودار Heatmap شباهت ویژگی های استخراج از تصاویر با استفاده از مدل ها است، هدف اصلی این نمودار مقایسه شباهت بین تصاویر مختلف بر اساس ویژگی های استخراج شده توسط مدل است. جدول ۱ نشان دهنده یک مقایسه کامل بین این سه مدل است.

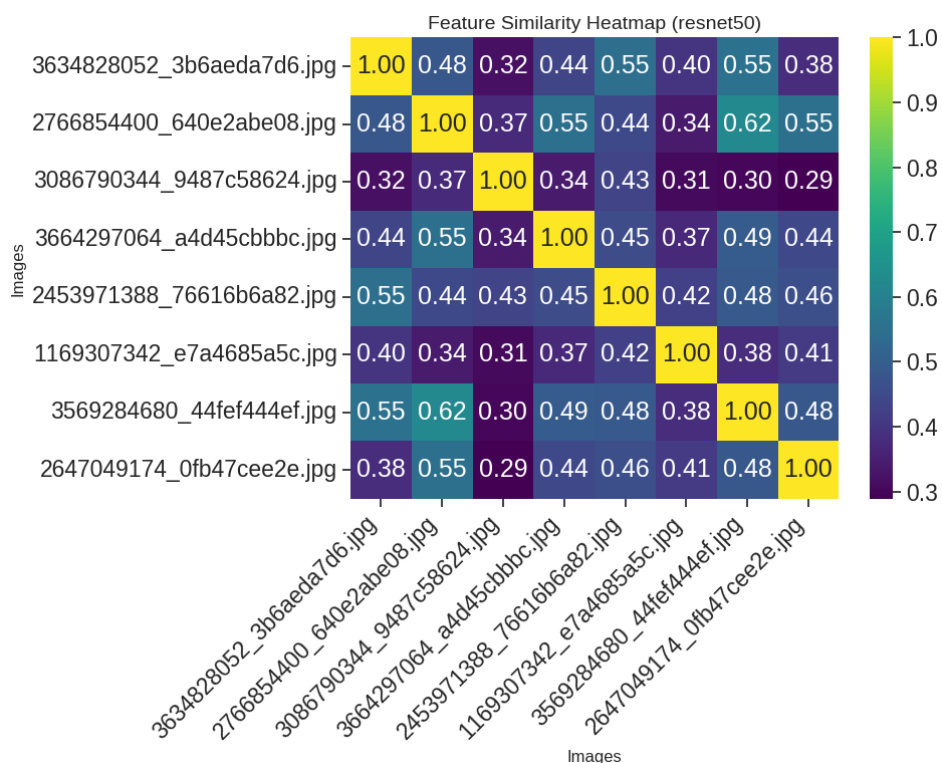
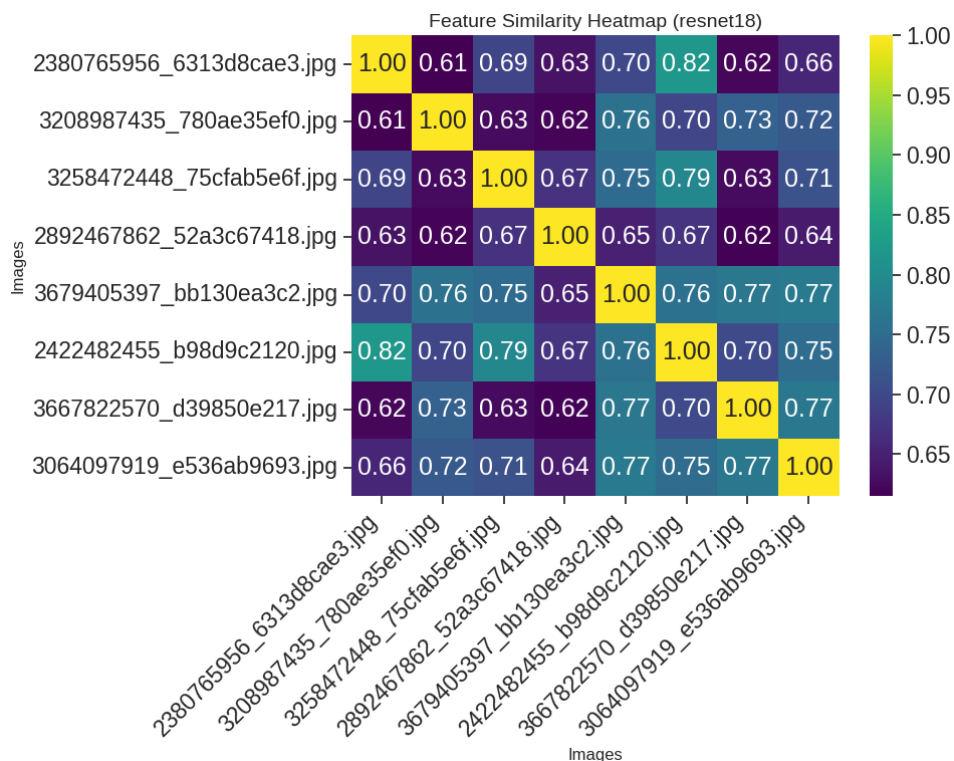


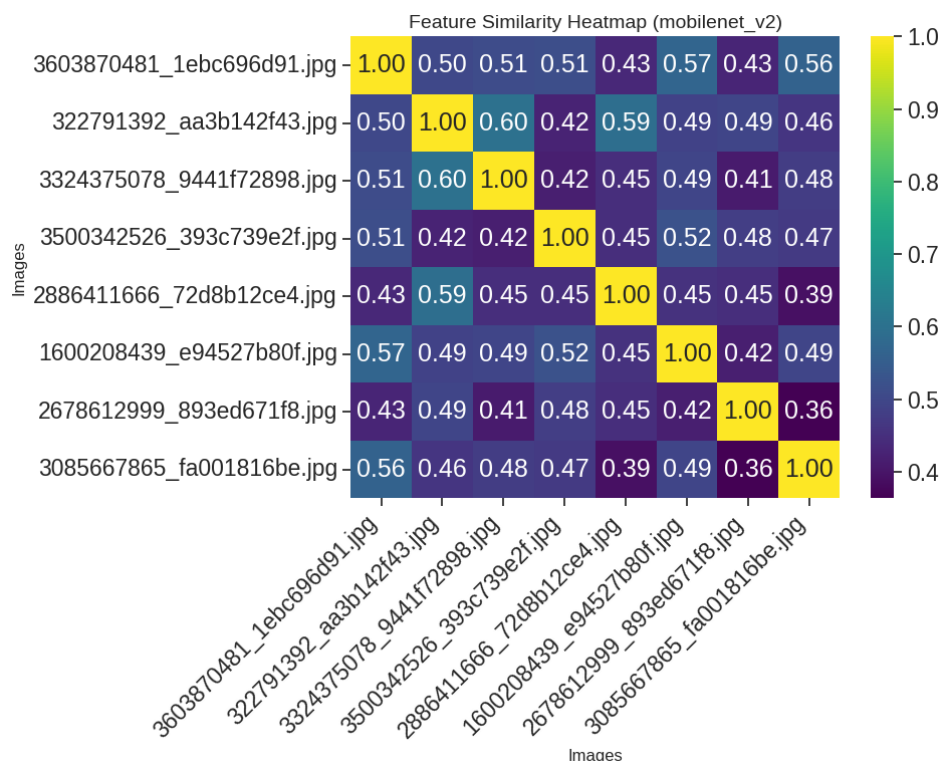
شکل ۹



شکل ۱۰







شکل ۱۱

Model	Parameters(M)	Feature Dimension	Feature Mean	Feature Std
ResNet18	11.7	512	0.2917	0.4335
ResNet50	25.6	2048	0.0669	0.1055
Mobilenet v2	3.5	1280	0.0676	0.0999s

جدول ۱

فرآیند آموزش به این صورت است که تصویر ورود

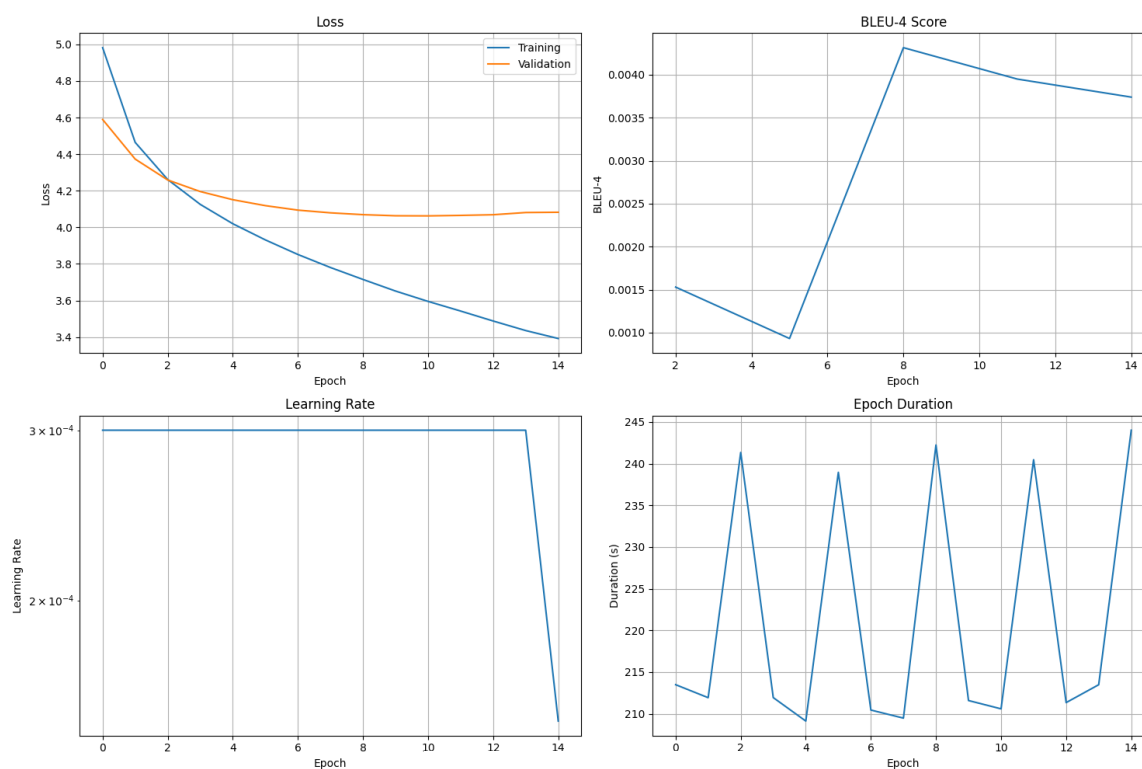
#### ۴. نتایج کمی و کیفی

برای آموزش مدل از Batch Size=32، Dropout Rate=0.5، Decoder Type=LSTM، Learning Rate=3e-4 و Number of Epoch=60 استفاده کرده ایم. اولین مدلی که به عنوان Encoder استفاده کردیم ResNet18 است، در مدل ترکیبی که با استفاده از این مدل ساخته می شود تعداد کل پارامترها قابل آموزش 5425322 است، تعداد پارامترهای Encoder برابر 131840 است و تعداد پارامترهای Decoder برابر 5161898 است. آموزش مدل پس از 16 Epoch متوقف شد، خطای آموزش نهایی 3.3019 است، خطای Validation نهایی 4.0808 است، بهترین خطای Validation ای که به دست آمد 4.0627 است و بهترین BLUE-4 Score ای که به دست آمد 0.0043 است. شکل ۱۲ نشان دهنده نمودارهای Loss، BLUE-4 Score، و Learning Rate و Epoch Duration این مدل است. بهترین Epoch آموزش Epoch 11 بود که در آن خطای Validation برابر

4.0627 بود و بهترین BLUE Epoch آموزش 9 Epoch بود که در آن BLUE Score برابر 0.0043 بود. جدول ۲ نشان دهنده یک مقایسه کامل از نتایج مدل ترکیبی با استفاده از Encoder های متفاوت است. شکل ۱۳ نشان دهنده نمودار های Loss، BLUE-، Score 4، Learning Rate و Epoch Duration از مدل با استفاده از ResNet50 Encoder و MobileNet Encoder است.

Endoer	Total Parameters	Encoder Parameters	Decoder Parameters	Final Train Loss	Final Validation Loss	Best Validation Loss	Best BLUE-4 Score	Best Epoch	Best BLUE Score Epoch
ResNet18	5425322	131840	5161898	3.3019	4.0808	4.0627	0.0043	11	9
ResNet50	5818538	525056	5161898	3.2292	4.0233	3.9971	0.0109	11	15
MobileNet V2	5621930	328448	5161898	3.2774	4.1010	4.0658	0.0030	11	10

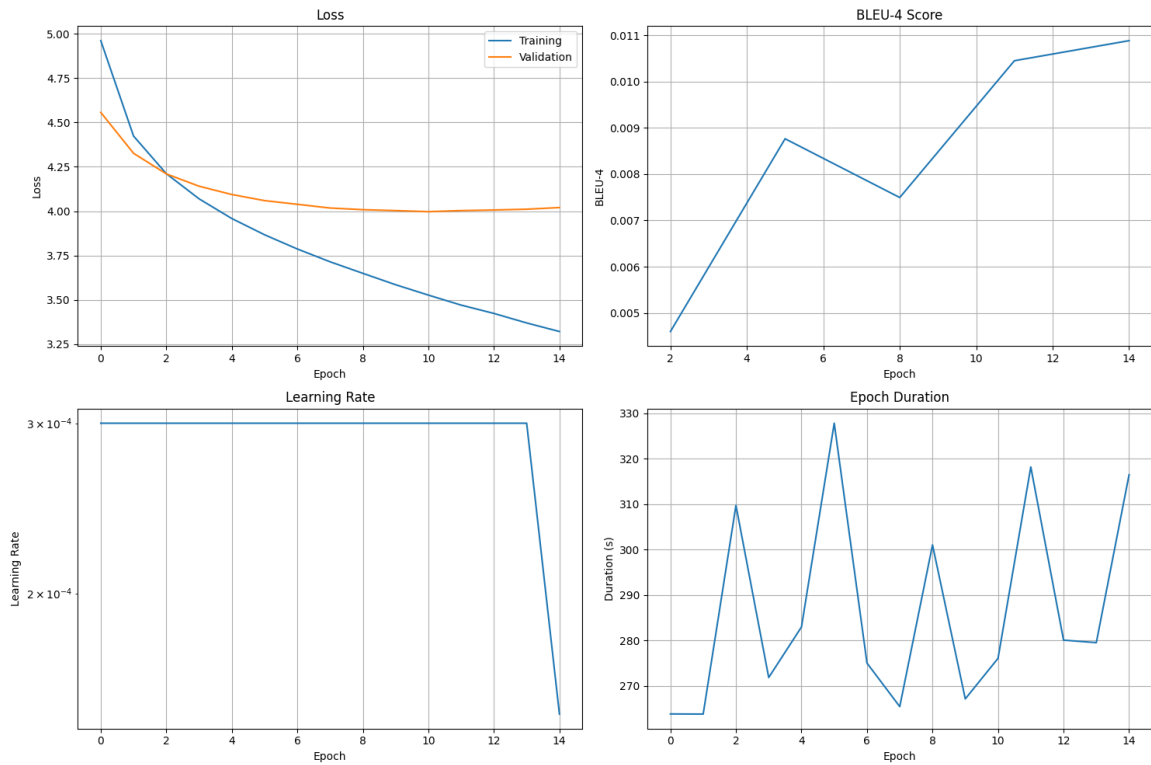
جدول ۲



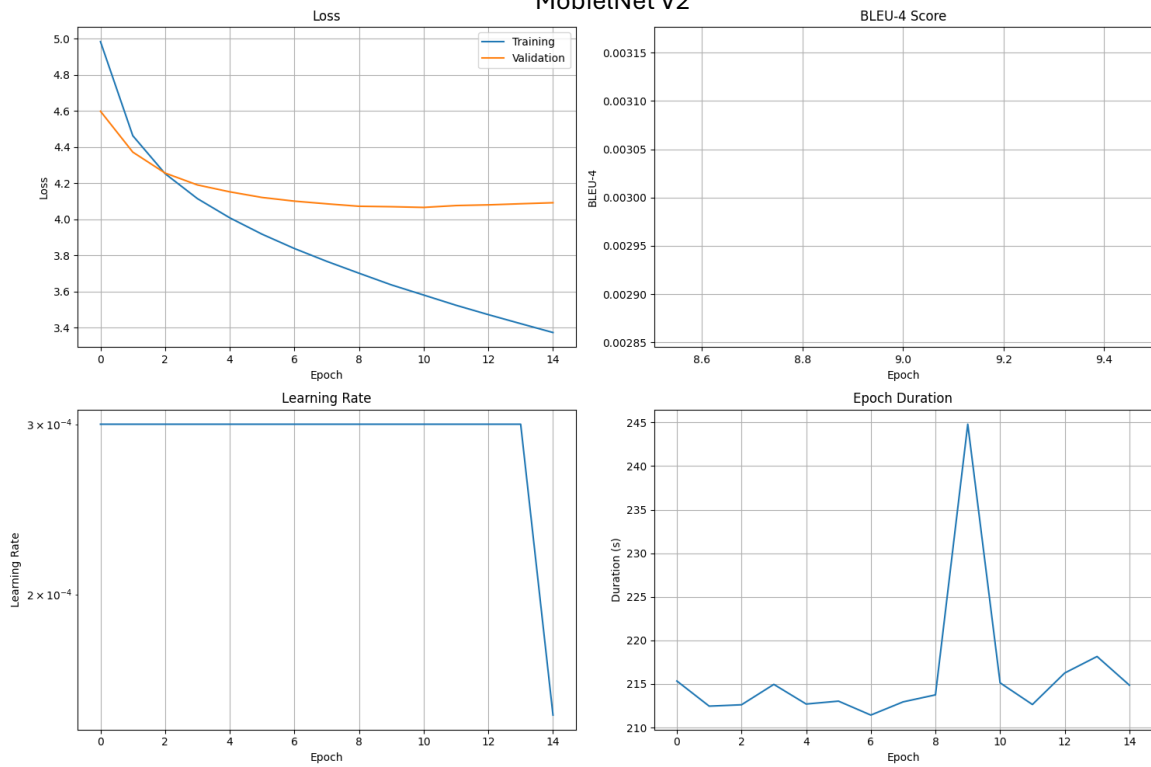
ResNet18

شکل ۱۲

## ResNet50



## MobielNet V2

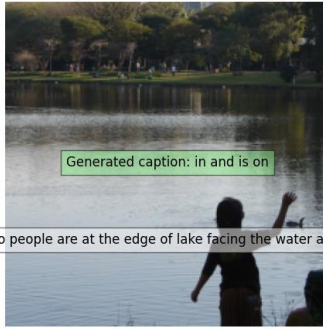


## شکل ۱۳

با بررسی نتایج می توان نتیجه گرفت که ResNet50 بهترین Encoder برای استفاده در مدل ترکیبی است. شکل ۱۴ نشان دهنده تصویر، Caption های واقعی آن تصویر و Caption های تولید شده توسط مدل است.

## ResNet18

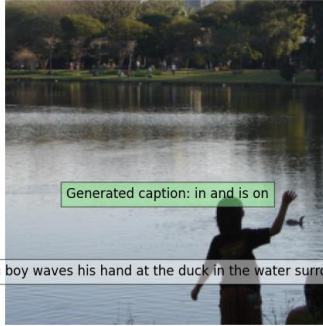
Image: 1022454332\_6af2c1449a.jpg



Generated caption: in and is on

True caption: two people are at the edge of lake facing the water and the city skyline

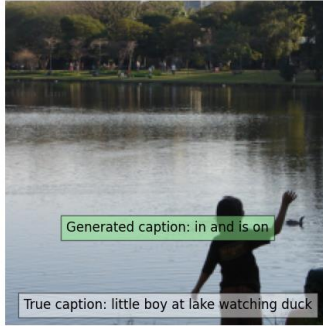
Image: 1022454332\_6af2c1449a.jpg



Generated caption: in and is on

True caption: young boy waves his hand at the duck in the water surrounded by green park

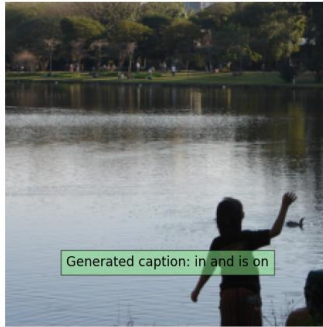
Image: 1022454332\_6af2c1449a.jpg



Generated caption: in and is on

True caption: little boy at lake watching duck

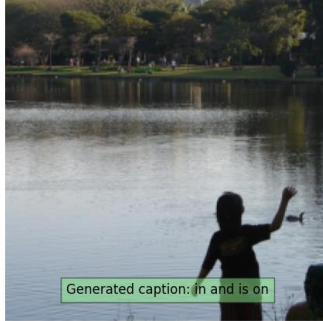
Image: 1022454332\_6af2c1449a.jpg



Generated caption: in and is on

True caption: large lake with lone duck swimming in it with several people around the edge of it

Image: 1022454332\_6af2c1449a.jpg

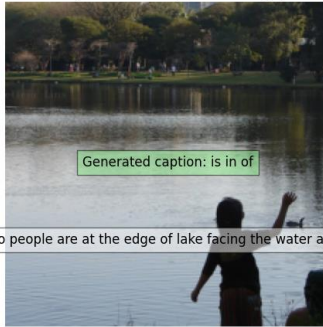


Generated caption: in and is on

True caption: child and woman are at waters edge in big city

## ResNet50

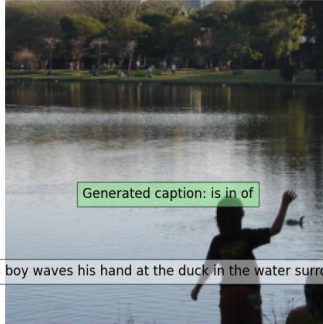
Image: 1022454332\_6af2c1449a.jpg



Generated caption: is in of

True caption: two people are at the edge of lake facing the water and the city skyline

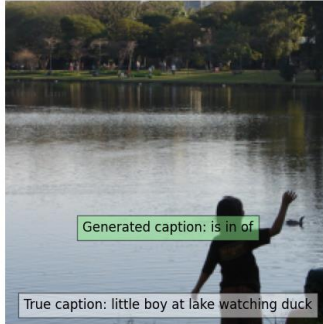
Image: 1022454332\_6af2c1449a.jpg



Generated caption: is in of

True caption: young boy waves his hand at the duck in the water surrounded by green park

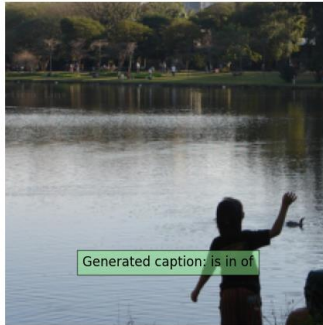
Image: 1022454332\_6af2c1449a.jpg



Generated caption: is in of

True caption: little boy at lake watching duck

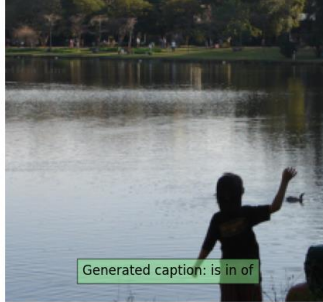
Image: 1022454332\_6af2c1449a.jpg



Generated caption: is in of

True caption: large lake with lone duck swimming in it with several people around the edge of it

Image: 1022454332\_6af2c1449a.jpg



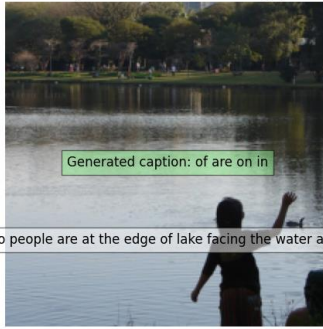
Generated caption: is in of

True caption: child and woman are at waters edge in big city



## MobielNet V2

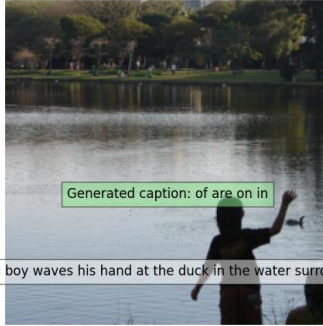
Image: 1022454332\_6af2c1449a.jpg



Generated caption: of are on in

True caption: two people are at the edge of lake facing the water and the city skyline

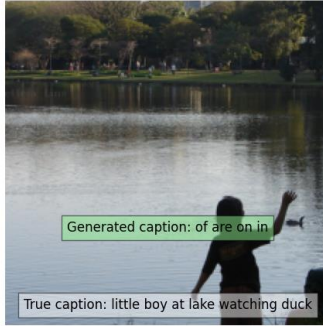
Image: 1022454332\_6af2c1449a.jpg



Generated caption: of are on in

True caption: young boy waves his hand at the duck in the water surrounded by green park

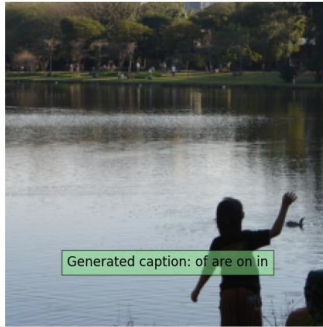
Image: 1022454332\_6af2c1449a.jpg



Generated caption: of are on in

True caption: little boy at lake watching duck

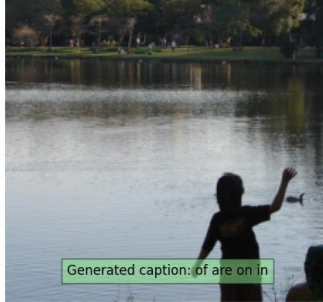
Image: 1022454332\_6af2c1449a.jpg



Generated caption: of are on in

True caption: large lake with lone duck swimming in it with several people around the edge of it

Image: 1022454332\_6af2c1449a.jpg

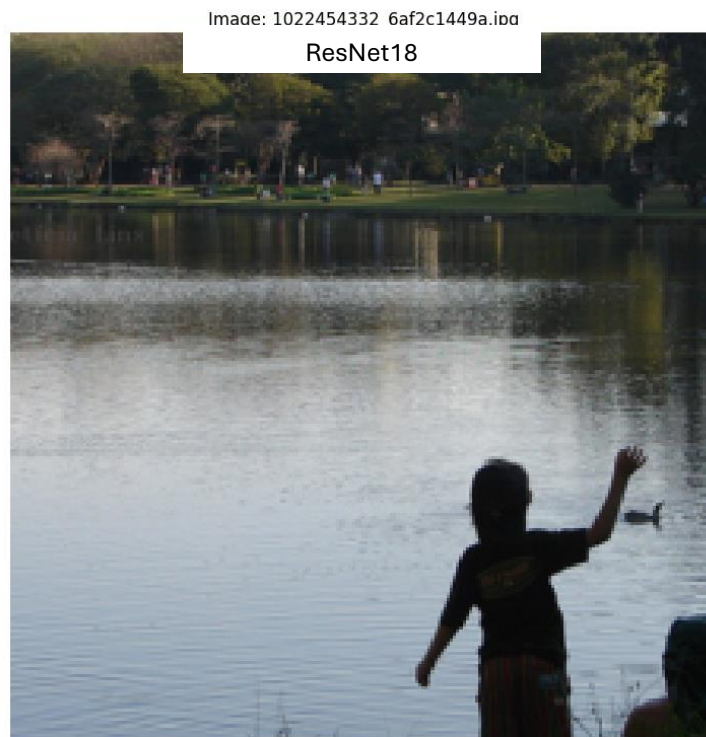


Generated caption: of are on in

True caption: child and woman are at waters edge in big city

## شکل ۱۴

شکل ۱۵ نشان دهنده مقایسه بین Greedy Decoding و Beam Search در مدل ترکیبی با سه Encoder مختلف است.

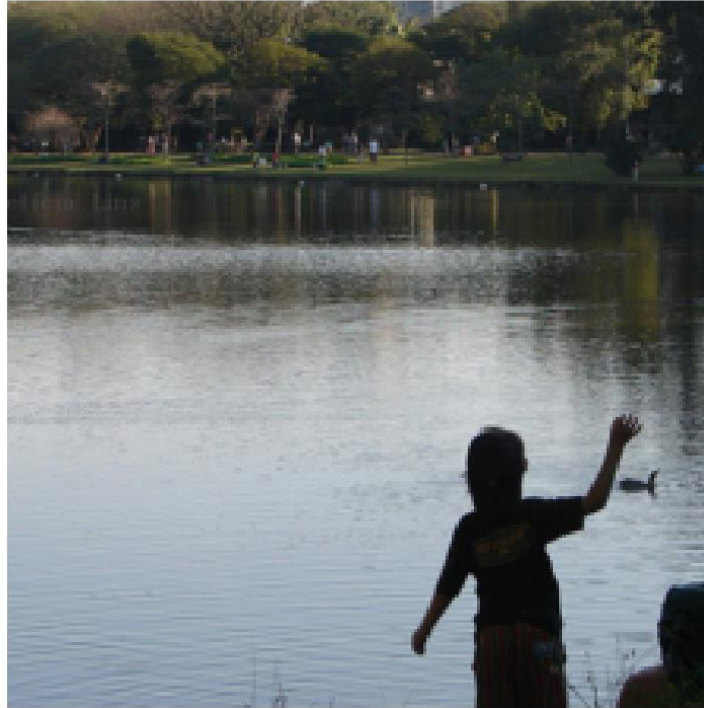


True caption: two people are at the edge of lake facing the water and the city skyline  
Greedy search: in and is on  
Beam search (k=3): in of is on  
Beam search (k=5): in is on

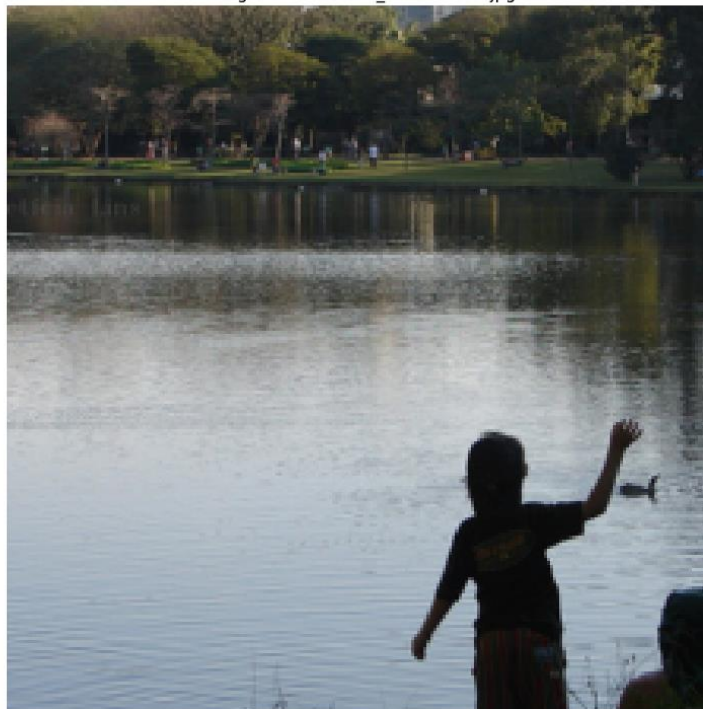


## ResNet50

Image: 1022454332\_6af2c1449a.jpg



True caption: little boy at lake watching duck  
Greedy search: is in of  
Beam search (k=3): is on  
Beam search (k=5): is on



True caption: large lake with lone duck swimming in it with several people around the edge of it  
 Greedy search: of are on in  
 Beam search (k=3): is in of and is on  
 Beam search (k=5): is in of

شکل ۱۵

## ۵. خلاصه مقاله Show and Tell: A Neural Image Caption Generator

- چکیده: این مقاله یک مدل End-to-End برای تولید خودکار Caption تصاویر معرفی می کند. این مدل از دو بخش اصلی تشکیل شده است، یک شبکه عصبی کانولوشنی برای استخراج ویژگی های تصویر و یک شبکه عصبی بازگشتی برای تولید Caption. این مدل با عنوان NIC(Neural Image Caption) قادر است با مشاهده یک تصویر بدون متن یک جمله به زبان طبیعی تولید کند که موضوع یا عملکرد اصلی تصویر را توضیح دهد.
- هدف مقاله: هدف این مقاله توسعه یک سیستم است که بتواند به صورت خودکار و بدون دخالت انسان یک Caption قابل فهم از تصاویر تولید کند. این کار شامل درک تصویر با استفاده از شبکه عصبی کانولوشنی و سپس تبدیل آن به یک توالی مناسب از کلمات با استفاده از شبکه عصبی بازگشتی است.
- معماری مدل: اولین بخش معماری این مدل Encoder است، در این بخش از یک شبکه عصبی کانولوشنی از پیش آموزش داده شده مانند Inception-V3 برای استخراج ویژگی های تصویر استفاده شده است و خروجی این شبکه به عنوان ورودی

اولیه برای شبکه عصبی بازگشتی استفاده می شود. دومین بخش معماری مدل Decoder است، در این بخش از یک شبکه عصبی بازگشتی مثل LSTM برای تولید کلمات استفاده شده است و این شبکه عصبی با استفاده از تابع هزینه  $\log$ -likelihood آموزش داده شده است. ورودی مدل یک تصویر است و خروجی آن یک دنباله از کلمات که کپشن را تشکیل می دهند.

- آموزش مدل: مدل با بیشینه سازی احتمال جمله صحیح S با توجه به تصویر I آموزش می بیند، برای آموزش از روش Maximum Likelihood Estimation و الگوریتم SGD استفاده می شود و هدف تولید یک توالی از کلمات است که احتمال شرطی  $P(S|I)$  را بیشینه کند. برای تولید جملات در مرحله تست از روش هایی مانند Beam Search و Greedy Decoding استفاده شده است. این مدل از Teacher Forcing در زمان آموزش استفاده می کند و در زمان Inferencing از Sampling یا Beam Search برای تولید کپشن استفاده می کند.
- معیار های ارزیابی: BLUE Score یکی از معیار های ارزیابی در تولید متن است. در مقایسه با مدل های قبلی NIC نتایج بسیار بهتری داشته است. BLUE-1 در Pascal VOC از 25 به 59 افزایش یافته است، BLUE-1 در Flickr30k از 56 به 66 افزایش یافته است، BLUE-1 در SBU از 19 به 28 افزایش یافته است و BLUE-4=27.7 در MSCOCO به دست آمده است. معیار های METEOR و CIDEr همبستگی بالاتری با ارزیابی انسانی دارند، نتایج NIC با این معیار ها هم بهتر از مدل های قبلی است و در بسیاری موارد توصیف های تولید شده قابل تشخیص از توصیف های انسانی نبوده اند.
- دیتاست ها: در این مقاله از دیتاست های Pascal VOC، Flickr8k، Flickr30k، SBU و MSCOCO استفاده شده است.
- ایده کلی: ایده اصلی این مقاله از Statistical Machine Translation الهام گرفته شده است، در SMT مدل یک جمله را به جمله دیگری ترجمه می کند اما در اینجا مدل یک تصویر را به یک جمله تبدیل می کند.
- نتایج: NIC یک رویکرد End-to-End است که نیازی به مراحل جداگانه مثل Detection + Template-Based-Captioning ندارد. در تمامی دیتاست ها عملکرد State-of-the-Art را دارد. کپشن های تولید شده اغلب طبیعی، قابل خواندن و دقیق هستند. مدل به خوبی از Beam Search برای تولید کپشن های بهتر استفاده می کند. نتایج نشان می دهند که افزایش اندازه دیتاست به بهبود چشم گیر عملکرد کمک می کند. در دیتاست SBU که لیبل ها خودکار و غیر دقیق هستند عملکرد مدل ضعیف تر است. در ارزیابی انسانی مدل هنوز دور از نتیجه انسان ها است.

## ۶. خلاصه مقاله Show, Attend and Tell

- هدف مقاله: این مقاله یک مدل جدید در حوزه تولید خودکار Caption تصویر با استفاده از Mechanism Attention معرفی می کند. این مدل قادر است بدون نیاز به تشخیص شیء جداگانه یک توالی از کلمات را تولید کند که Caption

معنایی تصویر را فراهم کند. هدف اصلی پژوهش افزایش دقت Caption ها با استفاده از Attention ، ایجاد قابلیت Interpretability در مدل CNN-RNN و مقایسه Soft Attention و Hard Attention در عملکرد و کارایی است.

■ معماری مدل: قسمت Encoder مدل یک شبکه عصبی کانولوشنی است که برای استخراج Feature Map های فضایی استفاده شده است. قسمت Endcoer یک LSTM برای تولید کلمات بر اساس ویژگی های تصویر و یک ماژول Attention است که در هر مرحله از RNN وزن هایی را روی Feature Map محاسبه می کند و فقط به بخش های مهم تصویر توجه می کند.

■ انواع توجه: در Soft Attention مدل با وزن دهی به تمام Feature ها توجه می کند، Soft Attention قابل آموزش با Backpropagation است و دقیق تر و قابل تفسیر تر است. در Hard Attention مدل فقط به یک قسمت تصویر توجه می کند، آموزش تصادفی و با Policy Gradient است و سریع تر ولی غیر قابل پیش بینی تر است.

■ روش تحقیق: از دیتاست های معروف Flickr8k, Flickr30k, MSCOCO برای آموزش استفاده شده است. معیار های ارزیابی BLUE-1, BLUE-4, ROUGE-L, CIDEr-D, METEOR هستند. آموزش با استفاده از Teacher Forcing, Negative Log-Likelihood Loss, Adam Optimizer صورت می گیرد.

■ نتایج: جدول زیر نشان دهنده نتایج به دست آمده است. نتایج نشان دهنده بهبود چشم گیر مدل با استفاده از Attention نسبت به مدل های بدون Attention است.

Dataset	Model	BLEU				METEOR
		BLEU-1	BLEU-2	BLEU-3	BLEU-4	
Flickr8k	Google NIC(Vinyals et al., 2014) <sup>1,2</sup>	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) <sup>o</sup>	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	<b>67</b>	44.8	29.9	19.5	18.93
	Hard-Attention	<b>67</b>	<b>45.7</b>	<b>31.4</b>	<b>21.3</b>	<b>20.30</b>
Flickr30k	Google NIC <sup>1,2</sup>	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	<b>18.49</b>
	Hard-Attention	<b>66.9</b>	<b>43.9</b>	<b>29.6</b>	<b>19.9</b>	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) <sup>a</sup>	—	—	—	—	20.41
	MS Research (Fang et al., 2014) <sup>1a</sup>	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) <sup>o</sup>	64.2	45.1	30.4	20.3	—
	Google NIC <sup>1,2</sup>	66.6	46.1	32.9	24.6	—
	Log Bilinear <sup>o</sup>	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	<b>23.90</b>
	Hard-Attention	<b>71.8</b>	<b>50.4</b>	<b>35.7</b>	<b>25.0</b>	23.04

■ نحوه کار کرد توجه: Encoder یک Feature Map فضایی از تصویر می سازد. Decoder در هر مرحله یک Context Vector را با استفاده از Attention Weights از Feature Map ها می سازد. Attention Weights با استفاده از یک MLP محاسبه می شوند. مدل به صورت End-to-End آموزش داده می شود.

■ مزایای استفاده از توجه: افزایش دقت Caption، قابلیت Interpretability، عدم نیاز به object detection، کاهش overfitting و نمایش visual-semantic alignment مزیت های توجه هستند.

■ تحلیل Caption ها: مدل می تواند Caption های دقیق تولید کند مثل A woman is holding a remote control و می تواند Caption های غلط تولید کند ولی با استفاده از Attention Map می توان دلیل خطا را متوجه شد.

- نتیجه گیری: مقاله Neural Image Caption Generation with Visual Attention نشان می دهد که با استفاده از Attention می توان مدلی ساخت که قادر به تولید Caption های دقیق تری باشد. قابلیت تفسیر مدل با استفاده از Attention Map فراهم شده است