

CS197 Harvard: AI Research Experiences

Fall 2022: Lecture 3 – “Shoulders of Giants”
Reading AI Research Papers

Instructed by Pranav Rajpurkar. Website <https://cs197.seas.harvard.edu/>

Abstract

Maybe you're trying to get into AI research, maybe you are a graduate student that is considering joining a new research lab, maybe you're in industry trying to present the latest advances on an AI problem to your colleagues. Regardless, you'll be faced with the daunting task of understanding the state of progress on the problem topic, and the gaps that are left to fill. I go through this exercise continually in my career, and a structured approach to understanding the state and gaps in something can make the task less daunting. This lecture is set up as a real walkthrough of the steps you would take to learn about a new topic in AI. My hope is that by the end of the lecture, you will have a blueprint for the kind of workflow you can use while approaching the reading of AI research papers.



DALL-E Generation: A voracious reader surrounded by papers around them, digital art

Learning outcomes:

- Conduct a literature search to identify papers relevant to a topic of interest
- Read a machine learning research paper and summarize its contributions

Approach

I'm going to break down the process of reading AI research papers into two pieces: reading *wide*, and reading *deep*. When you start learning about a new topic, you typically get more out of reading *wide*: this means navigating through literature reading small amounts of individual research papers. Our goal when reading wide is to build and improve our mental model of a research topic. Once you have identified key works that you want to understand well in the first step, you will want to read *deep*: here, you are trying to read individual papers in depth. Both reading *wide* and *deep* are necessary and complimentary, especially when you're getting started.

I am going to walk you through how you would approach reading wide and reading deep very concretely using the example of "deep learning for image captioning" as the hypothetical topic we want to break into. Let's get started.

Reading Wide

Let's start with a simple google search for "image captioning".

We get a definition of image captioning from the first result that defines the task for us. It's the second link that's interesting for us: Papers with Code

Papers with Code

Papers with Code is a community project with the mission to create a free and open resource with Machine Learning papers, code, datasets, methods and evaluation tables. I like it a lot and use it quite frequently.

The screenshot shows the 'Image Captioning' task page on the Papers with Code website. At the top, there's a navigation bar with a logo, a search bar, and links for 'Browse State-of-the-Art', 'Datasets', 'Methods', and 'More'. On the right, there are social media icons for Twitter and LinkedIn, and a 'Sign In' button. Below the navigation, a sidebar on the left indicates 'Natural Language Processing' and shows a thumbnail for 'Image Captioning'. A large title 'Image Captioning' is centered above a summary: '382 papers with code • 27 benchmarks • 49 datasets'. A detailed description follows: 'Image Captioning is the task of describing the content of an image in words. This task lies at the intersection of computer vision and natural language processing. Most image captioning systems use an encoder-decoder framework, where an input image is encoded into an intermediate representation of the information in the image, and then decoded into a descriptive text sequence. The most popular benchmarks are nocaps and COCO, and models are typically evaluated according to a BLEU or CIDEr metric.' Below this is a note: '(Image credit: Reflective Decoding Network for Image Captioning, ICCV'19)'. To the right, there's a visual comparison between a 'Basis decoder' result (a black and white photo of a clock tower) and the 'Ours' result (a color photo of a bridge with a clock tower over a river), with arrows pointing to specific parts of the image. A section titled 'Content' lists links to 'Introduction', 'Benchmarks', 'Datasets', 'Subtasks', 'Libraries', 'Papers', 'Most implemented', and 'Social'. At the bottom of the page, there are tabs for 'Trend', 'Dataset', 'Best Model', 'Paper', 'Code', and 'Compare', along with buttons for 'Add a Result', 'See all', and a download icon.

We're going to need to start making notes, so let's open up a Google Doc, and create an entry. (This is the final version of my [notes](#)).

The screenshot shows a Google Doc with a single section titled 'Image Captioning Notes'. Below the title is a link: <https://paperswithcode.com/task/image-captioning>. Underneath the link is a bulleted list of key points about the task:

- Task: describing the content of an image in words.
- Methods: encoder-decoder framework – input image is encoded into an intermediate representation of the information in the image, and then decoded into a descriptive text sequence.
- Datasets: nocaps and COCO
- Evaluation: BLEU or CIDEr metric.

Let's scroll to the benchmarks section on papers with code.

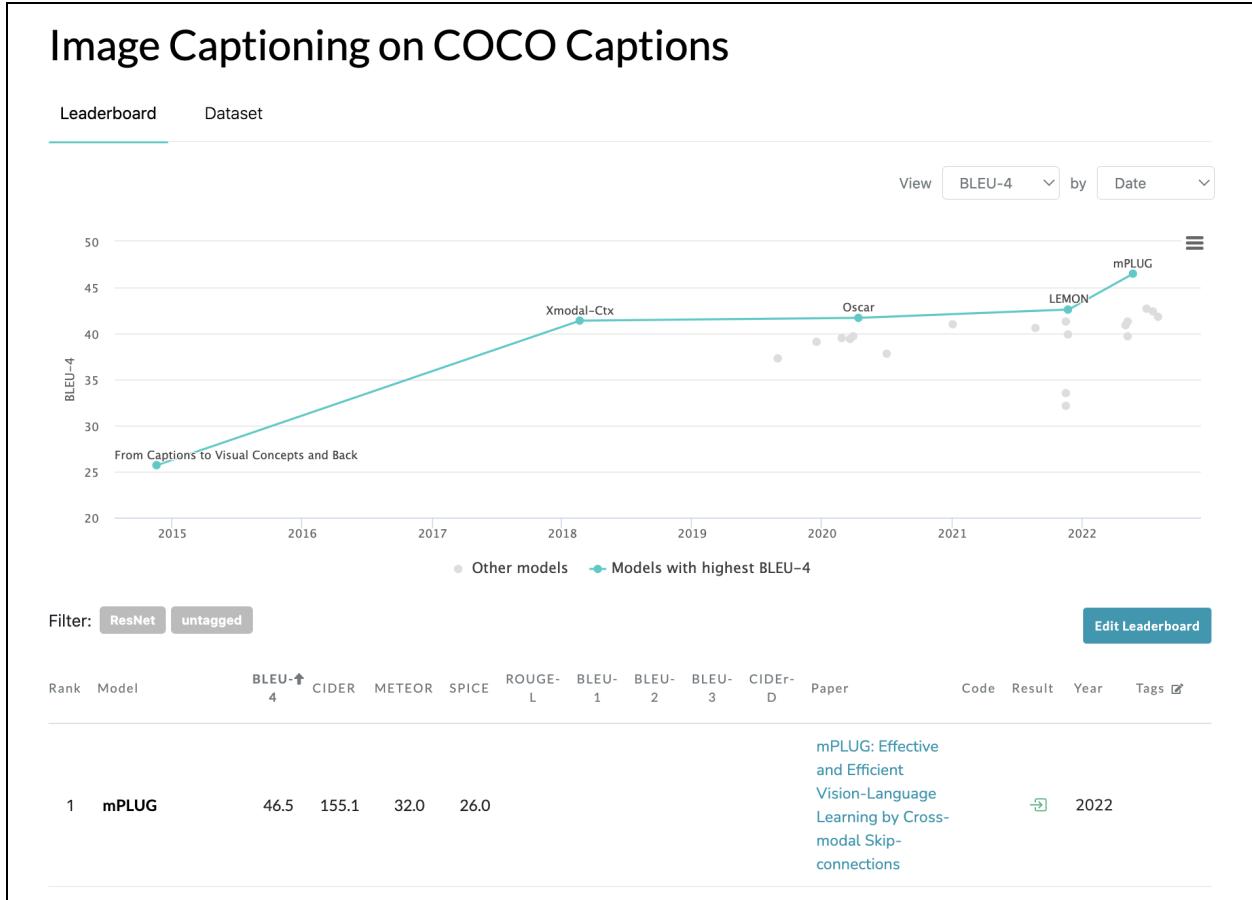
Benchmarks

[Add a Result](#)

These leaderboards are used to track progress in Image Captioning

| Trend | Dataset | Best Model | Paper | Code | Compare |
|-------|----------------------|--|-----------------------|----------------------|-------------------------|
| | COCO Captions | 🏆 mPLUG | Paper | Code | See all |
| | SCICAP | 🏆 CNN+LSTM (Vision only, First sentence) | Paper | Code | See all |
| | nocaps in-domain | 🏆 GIT, Single Model | Paper | Code | See all |
| | nocaps out-of-domain | 🏆 GIT, Single Model | Paper | Code | See all |
| | nocaps near-domain | 🏆 Microsoft Cognitive Services team | Paper | Code | See all |
| | nocaps entire | 🏆 Microsoft Cognitive Services team | Paper | Code | See all |
| | nocaps-val-in-domain | 🏆 LEMON_large | Paper | Code | See all |
| | COCO | 🏆 ExpansionNet v2 | Paper | Code | See all |

We see a few benchmarks, with the words “COCO” and “nocaps” repeating. Let’s try to look at the top benchmark first: COCO captions.



This shows us that as recently as in mid 2022, there has been a state-of-the-art (you might hear SOTA (pronounced like soda with a 't') with a method called mPLUG. The leaderboard is quite useful: it shows us metrics including BLEU-4, CIDEr, METEOR, SPICE, and similar variants. Soon, we'll need to understand what these are.

Let's click on the mPLUG paper link:

mPLUG: Effective and Efficient Vision-Language Learning by Cross-modal Skip-connections

24 May 2022 · Chenliang Li, Haiyang Xu, Junfeng Tian, Wei Wang, Ming Yan, Bin Bi, Jiabo Ye, Hehong Chen, Guohai Xu, Zheng Cao, Ji Zhang, Songfang Huang, Fei Huang, Jingren Zhou, Luo Si · [Edit social preview](#)

Large-scale pretrained foundation models have been an emerging paradigm for building artificial intelligence (AI) systems, which can be quickly adapted to a wide range of downstream tasks. This paper presents mPLUG, a new vision-language foundation model for both cross-modal understanding and generation. Most existing pre-trained models suffer from the problems of low computational efficiency and information asymmetry brought by the long visual sequence in cross-modal alignment. To address these problems, mPLUG introduces an effective and efficient vision-language architecture with novel cross-modal skip-connections, which creates inter-layer shortcuts that skip a certain number of layers for time-consuming full self-attention on the vision side. mPLUG is pre-trained end-to-end on large-scale image-text pairs with both discriminative and generative objectives. It achieves state-of-the-art results on a wide range of vision-language downstream tasks, such as image captioning, image-text retrieval, visual grounding and visual question answering. mPLUG also demonstrates strong zero-shot transferability when directly transferred to multiple video-language tasks.

 PDF

 Abstract

Let's read the abstract. We are new to image captioning, so don't expect to understand most of the sentences. I am going to make note of key information here in our Google Doc.

mPLUG: Effective and Efficient Vision-Language Learning by Cross-modal Skip-connections

<https://arxiv.org/pdf/2205.12005v2.pdf>

- mPLUG is a new vision-language foundation model for both cross-modal understanding and generation.
- Solves problem that pre-trained models suffer from the problems of low computational efficiency and information asymmetry brought by the long visual sequence in cross-modal alignment.
- Solves problem by introducing an effective and efficient vision-language architecture with novel cross-modal skip-connections.
- Achieves top performance on at least COCO.

Again, it's okay that we don't understand what this all means just yet; this is simply an attempt to extract the key idea first.

I am going to go back to Papers with Code and repeat this exercise for the winner on 2 more benchmarks: I picked nocaps-in domain and nocaps near-domain because they both had very recent submissions that had high performance.

Thus we come across the GIT method.

GIT: A Generative Image-to-text Transformer for Vision and Language

27 May 2022 · JianFeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, Lijuan Wang · [Edit social preview](#)

In this paper, we design and train a Generative Image-to-text Transformer, GIT, to unify vision-language tasks such as image/video captioning and question answering. While generative models provide a consistent network architecture between pre-training and fine-tuning, existing work typically contains complex structures (uni/multi-modal encoder/decoder) and depends on external modules such as object detectors/taggers and optical character recognition (OCR). In GIT, we simplify the architecture as one image encoder and one text decoder under a single language modeling task. We also scale up the pre-training data and the model size to boost the model performance. Without bells and whistles, our GIT establishes new state of the arts on 12 challenging benchmarks with a large margin. For instance, our model surpasses the human performance for the first time on TextCaps (138.2 vs. 125.5 in CIDEr). Furthermore, we present a new scheme of generation-based image classification and scene text recognition, achieving decent performance on standard benchmarks. Codes are released at \url{https://github.com/microsoft/GenerativeImage2Text}.

 PDF

 Abstract

You can follow my above example to make similar notes using this abstract. Note that we haven't yet opened the papers: that's intentional. Here are the notes I made:

GIT: A Generative Image-to-text Transformer for Vision and Language

<https://arxiv.org/pdf/2205.14100v4.pdf>

- Generative Image-to-text Transformer, GIT, unifies vision-language tasks such as image/video captioning and question answering.
- Solves the problem that existing work typically contains complex structures (uni/multi-modal encoder/decoder) and depends on external modules such as object detectors/taggers and optical character recognition (OCR).
- Solves problem by simplifying the architecture as one image encoder and one text decoder under a single language modeling task
- GIT establishes a new state of the arts on 12 challenging benchmarks with a large margin. including surpassing human performance for the first time on TextCaps.

I haven't used any external knowledge here; it's simply organizing key information from the abstract.

Exercise: Perform a similar summarization of the abstract for another method using the benchmarks. Look for active image captioning benchmarks where you see recent submissions, especially if the recent submissions are setting new SOTA.

Now that we've got an understanding of a couple of methods, let's try to understand the datasets a little better. We can still stick to papers with code. Let's find our way to the image captioning datasets:

The screenshot shows the 'Datasets' section of the paperswithcode.com website. The main heading is 'Datasets' with a subtitle '6,797 machine learning datasets'. A search bar and navigation links for 'Browse State-of-the-Art', 'Datasets', 'Methods', and 'More' are visible. A call-to-action button 'Share your dataset with the ML community!' is present. The main content area displays '49 dataset results for Image Captioning'. On the left, there are filters for 'Search for datasets' (with a search icon), 'Filter by Modality' (listing Images: 36, Texts: 25, 3D: 1, Audio: 1, Videos: 1), and 'Filter by Task' (with 'Image Captioning' selected and 'Question Answering' at 292). The results list includes:

- COCO (Microsoft Common Objects in Context)**: Description: The MS COCO (Microsoft Common Objects in Context) dataset is a large-scale object detection, segmentation, key-point detection, and captioning dataset. The dataset consists of... 6,318 PAPERS + 77 BENCHMARKS
- Flickr30k**: Description: The Flickr30k dataset contains 31,000 images collected from Flickr, together with 5 reference sentences provided by human annotators. 451 PAPERS + 9 BENCHMARKS
- Conceptual Captions**: Description: Automatic image captioning is the task of producing a natural-language utterance (usually a sentence) that correctly reflects the visual content of an image. Up to this point, the resour... 183 PAPERS + 1 BENCHMARK
- COCO Captions**: Description: COCO Captions contains over one and a half million captions describing over 330,000 images. For the training and validation images, five independent human generated captions ar... 105 PAPERS + 6 BENCHMARKS

We'll look into the 2 datasets we have come across before: nocaps, and COCO captions. Let's start with nocaps.

nocaps

Introduced by Agrawal et al. in [nocaps: novel object captioning at scale](#)

The nocaps benchmark consists of 166,100 human-generated captions describing 15,100 images from the OpenImages validation and test sets.

We can click the link to get to the abstract of the paper introducing the dataset.

[nocaps: novel object captioning at scale](#)

ICCV 2019 · Harsh Agrawal, Karan Desai, YuFei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, Peter Anderson · [Edit social preview](#)

Image captioning models have achieved impressive results on datasets containing limited visual concepts and large amounts of paired image-caption training data. However, if these models are to ever function in the wild, a much larger variety of visual concepts must be learned, ideally from less supervision. To encourage the development of image captioning models that can learn visual concepts from alternative data sources, such as object detection datasets, we present the first large-scale benchmark for this task. Dubbed 'nocaps', for novel object captioning at scale, our benchmark consists of 166,100 human-generated captions describing 15,100 images from the OpenImages validation and test sets. The associated training data consists of COCO image-caption pairs, plus OpenImages image-level labels and object bounding boxes. Since OpenImages contains many more classes than COCO, nearly 400 object classes seen in test images have no or very few associated training captions (hence, nocaps). We extend existing novel object captioning models to establish strong baselines for this benchmark and provide analysis to guide future work on this task.

[PDF](#)

[Abstract](#)

[ICCV 2019 PDF](#)

[ICCV 2019 Abstract](#)

As we did for the models, we can similarly make notes for the datasets.

nocaps: novel object captioning at scale

<https://arxiv.org/pdf/1812.08658v3.pdf>

- Dubbed 'nocaps', for **novel object captioning at scale**, benchmark consists of 166,100 human-generated captions describing 15,100 images from the OpenImages validation and test sets.
- Introduced to encourage the development of image captioning models that can learn visual concepts from alternative data sources, such as object detection datasets.
- The associated training data consists of COCO image-caption pairs, plus OpenImages image-level labels and object bounding boxes.
- Nearly 400 object classes seen in test images have no or very few associated training captions.
 - I came across *nocaps out-of-domain and nocaps in-domain earlier; maybe that's associated with above?*

Notice that I am leaving myself comments to come back to in my notes. Like for the methods, some words here might be new, but most words here are in typical English, so it's more easy to understand. We can do the same for the other dataset, COCO captions.

Microsoft COCO Captions: Data Collection and Evaluation Server

1 Apr 2015 · Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, C. Lawrence Zitnick · [Edit social preview](#)

In this paper we describe the Microsoft COCO Caption dataset and evaluation server. When completed, the dataset will contain over one and a half million captions describing over 330,000 images. For the training and validation images, five independent human generated captions will be provided. To ensure consistency in evaluation of automatic caption generation algorithms, an evaluation server is used. The evaluation server receives candidate captions and scores them using several popular metrics, including BLEU, METEOR, ROUGE and CIDEr. Instructions for using the evaluation server are provided.

Microsoft COCO Captions: Data Collection and Evaluation Server

<https://arxiv.org/pdf/1504.00325v2.pdf>

- Dataset (hopefully now done) contains over one and a half million captions describing over 330,000 images.
- For the training and validation images, five independent human generated captions are provided.
- Uses an evaluation server – *not sure what that means?*
- Server scores captions using several popular metrics, including BLEU, METEOR, ROUGE and CIDEr.

At this point, we have a good collection of 4 papers: 2 recent SOTA methods and 2 datasets. Looking at SOTA methods might especially be useful if we're looking to improve on the methods, but if we're trying to understand a problem domain more broadly, we need to pick up some more mature and influential works here.

Google Scholar

We're going to turn next to Google Scholar. Let's start by entering our search term: image captioning:

The screenshot shows the Google Scholar interface with the search term "image captioning" entered in the search bar. The results are sorted by relevance, showing four main entries:

- A comprehensive survey of deep learning for image captioning** [PDF] acm.org
MDZ Hossain, F Sohel, MF Shiratuddin... - ACM Computing Surveys ..., 2019 - dl.acm.org
... In this section, we review and describe the main categories of existing **image captioning** methods, including template-based **image captioning**, retrieval-based **image captioning**, and ...
☆ Save 99 Cite Cited by 455 Related articles All 7 versions
- Show and tell: Lessons learned from the 2015 mscoco image captioning challenge** [PDF] ieee.org
O Vinyals, A Toshev, S Bengio... - IEEE transactions on ..., 2016 - ieexplore.ieee.org
... difference of these metrics is on how humans rank on it versus several automatic **image captioning** systems (such as the one we propose). Interestingly, BLEU score seems to be quite ...
☆ Save 99 Cite Cited by 810 Related articles All 20 versions
- Boosting image captioning with attributes** [PDF] thecvf.com
T Yao, Y Pan, Y Li, Z Qiu, T Mei - Proceedings of the IEEE ..., 2017 - openaccess.thecvf.com
... /video captioning [20, 21, 31], we aim to formulate our **image captioning** models in an end-to-end fashion based on RNNs which encode the given **image** and/or its detected attributes ...
☆ Save 99 Cite Cited by 620 Related articles All 12 versions
- Image captioning with semantic attention** [PDF] thecvf.com
Q You, H Jin, Z Wang, C Fang... - Proceedings of the IEEE ..., 2016 - openaccess.thecvf.com
... and objects in an **image**. In this paper, we propose a new **image captioning** approach that ... Our definition for semantic attention in **image captioning** is the ability to provide a detailed, ...
☆ Save 99 Cite Cited by 1591 Related articles All 11 versions

Google scholar sorts by relevance and includes a few useful details, including the number of citations that the paper has received. We're in luck – we have a survey paper at the top of our search results. A survey paper typically reviews and describes the state of a problem space, and often includes challenges and opportunities. Reviews/Surveys may not always be up-to-date, comprehensive, or completely accurate, but especially if we're new to a space, can get us up to speed.

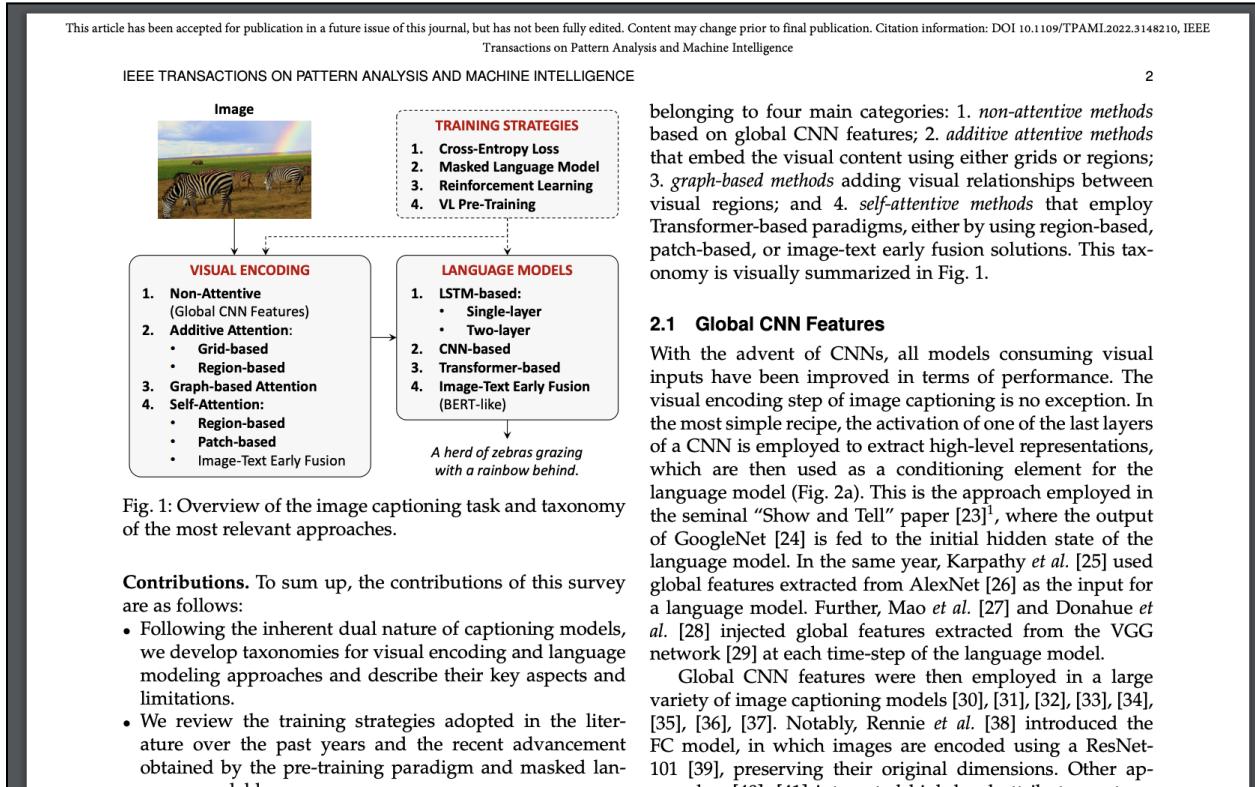
We might suspect that since we've seen some fairly recent papers achieve SOTA (2022), we might not want to look at a 2019 survey. Let's see if we can find a more recent one by explicitly searching for a survey and using the left timeline sidebar to filter results at least as recent as 2021.

The screenshot shows a Google Scholar search results page. The search query is "image captioning deep learning survey". The results are filtered to show "Articles" and there are approximately 15,500 results. The first result is a survey paper by Stefaniini, M. Cornia, L. Baraldi, et al., titled "From show to tell: a survey on deep learning-based image captioning". It was published in Pattern Analysis and Machine Intelligence in 2022. The second result is a paper by Oluwasammi, MU Aftab, Z Qin, ST Ngoo, TV Doan, et al., titled "[HTML] Features to text: a comprehensive survey on semantic segmentation and image captioning", published in Complexity in 2021. The third result is a paper by Zohouri-Shahzadi, JK Kalita, et al., titled "Neural attention for image captioning: review of outstanding methods", published in Artificial Intelligence Review in 2021. The fourth result is another version of the survey paper by Stefaniini et al., titled "From show to tell: A survey on image captioning", published in arXiv preprint arXiv in 2021.

Let's hit the first result.

| Abstract | Abstract: |
|-----------|--|
| Authors | Connecting Vision and Language plays an essential role in Generative Intelligence. For this reason, large research efforts have been devoted to image captioning, i.e. describing images with syntactically and semantically meaningful sentences. Starting from 2015 the task has generally been addressed with pipelines composed of a visual encoder and a language model for text generation. |
| Citations | |
| Keywords | |
| Metrics | During these years, both components have evolved considerably through the exploitation of object regions, attributes, the introduction of multi-modal connections, fully-attentive approaches, and BERT-like early-fusion strategies. However, regardless of the impressive results, research in image captioning has not reached a conclusive answer yet. This work aims at providing a comprehensive overview of image captioning approaches, from visual encoding and text generation to training strategies, datasets, and evaluation metrics. In this respect, we quantitatively compare many relevant state-of-the-art approaches to identify the most impactful technical innovations in architectures and training strategies. Moreover, many variants of the problem and its open challenges are discussed. The final goal of this work is to serve as a tool for understanding the existing literature and highlighting the future directions for a research area where Computer Vision and Natural Language Processing can find an optimal synergy. |
| Media | |

You might find that for survey papers, the abstract does not typically provide the same level of specificity as a typical research article. Survey papers, however, are typically more accessible (at least in some parts) because they include more background about a topic. Let's open this one up.



How should we read through a 15-page 2 column review article? In this stage, where we're reading *wide*, we will be very selective in what we read. For this paper, I read:

- Figure 1 on page 2 of the review. This typically gives a good visual organization of the key point of the review. Organization of section headings in this paper
- The "Contributions" on the second page. I didn't find this super useful in this paper, since it didn't make clear the takeaways. On other papers, I would try to find the takeaways.
- The "Conclusions and Future Directions" on the last page.

These rather short pieces of the paper should be sufficient for now. Take some time on your own to read through these sections and see if you can come up with 8-10 bullet points of notes.

Here are the notes I made using the above sections.

From Show to Tell: A Survey on Deep Learning-based Image Captioning

<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9706348>

- Divided into visual encoding, training strategies, language models. Attention shows up a lot in visual encoding.
- Variety of visual encoding strategies (pre-training from scratch, using detections, using features from multi-modal models) all appear to be on par in terms of performance.
- Training strategies have a recent advancement obtained by the pre-training paradigm and masked language model losses.
 - Growing size of pre-training models is a concern for equality in the community.
 - Growing dichotomy between early-fusion strategies and the encoder-decoder paradigm is an open issue.
- Specializing in particular domains and generating captions with different styles and aims is still among the main open challenges for image captioning.
 - variants such as novel objects captioning or controllable captioning, or subword-based tokenization techniques might handle issues.
- Fairness and bias need for designing specific evaluation metrics and focusing on the robustness to unwanted correlations.
- Development of scores that do not need reference captions for assessing the performance would be key for a shift towards unsupervised image captioning.

At this point in time, we already have 2 pages worth of notes!

Image Captioning Notes

<https://paperswithcode.com/task/image-captioning>

- Task: describing the content of an image in words.
- Methods: encoder-decoder framework – input image is encoded into an intermediate representation of the information in the image, and then decoded into a descriptive text sequence.
- Datasets: nocaps and COCO
- Evaluation: BLEU or CIDEr metric.

mPLUG: Effective and Efficient Vision-Language Learning by Cross-modal Skip-connections

<https://arxiv.org/pdf/2205.12005v2.pdf>

- mPLUG is a new vision-language foundation model for both cross-modal understanding and generation.
- Solves problem that pre-trained models suffer from the problems of low computational efficiency and information asymmetry brought by the long visual sequence in cross-modal alignment.
- Solves problem by introducing an effective and efficient vision-language architecture with novel cross-modal skip-connections.
- Achieves top performance on at least COCO.

GIT: A Generative Image-to-text Transformer for Vision and Language

<https://arxiv.org/pdf/2205.14100v4.pdf>

- Generative Image-to-text Transformer, GIT, unifies vision-language tasks such as image/video captioning and question answering.
- Solves the problem that existing work typically contains complex structures (uni/multi-modal encoder/decoder) and depends on external modules such as object detectors/taggers and optical character recognition (OCR).
- Solves problem by simplifying the architecture as one image encoder and one text decoder under a single language modeling task
- GIT establishes a new state of the arts on 12 challenging benchmarks with a large margin. including surpassing human performance for the first time on TextCaps.

nocaps: novel object captioning at scale

<https://arxiv.org/pdf/1812.08658v3.pdf>

- Dubbed 'nocaps', for novel object **captioning at scale**, benchmark consists of 166,100 human-generated captions describing 15,100 images from the OpenImages validation and test sets.
- Introduced to encourage the development of image captioning models that can learn visual concepts from alternative data sources, such as object detection datasets.
- The associated training data consists of COCO image-caption pairs, plus OpenImages image-level labels and object bounding boxes.
- Nearly 400 object classes seen in test images have no or very few associated training captions.
- I came across *nocaps out-of-domain* and *nocaps in-domain* earlier; maybe that's associated with above?

Microsoft COCO Captions: Data Collection and Evaluation Server

<https://arxiv.org/pdf/1504.00325v2.pdf>

- Dataset (hopefully now done) contains over one and a half million captions describing over 330,000 images.
- For the training and validation images, five independent human generated captions are provided.
- Uses an evaluation server – not sure what that means?
- Server scores captions using several popular metrics, including BLEU, METEOR, ROUGE and CIDEr.

From Show to Tell: A Survey on Deep Learning-based Image Captioning

<https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=9706348>

- Divided into visual encoding, training strategies, language models. Attention shows up a lot in visual encoding.
- Variety of visual encoding strategies (pre-training from scratch, using detections, using features from multi-modal models) all appear to be on par in terms of performance.
- Training strategies have a recent advancement obtained by the pre-training paradigm and masked language model losses.
 - Growing size of pre-training models is a concern for equality in the community.
 - Growing dichotomy between early-fusion strategies and the encoder-decoder paradigm is an open issue.
- Specializing in particular domains and generating captions with different styles and aims is still among the main open challenges for image captioning.
 - variants such as novel objects captioning or controllable captioning, or subword-based tokenization techniques might handle issues.
- Fairness and bias need for designing specific evaluation metrics and focusing on the robustness to unwanted correlations.
- Development of scores that do not need reference captions for assessing the performance would be key for a shift towards unsupervised image captioning.

That's great. Your notes, like mine, probably contain many terms you are encountering for the first time... We've seen terms like "encoder-decoder architecture," "language modeling task,"

“cross-modal skip connections,” or “subword-based tokenization techniques” that might be beyond our reach of current understanding, but that’s okay because when we’re reading *wide*.

At this point, we can still put together a fairly neat summary of what we’ve learnt! Before you read mine, try to compile your learnings – from reading the 1 overview page, 2 methods paper abstracts, 2 datasets paper abstracts, and small sections from 1 survey paper – all into a summary paragraph.

Learnings

- The task of image captioning is to describe the content of an image in words. Popular datasets to tackle this question include nocaps and coco, which each contain hundreds of thousands of images. Recent SOTA methods for image captioning include mPLUG and GIT. mPLUG proposes an effective and efficient vision-language architecture and GIT simplifies typically complex structures into one image encoder and one text decoder under a single language modeling task. Specializing image captioning for particular domains and generating captions with different styles and aims is possibly a big open challenge.

This is great!

Between Wide and Deep

You may have noticed that your notes are likely already diverging from mine, especially towards the last paper! Even if we have read the same content, you are starting to develop your own mental model of the problem space, and your interest will be piqued by a different set of words. That's good – you're developing a taste for the kind of research questions you will find interesting.

At this point, I find myself particularly intrigued by the SOTA methods: Why are they achieving high performance? According to the review paper, it looks like "training strategies using pre-training" have been an advance. Maybe that's worth keeping an eye out for!

Related Work

At this stage, I'll find it effective to read the related works sections of these papers: they often make it clear how researchers in the field have traditionally approached the problems and what the emerging trends are. It's important to pick papers that are recently published.

Let's dive in: both mPLUG (<https://arxiv.org/pdf/2205.12005v2.pdf>) and GIT (<https://arxiv.org/pdf/2205.14100v4.pdf>) are recently published methods that achieve SOTA.

The mPLUG related work starts towards the end of page 2. Let's read the related works subsection (I've only included one paragraph). I will focus particularly on the vision-language pre-training subsection.

two separate Transformer networks.

In this work, we propose mPLUG, a unified Multi-modal Pre-training framework for both vision-Language Understanding and Generation. mPLUG performs effective and efficient vision-language learning with novel cross-modal skip-connections to address the fundamental information asymmetry problem. Instead of fusing visual and linguistic representations at the same levels, the cross-modal skip-connections enables the fusion to occur at disparate levels in the abstraction hierarchy across the modalities. It creates inter-layer shortcuts that skip a certain number of layers for visual representations to reflect the semantic richness of language compared to vision. As shown

2 Related Work

2.1 Vision-Language Pre-training

Vision-Language pre-training (VLP) has recently received tremendous success and achieved state-of-the-art results across a variety of vision-language tasks [14, 15, 16]. In terms of how information from different modalities are aggregated, typical approaches to VLP [1, 2, 3, 5, 6, 17, 18] can be roughly divided into two categories: *dual encoder* and *fusion encoder*. Dual encoder approach utilizes two single-modal encoders to encode images and text separately, and then uses simple functions such as dot product to model the instance-level cross-modal interaction between image and text. The

Now, attempt to update your previous notes of mPLUG with your understanding of this related work section. Here are my own notes, unfiltered: I've used shorthand and left in spelling/grammar errors :)

mPLUG: Effective and Efficient Vision-Language Learning by Cross-modal Skip-connections

<https://arxiv.org/pdf/2205.12005v2.pdf>

- mPLUG is a new vision-language foundation model for both cross-modal understanding and generation.
- Solves problem that pre-trained models suffer from the problems of low computational efficiency and information asymmetry brought by the long visual sequence in cross-modal alignment.
- Solves problem by introducing an effective and efficient vision-language architecture with novel cross-modal skip-connections.
- Achieves top performance on at least COCO.
- Related work:
 - For vision language pre training, there are two approaches to how info from different modalities is aggregate: dual encoder, and fusion.
 - Dual encoder uses 2 encoders, one vision, one text, and simple functions to model interaction between image and text.
 - Plus: computationally efficient.
 - Minus: can't handle complicated vision-language tasks.
 - Examples: CLIP, ALIGN
 - Fusion encoders uses attention to model interaction between image and text.
 - Plus: better capture association between image and text for vision-language understanding tasks.
 - Minus: inference is slow, with exceptions e.x. PixelBERT.
 - Examples: Single-stream: UNITER, two-stream: LXMERT.
 - mPLUG enables the fusion to occur at disparate levels in the abstraction hierarchy across the modalities.

The examples we have collected will serve to populate our reading list when we start to read through individual papers.

Exercise: Repeat this process, this time for the GIT paper.

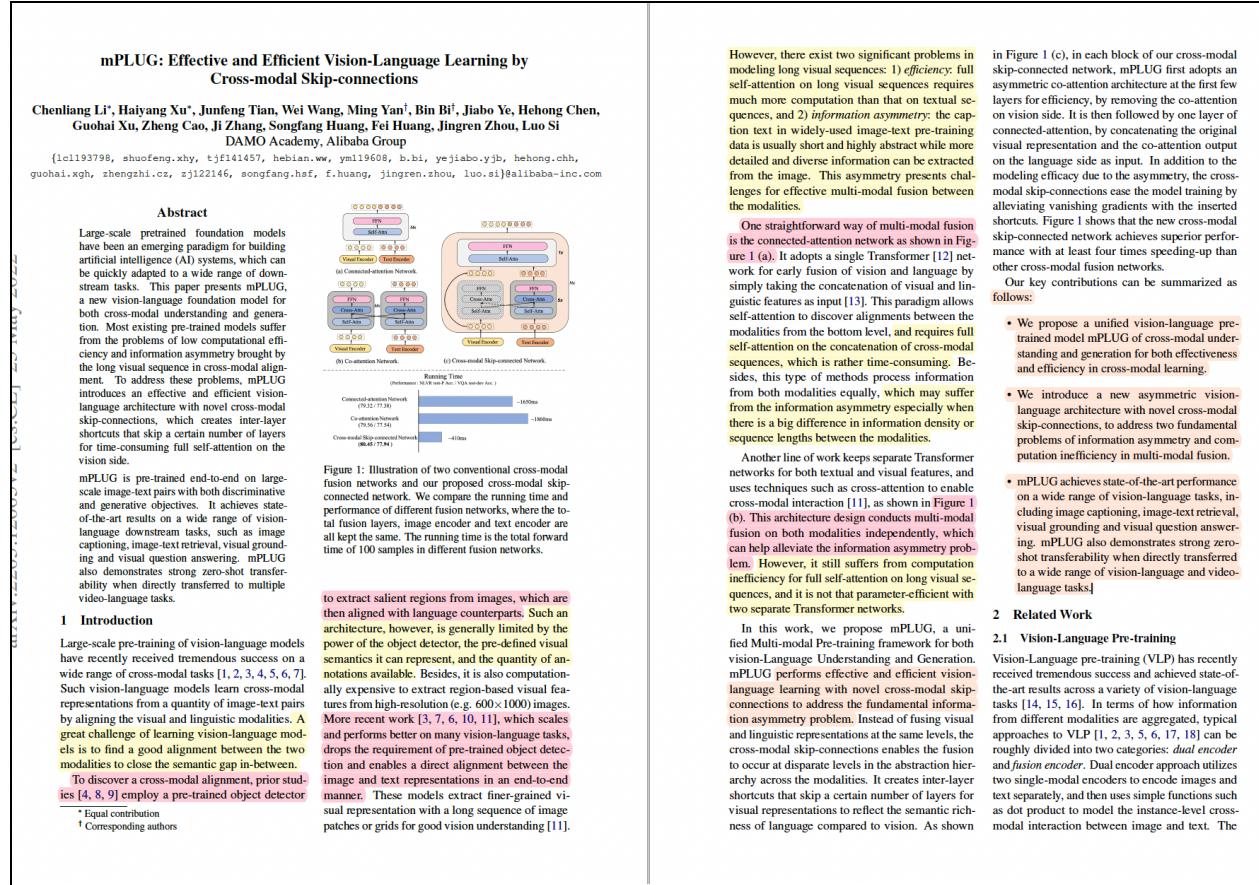
Reading Deep

We've identified a few key works for our topic in the first stage and gotten a mental model of the space of image captioning. We're now going to compliment wide reading with deep reading, covering the way to read an individual paper. Most papers are written for an audience that shares a common foundation: that's what allows for the papers to be relatively concise. Building that foundation takes time, in the span of months, if not years. Thus reading a first paper on a topic can easily take over 10 hours (some papers have definitely taken me 20 or 30 hours) and leave one feeling overwhelmed.

So I would like you to take an incremental approach here. Understand that, in your first pass, you will not understand more than 10% of the research paper. The paper may require us to read another more fundamental paper (which might require reading a third paper and so on; it could be [turtles all the way down](#))! Then, in your second pass, you might understand 20% of the paper. Understanding 100% of a paper might require a significant leap, maybe because it's poorly written, insufficiently detailed, or simply too technically/mathematically advanced. We thus want to aim to build up to understanding as much of the paper as possible – I'll bet that 70-80% of the paper is a good target.

We will go through the mPLUG paper, and I'll walk you through my first read of it. It's a good idea to be able to highlight papers as you go through them: or make comments. You can use Adobe or Preview (on Mac) to highlight PDFs (or a web-based solution like <https://hypothes.is/> for the annotation).

I'm going to read the Introduction section first. In a later lecture, I will share with you how to write good introductions. Introductions are a good way to start a paper because they are typically written for a more general audience than the rest of the paper.



mPLUG: Effective and Efficient Vision-Language Learning by Cross-modal Skip-connections

Chenliang Li¹, Haiyang Xu¹, Junfeng Tian, Wei Wang, Ming Yan¹, Bin Bi¹, Jiaob Ye, Hehong Chen, Guohai Xu, Zheng Cao, Ji Zhang, Songfang Huang, Fei Huang, Jingren Zhou, Luo Si
DAMO Academy, Alibaba Group
{lc1193798, shufeng.xhy, tji141457, hebian.ww, yml19608, b.bi, yejiaobo.yjb, hehong.chh, guohai.xgh, zhengzhi.cz, zjl22146, songfang.hsf, f.huang, jingren.zhou, luo.si}@alibaba-inc.com

Abstract

Large-scale pretrained foundation models have been an emerging paradigm for building artificial intelligence (AI) systems, which can be quickly adapted to a wide range of downstream tasks. This paper presents mPLUG, a new vision-language foundation model for both cross-modal understanding and generation. Most existing pre-training models suffer from the problem of low computational efficiency and information asymmetry brought by the long visual sequence in cross-modal alignment. To address these problems, mPLUG introduces an effective and efficient vision-language architecture with novel cross-modal skip-connections, which creates inter-layer shortcuts that skip a certain number of layers for time-consuming full self-attention on the visual side.

mPLUG is pre-trained end-to-end on large-scale image-text pairs with both discriminative and generative objectives. It achieves state-of-the-art results on a wide range of vision-language downstream tasks, such as image captioning, image-text retrieval, visual grounding and visual question answering. mPLUG also demonstrates strong zero-shot transferability when directly transferred to multiple video-language tasks.

1 Introduction

Large-scale pre-training of vision-language models have recently received tremendous success on a wide range of cross-modal tasks [1, 2, 3, 4, 5, 6, 7]. Such vision-language models learn cross-modal representations from a quantity of image-text pairs by aligning the visual and linguistic modalities. A great challenge of learning vision-language models is to find a good alignment between the two modalities to close the semantic gap in-between.

To discover a cross-modal alignment, prior studies [4, 8, 9] employ a pre-trained object detector

Figure 1: Illustration of two conventional cross-modal fusion networks and our proposed cross-modal skip-connected network. We compare the running time and performance of different fusion networks, where the total fusion layers, image encoder and text encoder are all kept the same. The running time is the total forward time of 100 samples in different fusion networks.

to extract salient regions from images, which are then aligned with language counterparts. Such an architecture, however, is generally limited by the power of the object detector, the pre-defined visual semantics it can represent, and the quantity of annotations available. Besides, it is also computationally expensive to extract region-based visual features from high-resolution (e.g. 600×1000) images. More recent work [3, 7, 6, 10, 11], which scales and performs better on many vision-language tasks, drops the requirement of pre-trained object detection and enables a direct alignment between the image and text representations in an end-to-end manner. These models extract finer-grained visual representation with a long sequence of image patches or grids for good vision understanding [11].

However, there exist two significant problems in modeling long visual sequences: 1) *efficiency*: full self-attention on long visual sequences requires much more computation than that on textual sequences, and 2) *information asymmetry*: the caption text is widely-used image-text pre-training data is usually short and highly abstract while more detailed and diverse information can be extracted from the image. This asymmetry presents challenges for effective multi-modal fusion between the modalities.

One straightforward way of multi-modal fusion is the connected-attention network as shown in Figure 1 (a). It adopts a single Transformer [12] network for early fusion of vision and language by simply taking the concatenation of visual and linguistic features as input [13]. This paradigm allows self-attention to discover alignments between the modalities from the bottom level, and requires full self-attention on the concatenation of cross-modal sequences, which is rather time-consuming. Besides, this type of methods process information from both modalities equally, which may suffer from the information asymmetry especially when there is a big difference in information density or sequence lengths between the modalities.

Another line of work keeps separate Transformer networks for both textual and visual features, and uses techniques such as cross-attention to enable cross-modal interaction [11], as shown in Figure 1 (b). This architecture design conducts multi-modal fusion on both modalities independently, which can help alleviate the information asymmetry problem. However, it still suffers from computation inefficiency for full self-attention on long visual sequences, and it is not parameter-efficient with two separate Transformer networks.

In this work, we propose mPLUG, a unified Multi-modal Pre-training framework for both vision-Language Understanding and Generation. mPLUG performs effective and efficient vision-language learning with novel cross-modal skip-connections to address the fundamental information asymmetry problem. Instead of fusing visual and linguistic representations at the same levels, the cross-modal skip-connections enables the fusion to occur at disparate levels in the abstraction hierarchy across the modalities. It creates inter-layer shortcuts that skip a certain number of layers for visual representations to reflect the semantic richness of language compared to vision. As shown

in Figure 1 (c), in each block of our cross-modal skip-connected network, mPLUG first adopts an asymmetric co-attention architecture at the first few layers for efficiency, by removing the co-attention on vision side. It is then followed by one layer of connected-attention, by concatenating the original visual representation and the co-attention output the language side as input. In addition to the modeling efficacy due to the asymmetry, the cross-modal skip-connections ease the model training by alleviating vanishing gradients with the inserted shortcuts. Figure 1 shows that the new cross-modal skip-connected network achieves superior performance with at least four times speeding-up than other cross-modal fusion networks.

Our key contributions can be summarized as follows:

- We propose a unified vision-language pre-trained model mPLUG of cross-modal understanding and generation for both effectiveness and efficiency in cross-modal learning.
- We introduce a new asymmetric vision-language architecture with novel cross-modal skip-connections, to address two fundamental problems of information asymmetry and computation inefficiency in multi-modal fusion.
- mPLUG achieves state-of-the-art performance on a wide range of vision-language tasks, including image captioning, image-text retrieval, visual grounding and visual question answering. mPLUG also demonstrates strong zero-shot transferability when directly transferred to a wide range of vision-language and video-language tasks.]

2 Related Work

2.1 Vision-Language Pre-training

Vision-Language pre-training (VLP) has recently received tremendous success and achieved state-of-the-art results across a variety of vision-language tasks [14, 15, 16]. In terms of how information from different modalities are aggregated, typical approaches to VLP [1, 2, 3, 5, 6, 17, 18] can be roughly divided into two categories: *dual encoder* and *fusion encoder*. Dual encoder approach utilizes two single-modal encoders to encode images and text separately, and then uses simple functions such as dot product to model the instance-level cross-modal interaction between image and text. The

I have highlighted what I found to be important parts of the introduction. The yellow highlights are the problems/challenges, the pink highlights are the solutions to the challenges, and the orange highlights are the main contributions of the work we're reading.

Notice the alternating yellow/pink highlights. The paper is introducing a general problem, talking about a solution to that problem, then a problem with the solution, and another solution to that problem. Four levels deep, the paper specifies the problem it solves.

Notice then that the contribution of the paper is a specific solution for a specific problem of a more general solution for a more general problem of an even more general solution to an even more general problem etc. This is typical. We can summarize our understanding of this problem-solution chain:

- Introduction:
 - Problem 1: How to find alignment between image and text modalities?
 - Solution 1: Pre-trained object detectors to find salient regions from images.
 - Problem 2: Limited by power of object detector & available annotations.
 - Solution 2: Direct alignment without object detectors
 - Problem 2a: Efficiency because of lot of computation of self-attention on visual sequences
 - Problem 2b: information asymmetry because text is short compared to info in image.
 - Solution 2a: connected attention network, using single transformer for early fusion. Has problems 2a and 2b.
 - Solution 2b: cross attention network, does fusion on both modalities independently. No longer has problem 2b, but still has 2a.
 - Solution 3 (proposed solution): cross-modal skip connections. Solves problem 2a and 2b.

This is neat – we've built a mental model of how the different pieces fit in. A well written introduction should allow you to extract such a problem–solution chain, but not every paper will make this explicit. We can often refer to the figure 1 to understand the main idea of the paper better:

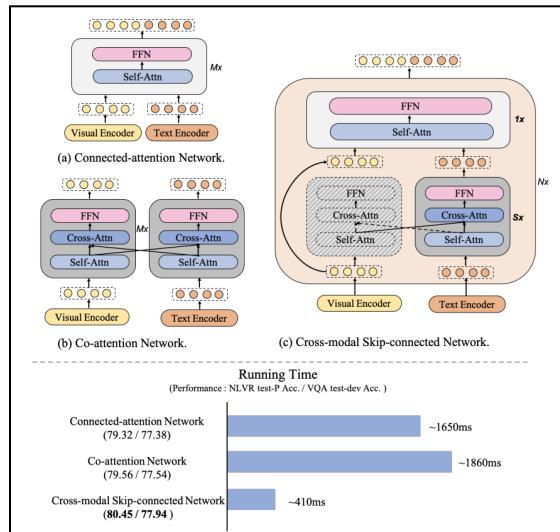


Figure 1: Illustration of two conventional cross-modal fusion networks and our proposed cross-modal skip-connected network. We compare the running time and performance of different fusion networks, where the total fusion layers, image encoder and text encoder are all kept the same. The running time is the total forward time of 100 samples in different fusion networks.

We can see the proposed solution (c), and the comparison to solutions in (a) and (b), which correspond to solutions 2a and 2b in our notes.

What's next? We've already read the abstract previously, we've read the introduction, we've seen Figure 1. On this paper, we've also had the opportunity to read the related work.

Now we read the methods, right?

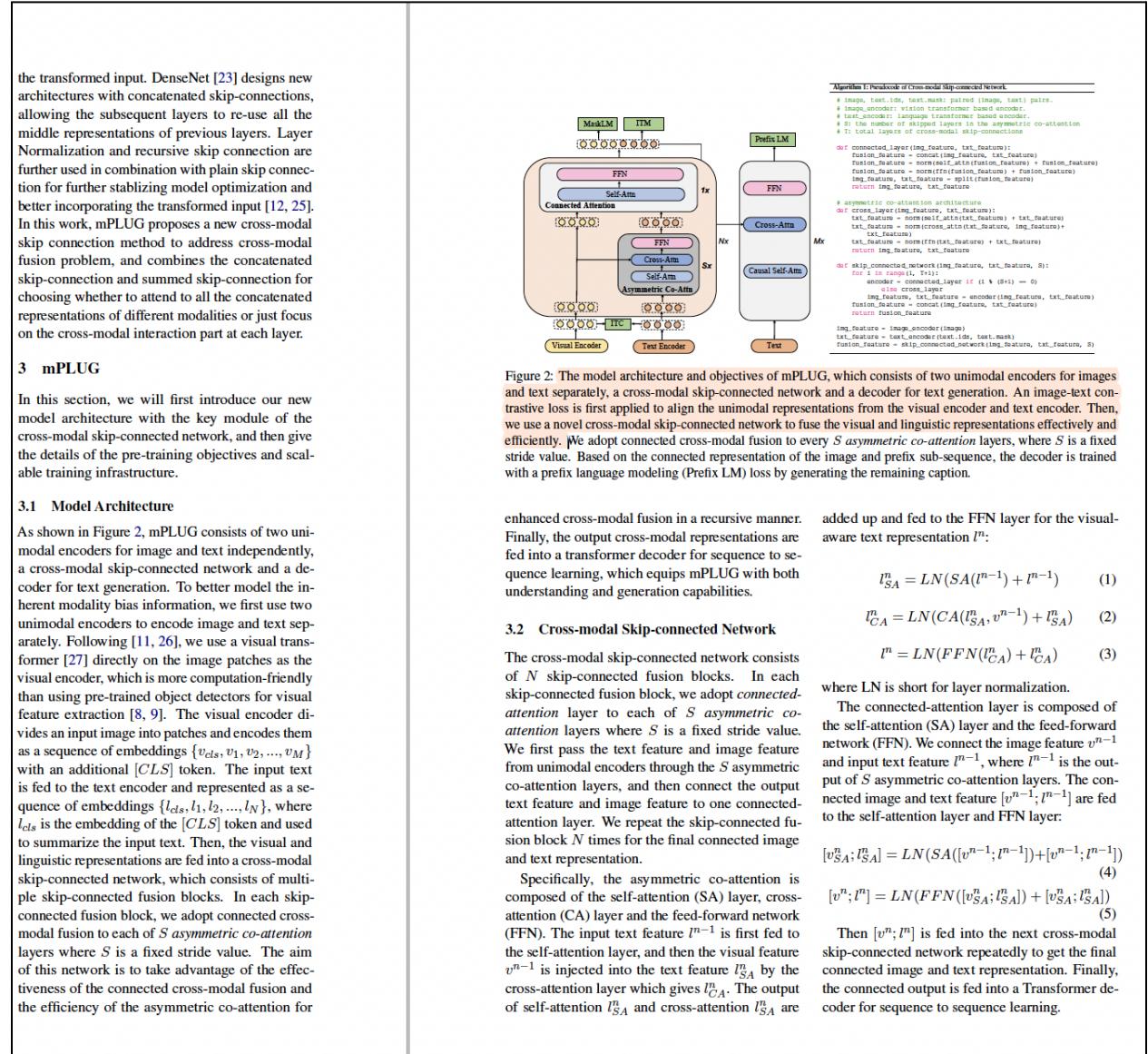


Figure 2: The model architecture and objectives of mPLUG, which consists of two unimodal encoders for image and text separately, a cross-modal skip-connected network and a decoder for text generation. An image-text contrastive loss is first applied to align the unimodal representations from the visual encoder and text encoder. Then, we use a novel cross-modal skip-connected network to fuse the visual and linguistic representations effectively and efficiently. We adopt connected cross-modal fusion to every S asymmetric co-attention layers, where S is a fixed stride value. Based on the connected representation of the image and prefix sub-sequence, the decoder is trained with a prefix language modeling (Prefix LM) loss by generating the remaining caption.

the transformed input. DenseNet [23] designs new architectures with concatenated skip-connections, allowing the subsequent layers to re-use all the middle representations of previous layers. Layer Normalization and recursive skip connection are further used in combination with plain skip connection for further stabilizing model optimization and better incorporating the transformed input [12, 25]. In this work, mPLUG proposes a new cross-modal skip connection method to address cross-modal fusion problem, and combines the concatenated skip-connection and summed skip-connection for choosing whether to attend to all the concatenated representations of different modalities or just focus on the cross-modal interaction part at each layer.

3 mPLUG

In this section, we will first introduce our new model architecture with the key module of the cross-modal skip-connected network, and then give the details of the pre-training objectives and scalable training infrastructure.

3.1 Model Architecture

As shown in Figure 2, mPLUG consists of two unimodal encoders for image and text independently, a cross-modal skip-connected network and a decoder for text generation. To better model the inherent modality bias information, we first use two unimodal encoders to encode image and text separately. Following [11, 26], we use a visual transformer [27] directly on the image patches as the visual encoder, which is more computation-friendly than using pre-trained object detectors for visual feature extraction [8, 9]. The visual encoder divides an input image into patches and encodes them as a sequence of embeddings $\{v_{cls}, v_1, v_2, \dots, v_M\}$ with an additional $[CLS]$ token. The input text is fed to the text encoder and represented as a sequence of embeddings $\{l_{cls}, l_1, l_2, \dots, l_N\}$, where l_{cls} is the embedding of the $[CLS]$ token and used to summarize the input text. Then, the visual and linguistic representations are fed into a cross-modal skip-connected network, which consists of multiple skip-connected fusion blocks. In each skip-connected fusion block, we adopt connected cross-modal fusion to each of S asymmetric co-attention layers where S is a fixed stride value. The aim of this network is to take advantage of the effectiveness of the connected cross-modal fusion and the efficiency of the asymmetric co-attention for

enhanced cross-modal fusion in a recursive manner. Finally, the output cross-modal representations are fed into a transformer decoder for sequence to sequence learning, which equips mPLUG with both understanding and generation capabilities.

3.2 Cross-modal Skip-connected Network

The cross-modal skip-connected network consists of N skip-connected fusion blocks. In each skip-connected fusion block, we adopt *connected-attention* layer to each of S *asymmetric co-attention* layers where S is a fixed stride value. We first pass the text feature and image feature from unimodal encoders through the S asymmetric co-attention layers, and then connect the output text feature and image feature to one connected-attention layer. We repeat the skip-connected fusion block N times for the final connected image and text representation.

Specifically, the asymmetric co-attention is composed of the self-attention (SA) layer, cross-attention (CA) layer and the feed-forward network (FFN). The input text feature l^{n-1} is first fed to the self-attention layer, and then the visual feature v^{n-1} is injected into the text feature l_{SA}^n by the cross-attention layer which gives l_{CA}^n . The output of self-attention l_{SA}^n and cross-attention l_{CA}^n are

added up and fed to the FFN layer for the visual-aware text representation l^n :

$$l_{SA}^n = LN(SA(l^{n-1}) + l^{n-1}) \quad (1)$$

$$l_{CA}^n = LN(CA(l_{SA}^n, v^{n-1}) + l_{SA}^n) \quad (2)$$

$$l^n = LN(FFN(l_{CA}^n) + l_{CA}^n) \quad (3)$$

where LN is short for layer normalization.

The connected-attention layer is composed of the self-attention (SA) layer and the feed-forward network (FFN). We connect the image feature v^{n-1} and input text feature l^{n-1} , where l^{n-1} is the output of S asymmetric co-attention layers. The connected image and text feature $[v^{n-1}; l^{n-1}]$ are fed to the self-attention layer and FFN layer:

$$[v_{SA}^n; l_{SA}^n] = LN(SA([v^{n-1}; l^{n-1}]) + [v^{n-1}; l^{n-1}]) \quad (4)$$

$$[v^n; l^n] = LN(FFN([v_{SA}^n; l_{SA}^n]) + [v_{SA}^n; l_{SA}^n]) \quad (5)$$

Then $[v^n; l^n]$ is fed into the next cross-modal skip-connected network repeatedly to get the final connected image and text representation. Finally, the connected output is fed into a Transformer decoder for sequence to sequence learning.

As you can see, I have made no highlights on sections 3.1, and 3.2. We're missing the context of 10-20 papers that we will need to read to fill in the gaps. That's okay: this is the 20-30% of the paper we are not going to have an in-depth understanding of. We can still squeeze out some understanding from the methods section using Figure 2, which presents a simpler description of the architecture and objectives of mPlug. We also have python pseudocode that actually makes it easier to work through the operations and flow of the model. In general, well constructed figures and algorithm pseudocode can be of huge help to us!

What we can do for the methods section is maintain a list of concepts that we haven't quite understood: if there's a link to the paper references, we'll copy it over. Here's what that might look like:

To Understand:

- From mPLUG:
 - Self-attention + Cross-attention?
 - Detail: Layer normalization?
 - Image-Text Contrastive (ITC): follows "Align before fuse: Vision and language representation learning with momentum distillation"
 - Prefix Language Modeling (PrefixLM): task follows Palm: Pre-training an autoencoding & autoregressive language model for context-conditioned generation.
 - From related work, "Vlmo: Unified vision-language pretraining with mixture-of-modality-experts" using dual encoder and fusion encoder modules.
 - Cross modal interaction example: "An empirical study of training end-to-end vision-and-language transformers."

You would have thus created a list of concepts you need to learn about, and the relevant paper for each, if the paper specifies any.

Let's continue reading, making our way through the methods sections, and the experiments section, highlighting the parts that we consider relevant to our understanding of the method as it relates to image captioning.

| Models | Data | COCO Caption | | | | | | NoCaps | | | |
|-----------------------------|-------|----------------------------|------|-------|--------------------|------|------|--------|------|-------|------|
| | | Cross-entropy Optimization | | | CIDEr Optimization | | | | | | |
| | | B@4 | M | C | S | B@4 | M | C | S | C | S |
| Encoder-Decoder | CC12M | - | - | 110.9 | - | - | - | - | - | 90.2 | 12.1 |
| E2E-VLP [19] | 4M | 36.2 | - | 117.3 | - | - | - | - | - | - | - |
| VinVL [9] | 5.65M | 38.5 | 30.4 | 130.8 | 23.4 | 41.0 | 31.1 | 140.9 | 25.2 | 97.3 | 13.8 |
| OSCAR [4] | 6.5M | - | - | - | - | 41.7 | 30.6 | 140.0 | 24.5 | 83.4 | 11.4 |
| SimVLM _{large} [7] | 1.8B | 40.3 | 33.4 | 142.6 | 24.7 | - | - | - | - | - | - |
| LEMON _{large} [33] | 200M | 40.6 | 30.4 | 135.7 | 23.5 | 42.3 | 31.2 | 144.3 | 25.3 | 113.4 | 15.0 |
| BLIP [34] | 129M | 40.4 | - | 136.7 | - | - | - | - | - | 113.2 | 14.8 |
| OFA [35] | 18M | - | - | - | - | 43.5 | 31.9 | 149.6 | 26.1 | - | - |
| mPLUG | 14M | 43.1 | 31.4 | 141.0 | 24.2 | 46.5 | 32.0 | 155.1 | 26.0 | 114.8 | 14.8 |

Table 1: Evaluation Results on COCO Caption “Karpathy” test split and NoCaps validation set. B@4: BLEU@4, M: METEOR, C: CIDEr, S: SPICE.

3.3 Pre-training Tasks

We perform four pre-training tasks including three understanding tasks (Image-Text Contrastive Learning, Image-Text Matching, Masked Language Modeling) and one generation task (Prefix Language Modeling). These pre-training tasks are optimized jointly.

Image-Text Contrastive (ITC): Following [6], we employ the task to align the image features and the text features from the unimodal encoders. Specifically, we calculate the softmax-normalized image-to-text and text-to-image similarity, and take two dynamic memory queues (text, image) to increase the number of negative examples as MoCo [28].

Image-Text Matching (ITM): This task aims to predict whether an image and a sentence match with each other on the cross-modal representation. We also select hard negative image-text pairs based on the contrastive text-image similarity as [6].

Masked Language Modeling (MLM): The task setup is basically the same as in BERT [29], where we randomly mask 15% of tokens in text and the model is asked to predict these masked words with the cross-modal representations.

Prefix Language Modeling (PrefixLM): This task aims to generate the caption given an image and predict the text segment subsequent to the cross-modal context as [30]. It optimizes a cross-entropy loss by maximizing the likelihood of text in an autoregressive manner.

4 Distributed Learning on a Large Scale
Training a big model like mPLUG on large-scale datasets faces many efficiency challenges. We increase the throughput from the perspective of reducing memory usage and computation time, thereby accelerating the training of the model.

The memory usage during model training is mainly composed of two aspects: the static memory usage composed of parameters/optimized states/gradients, etc., and the runtime memory usage caused by intermediate variables like activation values. For static memory overhead, we use the ZeRO [31] technique to partition parameters/optimizer states/gradients into the entire data-parallel group, so that the static memory overhead of a single GPU can be approximately reduced to $1/N$, where N denotes the number of GPU cards. We use gradient checkpointing [32] for the runtime memory cost, which greatly reduces the runtime memory usage at the

expense of increasing forward time by recomputing part of the activation values during backward pass without keeping them in memory.

To reduce the computation time, we use BF16 precision training. BF16 is a new data type supported by NVIDIA’s new Ampere architecture GPU like A100. Compared with the previously widely used mixed-precision training of FP16 and FP32, BF16 has the same representation range as FP32, thereby reducing the risk of numerical overflow and improving model convergence stability, and at the same time has the same fast computing speed as FP16.

5 Experiments

5.1 Data & Setup

Following the previous work [6], we use the same pre-training dataset with 14M images with texts, which includes two in-domain datasets (MS COCO [36] and Visual Genome [37]), and three web out-domain datasets (Conceptual Captions [38], Conceptual 12M [39], SBU Caption [40]).

We pretrain the model for 30 epochs with the total batch size of 1024 on 16 NVIDIA A100 GPUs. We use a 6-layer Transformer for both the text encoder and the cross-modal skip-connected network, and a 12-layer Transformer for the decoder. The text encoder is initialized using the first 6 layers of the BERT_{base} [29] model and the skip-connected network is initialized using the last 6 layers of the BERT_{base}. We initialize the visual encoder by CLIP-ViT [17] pretrained on 400M noisy image-text pairs. The visual transformer with ViT-B/16 is used as our base architecture, the one with ViT-L/14 as the large architecture. We use the AdamW [41] optimizer with a weight decay of 0.02. The learning rate is warmed up to 1e-5 (ViT-B/16) and 1e-4 (BERT_{base}) for mPLUG_{Gvt}, and 5e-6 (ViT-L/14) and 5e-5 (BERT_{base}) for mPLUG_{mt}, in the first 1000 iterations, and decay to 1e-6 following a cosine schedule. During pre-training, we take random image crops of resolution 256×256 (ViT-B/16)/ 224×224 (ViT-L/14) as inputs and also apply RandomCrop [42] to improve the generalization of vision encoder. For VQA and image captioning tasks, we do an additional continue pre-training on 4M image-text pairs. We increase the image resolution during finetuning. For image-text contrastive learning, the queue size is set as 65,536 and the momentum coefficient is set as 0.995.

| Models | Data | COCO Caption | | | | | | NoCaps | | | |
|-----------------------------|-------|----------------------------|------|-------|--------------------|------|------|--------|------|-------|------|
| | | Cross-entropy Optimization | | | CIDEr Optimization | | | | | | |
| | | B@4 | M | C | S | B@4 | M | C | S | C | S |
| Encoder-Decoder | CC12M | - | - | 110.9 | - | - | - | - | - | 90.2 | 12.1 |
| E2E-VLP [19] | 4M | 36.2 | - | 117.3 | - | - | - | - | - | - | - |
| VinVL [9] | 5.65M | 38.5 | 30.4 | 130.8 | 23.4 | 41.0 | 31.1 | 140.9 | 25.2 | 97.3 | 13.8 |
| OSCAR [4] | 6.5M | - | - | - | - | 41.7 | 30.6 | 140.0 | 24.5 | 83.4 | 11.4 |
| SimVLM _{large} [7] | 1.8B | 40.3 | 33.4 | 142.6 | 24.7 | - | - | - | - | - | - |
| LEMON _{large} [33] | 200M | 40.6 | 30.4 | 135.7 | 23.5 | 42.3 | 31.2 | 144.3 | 25.3 | 113.4 | 15.0 |
| BLIP [34] | 129M | 40.4 | - | 136.7 | - | - | - | - | - | 113.2 | 14.8 |
| OFA [35] | 18M | - | - | - | - | 43.5 | 31.9 | 149.6 | 26.1 | - | - |
| mPLUG | 14M | 43.1 | 31.4 | 141.0 | 24.2 | 46.5 | 32.0 | 155.1 | 26.0 | 114.8 | 14.8 |

Table 1: Evaluation Results on COCO Caption “Karpathy” test split and NoCaps validation set. B@4: BLEU@4, M: METEOR, C: CIDEr, S: SPICE.

| Models | Data | Test-dev | | Test-std | | Methods Pretrained on More Data |
|----------------------|------|----------|---|----------|---|---------------------------------|
| | | | | | | |
| | | | | | | |
| ALBEF [6] | 14M | 75.84 | - | 76.04 | - | |
| BLIP [34] | 129M | 78.25 | - | 78.32 | - | |
| SimVLM [7] | 1.8B | 80.03 | - | 80.34 | - | |
| Florence [45] | 0.9B | 80.16 | - | 80.36 | - | |
| OFA [35] | 18M | 79.87 | - | 80.02 | - | |
| VLMo [20] | - | 79.94 | - | 79.98 | - | |
| mPLUG _{vt} | 14M | 79.79 | - | 79.81 | - | |
| mPLUG _{gvt} | 14M | 81.27 | - | 81.26 | - | |

Table 2: Evaluation Results on VQA test set.

5.2 Evaluation on Vision-Language Tasks

We compare our pre-trained model against other VLP models on the six downstream V+L tasks. We introduce each task and our fine-tuning strategy below. Details of the datasets and fine-tuning hyperparameters are in Appendix.

5.2.1 Visual Question Answering

The VQA task [14] requires the model to answer natural language questions given an image. Most

methods [1, 20, 4, 7] deal with visual question answering tasks as multi-label classification on predefined answer sets. This strategy achieves strong performance, but it is not suitable for real-world open scenarios. We treat VQA as an answer generation task and directly use unconstrained open-vocab generation during inference, which is different from constrained close-vocab generation models [6, 35]. Following [4, 35], we concatenate the question with the object labels and OCR tokens extracted from image. As shown in Table 2, mPLUG achieves 81.27 on Test-std split and outperforms the SOTA models including SimVLM and Florence, which use $100X$ and $60X$ more pre-training image-text pairs, respectively. Based on the same 4M pre-training data, mPLUG outperforms CLIP-VIL and METTER, which also use CLIP [17] as the visual encoder. Besides, under the same settings, mPLUG always significantly outperforms ALBEF and BLIP which only rely on co-attention from images to text for cross-modal fusion. The gain can derive from the network design of cross-modal skip-connections specifically for information asymmetry of the two modalities. Neither ALBEF nor BLIP addresses this problem well, with bias towards the language modality.

5.2.2 Image Captioning

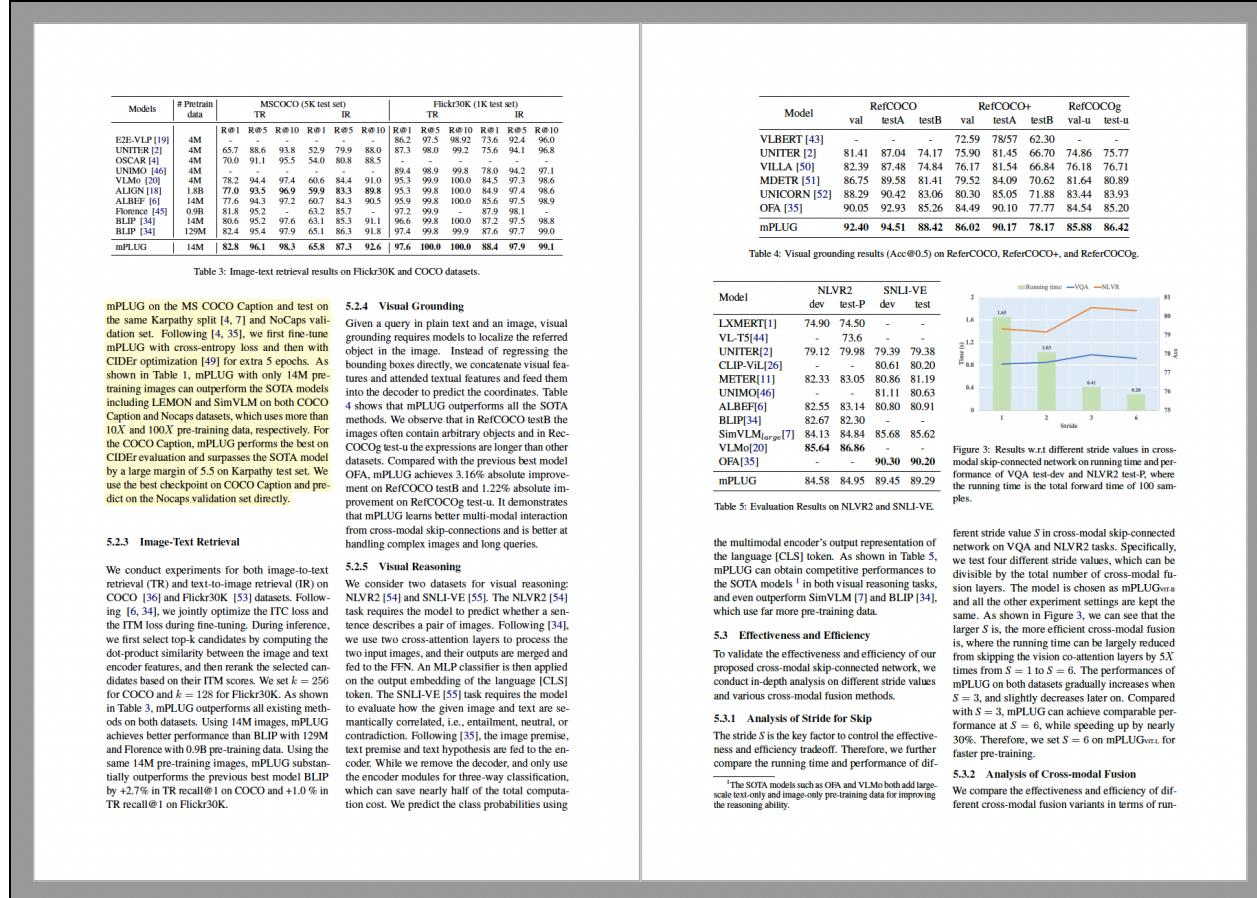
The image captioning task requires a model to generate an appropriate and fluent caption for a given image. We evaluate image captioning on two datasets COCO Caption [47] and NoCaps [48]. mPLUG finetuned with training data of COCO Caption is tested on both of the datasets. We train

Here, I have highlighted the data setup, the main table of results (Table 1) corresponding to image captioning, and the image captioning section. Table 1 includes both methods and evaluation metrics that may be unfamiliar to us, and we can make not of that in our “To Understand” list.

To Understand:

- From mPLUG:
 - Methods
 - Self-attention + Cross-attention?
 - Detail: Layer normalization?
 - Image-Text Contrastive (ITC): follows “Align before fuse: Vision and language representation learning with momentum distillation”
 - Prefix Language Modeling (PrefixLM): task follows Palm: Pre-training an autoencoding & autoregressive language model for context-conditioned generation.
 - From related work, “Vlmo: Unified vision-language pretraining with mixture-of-modality-experts” using dual encoder and fusion encoder modules.
 - Cross modal interaction example: “An empirical study of training end-to-end vision-and-language transformers.”
 - Experiments:
 - What are BLEU-4, METEOR, CIDEr, SPICE?
 - Method comparisons: SimVLM is a competitive method. Read “Simvlm: Simple visual language model pretraining with weak supervision”

We continue reading through the rest of the experiments section



While we can read through the results on tasks other than “image captioning,” I have ignored highlighting them. When reading papers, we are often being presented with a firehose of information, and it’s useful to be selective in what we pay attention to. Let’s continue reading.

¹The SOTA models such as OFA and VLMo both add large-scale text-only and image-only pre-training data for improving the reasoning ability.

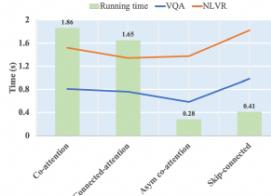


Figure 4: Results w.r.t different cross-modal fusions on running time and performance on VQA test-dev and NLVR2 test-P, where the running time is the total forward time of 100 samples.

| Model | Throughput (Samples/S) |
|-----------------------|------------------------|
| baseline | 124.0 |
| + BFLOAT16 | 182.7 |
| + Gradient Checkpoint | 238.2 |
| + ZeRO | 422.5 |

Table 6: Training Throughput

ning time and performances on VQA and NLVR2 tasks. Specifically, we pre-train mPLUG with different cross-modal fusion network based on the same image encoder and text encoder. All the pre-training settings and the number of fusion layers are kept the same as in the original mPLUG pre-training. As shown in Figure 4, the fusion methods of co-attention and connected-attention both requires much more running time due to long visual sequence. Compared with the two fusion methods, our proposed skip-connected network is $4X$ faster and obtain better performance on both datasets. We also compare it with the asymmetric co-attention used in BLIP [6, 34] which only relies on the co-attention layers from images to text. Despite running slightly faster than the skip-connected network does, the asymmetric co-attention performs worse in accuracy on both datasets. The performance degradation is attributed to the information asymmetry and bias towards language, as shown in Section 5.2.1.

5.3.3 Large-scale Training

Combining the techniques introduced in Section 4 has dramatically increased the training throughput. With the utilization of memory saving and acceler-

| Model | In | Near | Out | Overall |
|------------------------------|--------------|---------------|--------------|--------------|
| SimVLM _{base} [7] | 83.2 | 84.1 | 82.5 | 83.5 |
| SimVLM _{huge} [7] | 101.2 | 100.4 | 102.3 | 101.4 |
| Oscar† [4] | 85.4 | 84.0 | 80.3 | 83.4 |
| VinVL† [9] | 103.7 | 95.6 | 83.8 | 94.3 |
| SimVLM _{huge} † [7] | 113.7 | 110.9 | 115.2 | 112.2 |
| mPLUG | 86.34 | 81.5 | 90.49 | 84.02 |
| mPLUG† | 116.7 | 113.75 | 117.0 | 114.8 |

Table 7: Image captioning results on NoCaps validation split (zero-shot and finetuned), and {In, Near, Out} refer to in-domain, near-domain and out-of-domain respectively. † denotes the models finetuned on COCO Caption dataset.

| Model | TR | | IR | |
|------------------|-------------|-------------|-------------|-------------|
| | R@1 | R@5 | R@1 | R@5 |
| <i>Zero-Shot</i> | | | | |
| CLIP [17] | 88.0 | 98.7 | 68.7 | 90.6 |
| ALIGN [18] | 88.6 | 98.7 | 75.7 | 93.8 |
| FLIP [56] | 89.8 | 99.2 | 75.0 | 93.4 |
| Florence [45] | 90.9 | 99.1 | 76.7 | 93.6 |
| ALBEP† [6] | 94.1 | 99.5 | 82.8 | 96.3 |
| BLIP† [34] | 94.8 | 99.7 | 84.9 | 96.7 |
| mPLUG | 93.0 | 99.5 | 82.2 | 95.8 |
| mPLUG† | 95.8 | 99.8 | 86.4 | 97.6 |

Table 8: Zero-shot image-text retrieval results on Flickr30K. † denotes the models finetuned on COCO.

ated training techniques, the throughput of mPLUG improves $3X$ more from 124 samples per second to 422 samples per second, as shown in Table 6.

5.4 Zero-shot Transferability

In this section, we examine the generalization of mPLUG and compare the zero-shot result on two Vision-Language and three Video-Language tasks.

5.4.1 Zero-shot Vision-Language Tasks

The pretraining of mPLUG adopts image-text contrastive and prefix language modeling tasks on large-scale image-text pairs. Thus, mPLUG has zero-shot generalization ability in image-text retrieval and image captioning. **Image Caption:** First, we take the pretrained mPLUG model and directly decode on NoCaps validation set without further finetuning. Following [7, 34], we feed a prefix prompt “A picture of” into the text encoder to improve the quality of decoded captions.

| Model | # Pretrain data | MSRVTT-Retrieval | | |
|------------------|-----------------|------------------|-------------|-------------|
| | | R@1 | R@5 | R@10 |
| <i>Zero-Shot</i> | | | | |
| MIL-NCE [57] | How100M | 9.9 | 24.0 | 32.4 |
| VideoCLIP [58] | How100M | 10.4 | 22.2 | 30.0 |
| VATT [59] | How100M, AudSet | - | - | 29.7 |
| ALPRO [60] | W2M, C3M | 24.1 | 44.7 | 55.4 |
| VIOLET [61] | Y180M, W2M, C3M | 25.9 | 49.5 | 59.7 |
| CLIP [17] | WIT400M | 26.0 | 49.4 | 60.7 |
| Florence [45] | FLD900M | 37.6 | 63.8 | 72.6 |
| BLIP † [34] | 129M | 43.3 | 65.6 | 74.7 |
| mPLUG | 14M | 38.1 | 59.2 | 68.2 |
| mPLUG † | 14M | 44.3 | 66.4 | 75.4 |

Table 9: Zero-shot video-language results on text-to-video retrieval on the 1k test split of the MSRVTT dataset. † denotes the models finetuned on COCO. Video datasets include HowTo100M [62], WebVid-2M(W2M) [63], YT-Temporal-180M(Y180M) [64]. Image datasets include CC3McC3M) [38], FLD900M [45], WIT400M [17]. Audio datasets include AudioSet(AudSet) [65].

As shown in Table 7, the zero-shot performance of mPLUG is competitive with fully supervised baselines such like Oscar and VinVL. With further finetuning on MSCOCO dataset, mPLUG outperforms the SimVLM_{huge}, which use more pre-training image-text pairs and has larger model parameters. **Image-text Retrieval:** We perform zero-shot retrieval on Flickr30K. The result is shown in Table 8, where zero-shot mPLUG outperforms models (CLIP, ALIGN, Florence) pretrained with more image-text pairs. Following [34], we also evaluate zero-shot retrieval by the model finetuned on MSCOCO dataset. Table 8 shows that mPLUG achieves better performance than the previous SOTA models.

5.4.2 Zero-shot Transfer to Video-Language Tasks

To evaluate the generalization ability of mPLUG to Video-Language Tasks, we conduct zero-shot experiments on Video-text Retrieval, Video Caption and Video Question Answering. Following [34], we uniformly sample n frames for each video ($n = 8$ for Retrieval, $n = 16$ for QA, $n = 8$ for Caption), and concatenate the frame features into a single sequence. **Video-text Retrieval:** We evaluate the mPLUG models pretrained and further finetuned on the COCO-retrieval image-text dataset

You can now see that we have found another section of experiments. Zero-shot transferability may not be a concept we’re familiar with, so we likely will have to add that to our “to understand” list.

Finally, let’s read the conclusion of the paper.

6 Conclusion

This paper presents mPLUG, an effective and efficient VLP framework for both cross-modal understanding and generation. mPLUG introduces a new asymmetric vision-language architecture with novel cross-modal skip-connections, to address two fundamental problems of information asymmetry and computation efficiency in cross-modal alignment. Pretrained on large-scale image-text pairs, mPLUG achieves state-of-the-art performance on a wide range of vision-language tasks. mPLUG also demonstrates strong zero-shot transfer ability when directly applied to multiple video-language tasks. Our work explores the cross-modal alignment with a newly-designed VLP architecture and we hope it can help promote future research on image-text foundation models.

You might notice that although the conclusion was similar to the abstract and to the introduction, we already have a much better understanding of these words than we did towards the start of reading this paper. As I mentioned before, reading the paper the first time, we expect to really only understand ~10% of it, especially if we don't have the background of the papers which are being built on. With this same approach, we can start making our way through our "To Understand" list.

Conclusion

I hope this walkthrough of a first read of an AI paper in a new topic gives you the confidence to dive into a new problem topic. Both the process of reading wide and reading deep are iterative: we often need to re-search and re-read, and the act of making notes as you go can significantly help in building and cross-checking your mental model.