# BLIP paper summary

## BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation

# 1. High-Level Brief

### Introduction

The BLIP framework emerges as a significant advancement in the field of Vision-Language Pre-training (VLP), addressing critical challenges faced by its predecessors. Traditional VLP models often struggled with the integration of image and language modalities, particularly in efficiently handling the noisy data sourced from the web and being restricted to either understanding-based or generation-based tasks. BLIP's innovative approach lies in its ability to bridge these gaps, offering a more holistic and versatile solution to VLP.

### BLIP's Contributions

**BLIP introduces two key innovations that set it apart from existing methods:**

- **Multimodal Mixture of Encoder-Decoder (MED):** This novel architecture is the cornerstone of BLIP, enabling the model to adapt flexibly to a wide range of tasks. Unlike traditional models that are typically confined to either understanding or generation tasks, MED allows BLIP to excel in both. It achieves this by functioning in multiple modes – as a unimodal encoder for separate processing of image and text, an image-grounded text encoder for integrating visual information into text encoding, and an image-grounded text decoder for generating textual descriptions from images.
- **Captioning and Filtering (CapFilt):** Addressing the challenge of noisy web data, BLIP incorporates CapFilt, a sophisticated data refinement process. CapFilt enhances the quality of training data by generating synthetic captions for images and then filtering out those that are noisy or less informative. This process not only cleanses the dataset but also enriches it, contributing significantly to the model's training and subsequent performance.

BLIP's architecture and data processing strategy represent a substantial leap forward in VLP, offering a more effective and efficient way to leverage the vast amounts of web-sourced image-text pairs. The framework's ability to handle diverse

data and tasks positions it as a highly adaptable and robust solution in the realm of multimodal AI.

## 2. Datasets Used

Overview of Datasets

BLIP's training leverages a comprehensive blend of datasets, each contributing unique aspects to the model's development:

- Human-Annotated Datasets:
    - COCO (Common Objects in Context)
    - Visual Genome
    - Web-Crawled Datasets:
    - Conceptual Captions
    - SBU Captions
    - LAION Dataset

**Purpose of Diverse Data**

- **Diversity and Scale**: The combination of structured, human-annotated datasets with the vast, unstructured web-crawled data ensures that BLIP is exposed to a wide array of visual and linguistic contexts. This diversity is key to developing a model capable of understanding and generating a broad spectrum of human language and visual scenes.
- **Quality and Challenge:** While datasets like COCO and Visual Genome provide high-quality, reliable annotations, web-crawled datasets introduce the challenge of noisy, less reliable data. This mix allows BLIP to not only learn from the best-available data but also to develop robustness against lower-quality inputs, a common scenario in real-world applications.
- **Real-World Applicability:** By training on data that reflects a wide range of real-world scenarios, BLIP is better equipped to handle practical applications in various domains, making it a versatile and powerful tool in vision-language tasks.

## 3. GPUs and Training Setup

**Computational Resources**

The training of the BLIP model was a computationally intensive task, demanding significant hardware capabilities:

- **GPU Configuration:**
  - The model was trained on two nodes, each equipped with 16GB high-performance GPUs.

**Training Strategy**

BLIP's training strategy was carefully designed to maximize model performance and efficiency:

- **Training Parameters:**
  - **Batch Size:** Large batch sizes were used (2880 for ViT-B and 2400 for ViT-L).
  - **Learning Rate:** The model employed the AdamW optimizer with a learning rate initially warmed-up to 3e-4 for ViT-B / 2e-4 for ViT-L, followed by a linear decay. This strategy helped in stabilizing the training initially and then fine-tuning the model weights effectively.
  - **Weight Decay:** A weight decay of 0.05 was set to regularize the model and prevent overfitting.
- **Training Epochs:**
  - BLIP was trained for 20 epochs, ensuring sufficient exposure to the diverse training data for robust learning while balancing the computational costs.
- **Image Resolution:**
  - During pre-training, random image crops of 224×224 resolution were used, which were then increased to 384×384 during fine-tuning. This change allowed the model to first learn from more general features and then adapt to finer details for better performance.

Importance of Robust Training

The extensive GPU resources and the meticulously planned training strategy were instrumental in developing BLIP. This robust training setup ensured that the model could effectively learn from the diverse datasets, adapt to various tasks, and ultimately perform at a state-of-the-art level across different vision-language benchmarks.

# 4. Main Algorithms and Their Usage

**Nucleus Sampling**

- Purpose: Nucleus Sampling is employed for generating diverse synthetic captions, which are crucial for training the BLIP model with varied and rich language inputs.

- Process: It involves selecting the next word in a caption based on a cumulative probability distribution, ensuring that the generated captions are not just high probability but also diverse.
- Pseudo code for top-p or Nucleus sampling

```
function nucleus_sampling(token_probs, p_threshold):
    sorted_probs, sorted_tokens = sort(token_probs, descending=True)
    cumulative_probs = cumulative_sum(sorted_probs)
    selected_tokens = sorted_tokens[cumulative_probs <= p_threshold]
    return random_choice(selected_tokens)
```

## Image-Text Contrastive Loss (ITC)

- Purpose: Used to align the feature space of visual and textual inputs, essential for tasks like image-text retrieval.
- Method: It involves calculating similarities between image and text features and using these for contrastive learning.

```
function itc_loss(image_features, text_features, temperature):
    similarities = cosine_similarity(image_features, text_features) /
temperature
    loss = contrastive_loss(similarities)
    return loss
```

## Image-Text Matching Loss (ITM)

- Purpose: This loss helps the model learn finer alignments between image and text pairs, crucial for understanding the specific connections between visual and textual elements.
- Implementation: The model predicts whether an image-text pair matches and this prediction is used in a binary classification loss.

```
function itm_loss(image_text_pairs, labels):
    predictions = model.predict(image_text_pairs)
    loss = binary_cross_entropy(predictions, labels)
    return loss
```

## Language Modeling Loss (LM)

- Role: Used in the text decoding process, enabling the model to generate coherent and contextually relevant textual descriptions from images.
- Approach: The model maximizes the likelihood of generating the correct next word in a caption given an image.

```
function lm_loss(captions, model):
    predictions = model.generate(captions)
    loss = cross_entropy(predictions, captions)
    return loss
```

**Cross-Attention Mechanism**

- Function: Allows the model to integrate and focus on relevant features from both image and text modalities, enhancing the overall understanding and generation capabilities.
- Mechanism: Involves computing attention weights between image and text features and applying these to generate a combined representation.

```
function cross_attention(query, key, value):
    attention_scores = softmax(matmul(query, key.T) / sqrt(size_of_key))
    output = matmul(attention_scores, value)
    return output
```

**Multimodal Mixture of Encoder-Decoder (MED)**

- Design: A versatile architecture that can function as an encoder, grounded encoder, or grounded decoder, accommodating various types of vision-language tasks.
- Structure: Combines different aspects of transformer-based models tailored for both visual and textual processing.

```
class MED:
    def __init__(self, image_encoder, text_encoder, text_decoder):
        self.image_encoder = image_encoder
        self.text_encoder = text_encoder
        self.text_decoder = text_decoder

    def forward(self, images, text, task_type):
        if task_type == 'encode':
            return self.text_encoder.encode(text),
self.image_encoder.encode(images)
```

```
        elif task_type == 'grounded_encode':
            return self.text_encoder.grounded_encode(text, images)
        elif task_type == 'grounded_decode':
            return self.text_decoder.decode(images)
```

**Captioning and Filtering (CapFilt)**

- Objective: To refine the training dataset by generating high-quality synthetic captions and filtering out the noisy or irrelevant ones.
- Process: Uses a caption generator to create captions for images and a filter to remove suboptimal captions, ensuring the data used for training is of high quality.

```
function capfilt(images, captioner, filter):
   synthetic_captions = captioner.generate_captions(images)
   filtered_captions = filter.remove_noisy_captions(images,
synthetic_captions)
   return filtered_captions
```

# 5. Conclusion and Total Pipeline

- **Novel Architectural Approach:** The introduction of the Multimodal Mixture of Encoder-Decoder (MED) architecture has set new standards in flexibility and efficiency for VLP tasks, enabling the model to adeptly handle both understanding-based and generation-based tasks.
- **Enhanced Data Quality with CapFilt:** By implementing the Captioning and Filtering process, BLIP effectively improves the quality of its training data, a crucial factor in the success of any machine learning model. This approach not only filters out noise but also enriches the dataset with diverse and informative synthetic captions, leading to more robust learning.
- **State-of-the-Art Performance:** Across a range of vision-language tasks, including image-text retrieval, image captioning, and visual question answering, BLIP has demonstrated superior performance, thanks in part to its ability to learn from a varied and extensive dataset.