# Exploring Model Compression for Edge-Optimized Video Description
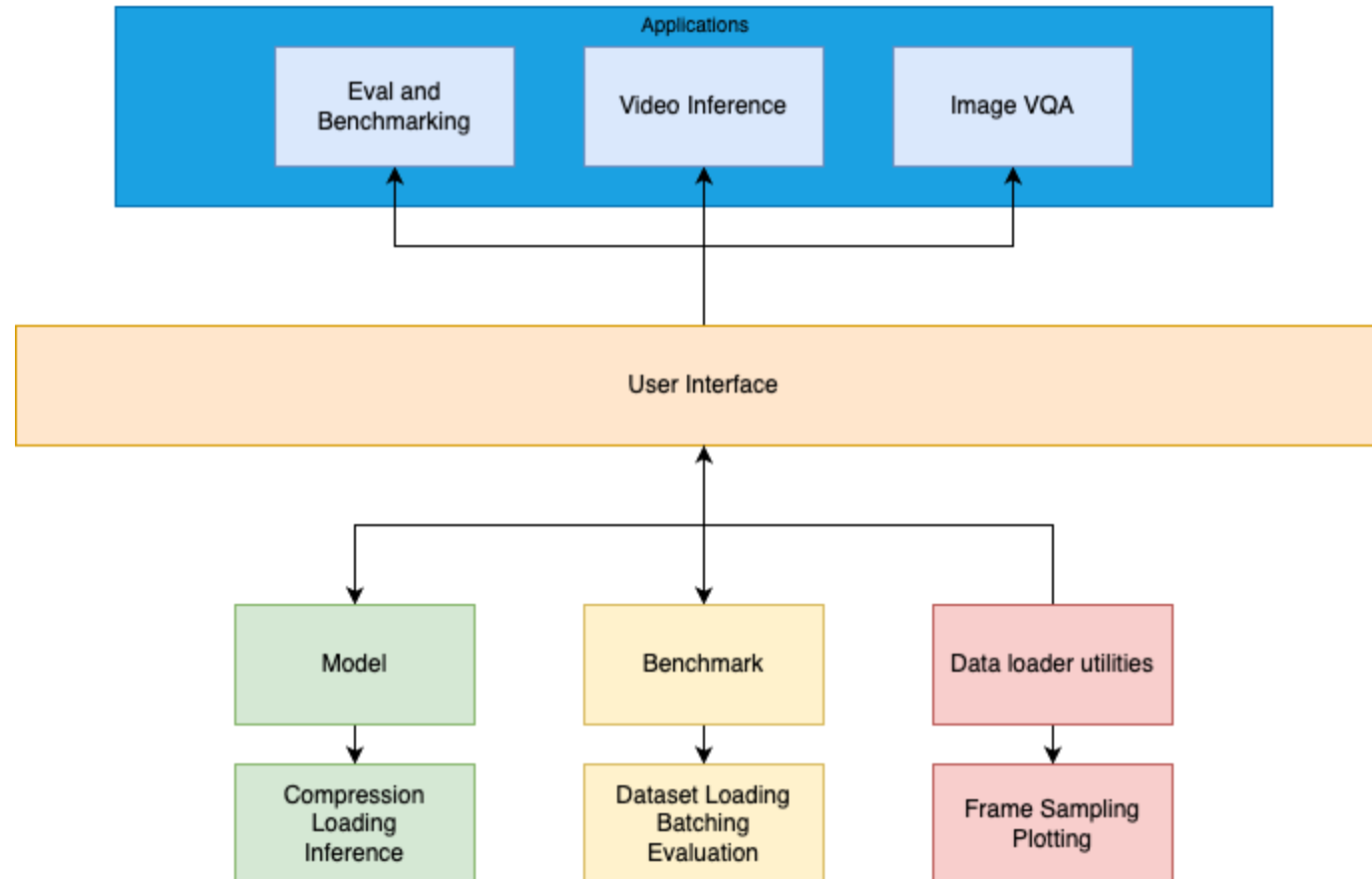
Members:

Amit Phabba, Raj Shah, Sukrit Kumar, Dhruv Rauthan, Dawei Qin

Georgia Tech

# System Design

- Fully plug and play design
- Allows great flexibility and ability to extend to new benchmarks, compression methods

# Distribution of work

| Task | Person Responsible |
| --- | --- |
| Model Quantization pipeline setup | Sukrit, Raj, Dhruv, Dawei, Amit |
| DocVQA and VQAv2 Benchmark setup + Eval pipeline | Sukrit (Slide 6,7,8) |
| Low-rank weight Factorization + fine-tuning | Raj (Refer Slide: 39 to 52) |
| Model Pruning + Post prune fine-tuning | Dawei, Sukrit (Refer Slide: 35 to 38) |
| ScienceQA Benchmark setup + Eval pipeline | Dawei (Slide 9) |
| Flickr30k Benchmark setup + Eval pipeline | Raj (Slide 10) |
| Video Inference Pipeline setup | Sukrit, Raj, Dhruv, Dawei, Amit (Slide 55 for demo, used throughout) |
| Benchmark Launcher, Eval.py | Sukrit, Raj, Dhruv, Dawei, Amit |
| Demo, Workshop Presentation | Sukrit, Raj, Dawei, Dhruv, Amit |
| Demo benchmark result analysis | Dhruv, Raj |

Georgia Tech

# Benchmark Examples

- We evaluate on 4 benchmarks :

  o DocVQA

  o VQAv2

  o ScienceVQA

  o Flickr30

- For all 4 benchmarks we compute the BERT Score as evaluation
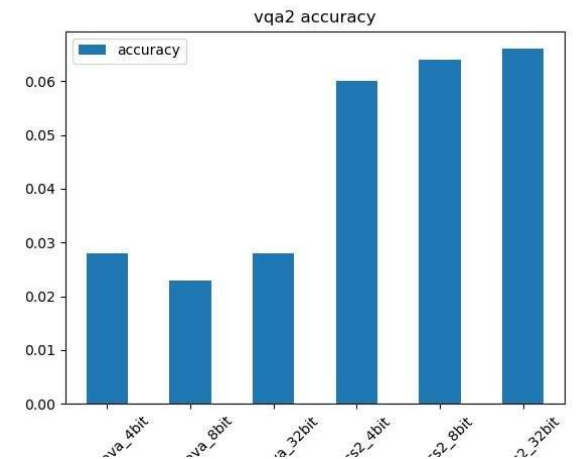
Georgia Tech.

# Evaluation Method

From Levenshtein Similarity to BERT Score:

We started with Levenshtein Similarity as the evaluation method but it did not work well:



Running VQA2 with levenshtein, as we can see, the output is good but receives poor score.

# Evaluation Method

Levenshtein similarity looks at the distance between the output and answer, so even if the answer is correct, it can get a bad score, as it does not know the meanings of the sentences.

Bert score, however, it uses pretrained contextual embeddings from BERT, which allows understanding of meanings.

BERT result:

# DocVQA

- [DocVQA](#) is a image-doc QA benchmark, where the VLM is provided with an image and a corresponding Question expecting an answer extracted from the image. The benchmark can take multiple possible answers and the score is calculated keeping this in mind



Query: what is the date mentioned in this letter?
Answer: 1/8/93



Query: Who is in cc in this letter?
Answer: T.F. Riehl



Query: What is the name of the company
Answer: ITC Limited

Query: what is the date mentioned in this letter? Answer: April 27, 1993





Query: What is the location of the NSDA
Answer: T.F. Riehl

# VQAv2

- VQAv2 is another popular vision QA benchmark. It consists of a single image along with a corresponding query and generally a single word answer



Query: How many sinks are in the bathroom
Answer: 4



Query: What is the source of the red lines in the picture
Answer: Tail Lights



Query: What is on top of the zebra's head and down the back of his neck?
Answer: Mane



Query: How many sheep are there?
Answer: 1



Query: Is anyone sitting on the bikes
Answer: Yes

Georgia Tech

# ScienceQA

Science QA is a dataset which includes images, queries and answers on the topic of natural science, language science, and social science. The questions are in MCQ formats which makes evaluation easier. The questions are collected from elementary and high school classes. According to the official page, some of the questions do not have image content, so we filter the invalid question by code:

```python
# Skip questions without images
valid_examples = [
    example for example in examples if example["image"] and example["image"][0] is not None
]

if not valid_examples:
    #print(f"Skipping batch {i} as no valid examples with images found.", flush=True)
    continue

# Process valid examples
self.answers_unique.extend([
    chr(ord('A') + int(example["answer"][0])) for example in valid_examples
])
```

Dataset URL: https://scienceqa.github.io/#dataset
Total data size is 21,208, but some are filtered during loading.

Georgia Tech

# Flickr30k

- Flickr30k is an image captioning benchmark, where the Vision-Language Model (VLM) is provided with an image and a corresponding caption. The benchmark measures the model's ability to generate accurate and descriptive captions based on the image content.



A man sits in a chair while holding a large stuffed animal of a lion.

Three people are standing outside near large pipes and a metal railing.

Two young guys with shaggy hair look at their hands while hanging out in the yard.

A girl is on roller skates talking on her cellphone standing in a parking lot.

# Quantization examples for Scienceqa:



Q1. Which of the following could Gordon's test show?



Q2. What is the name of the colony shown?



Q3. Which of these organisms contains matter that was once part of the lichen?



Q4. What is the expected ratio of offspring with a woolly fleece to offspring with a hairy fleece? Choose the most likely ratio.

Q5. Which property do these three objects have in common?



tortoise shell     crown     basketball

Georgia Tech

# Quantization examples for vqa2:



Q1. Where are the kids riding?



Q2. Is this boy a good pitcher?



Q3. What is the person wearing?



Q4. How many sinks are in this bathroom



Q5. What sport are the girls playing?

# Quantization examples for Docvqa:



Q1. What the location address of NSDA?

Q2. According to budget request summary what is total amount of other expenses??

Q3. Who is 'presiding' TRRF GENERAL SESSION (PART 1)?

Q4. How many nomination committee meetings has Y. C. Deveshwar attended?

Q5. How many nomination committee meetings has S. Banerjee attended?

# Quantization demo results [Qwen2-VL]

| Qwen/Scienceqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(B) | B | B | B |
| Q2(B) | B | B | B |
| Q3(B) | B | B | B |
| Q4(B) | D | D | D |
| Q5(C) | C | A | A |

| Qwen/vqa2 | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(carnival ride) | ride | carousel | ride |
| Q2(yes) | yes | no | yes |
| Q3(wetsuit) | wetsuit | wetsuit | wetsuit |
| Q4(4) | 2 | 1 | 2 |
| Q5(soccer) | soccer | soccer | Soccer |

Georgia Tech

# Quantization demo results [Qwen2-VL]

| Qwen/docvqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1('1128 SIXTEENTH ST., N. W., WASHINGTON, D. C. 20036', '1128 sixteenth st., N. W., washington, D. C. 20036') | 1000000 | 10000 S. 100th St. | 1000 N. Michigan Ave., Chicago, IL 60611 |
| Q2('975.00') | $972,000 | $ 972 | $ 975.00 |
| Q3('TRRF Vice President', 'lee a. waller') | TRRF Vice President | Ler A. Willner | Lee A. Waller |
| Q4('2') | 1 | 3 | 2 |
| Q5('2') | 1 | 3 | 2 |

Georgia Tech

# Quantization demo results [Idefics]

| Idefics/Scienceqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(B) | B | B | B |
| Q2(B) | B | B | B |
| Q3(B) | B | B | B |
| Q4(B) | C | C | C |
| Q5(C) | C | C | C |

| Idefics/vqa2 | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(carnival ride) | Fire truck | Fire truck | Fire truck |
| Q2(yes) | Yes | Yes | Yes |
| Q3(wetsuit) | Wetsuit | Wetsuit | Wetsuit |
| Q4(4) | There are two sinks in this bathroom | There are two sinks in this bathroom | There are two sinks in this bathroom |
| Q5(soccer) | Soccer | Soccer | Soccer |

Georgia Tech

# Quantization demo results [Idefics]

| Idefics/docvqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1('1128 SIXTEENTH ST., N. W., WASHINGTON, D. C. 20036', '1128 sixteenth st., N. W., washington, D. C. 20036') | 1128 Sixteenth St. N.W | 1128 Sixteenth St. N.W | 1128 Sixteenth St. N.W. |
| Q2('975.00') | $ 975.00 | $ 975.00 | $ 975.00 |
| Q3('TRRF Vice President', 'lee a. waller') | Lee a. waller | Lee a. waller | Lee a. waller |
| Q4('2') | 2 | 2 | 2 |
| Q5('2') | 2 | 2 | 2 |

# Quantization demo results [LlavaNext]

| LlavaNext/Scienceqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(B) | B | B | B |
| Q2(B) | B | B | B |
| Q3(B) | B | B | B |
| Q4(B) | D | D | D |
| Q5(C) | A | A | A |

| LlavaNext/vqa2 | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(carnival ride) | The kids are riding in a carousel at a fair or amusement park | Carousel | Carousel |
| Q2(yes) | Based on the image, it's difficult to determine if the boy is a good pitcher, as we can only see him in action during a game. However, he appears to be focused and has a proper stance, which is a good start. His form and technique could be improved, but it's | Based on the image, it's difficult to determine if the boy is a good pitcher as we can only see him in action for a single moment. However, he appears to be in a standard pitching stance with a baseball glove on his left hand, which suggests he is prepared to pitch. His | Based on the image, it's difficult to determine if the boy is a good pitcher as we can only see him in action for a single moment. However, he appears to be in a proper pitching stance with a baseball glove on his left hand, which suggests he might be familiar with the sport |
| Q3(wetsuit) | The person is wearing a wetsuit | The person is wearing a wetsuit | The person is wearing a wetsuit |
| Q4(4) | 4 | There are two sinks in this bathroom | There are two sinks in this bathroom |
| Q5(soccer) | The girls are playing soccer | Soccer | The girls are playing soccer |

# Quantization demo results [LlavaNext]

| LlavaNext/docvqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1('1128 SIXTEENTH ST., N. W., WASHINGTON, D. C. 20036', '1128 sixteenth st., N. W., washington, D. C. 20036') | The National Soft Drink Association is located at 1619 15th St. N.W., Washington, D.C. 2000006 | The National Soft Drink Association is located at 1619 13th St., N.W., Washington, D.C. 2000007 | The location address of NSDA (National Soft Drink Association) is 16191 Sixteenth Street, N.W., Washington, D.C. 2000006 |
| Q2('975.00') | The total amount of other expenses listed in the budget request summary is $9770.000 | The total amount of other expenses listed in the budget request summary is $975.000 | The total amount of other expenses listed in the budget request summary is $975.000 |
| Q3('TRRF Vice President', 'lee a. waller') | Lee A. Wallner is presiding the TRRF General Session (Part 1) | Lee A. Wallner is presiding the TRRF General Session (Part 1) | Lee A. Wallner is presiding the TRRF General Session (Part 1) |
| Q4('2') | Y. C. Deveshwar attended two nomination committee meetings during the financial year ended 31st March, 2013 | Y. C. Deveshwar attended two nomination committee meetings during the financial year 2012-13 | Y. C. Deveshwar attended two nomination committee meetings during the financial year 2012-13 |
| Q5('2') | S. Banerjee has attended two nomination committee meetings during the financial year 2012-13 | S. Banerjee has attended two nomination committee meetings during the financial year 2012-13 | S. Banerjee has attended two nomination committee meetings during the financial year 2012-13 |

Georgia Tech

# Quantization demo results [Llava]

| Llava/Scienceqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(B) | B | B | B |
| Q2(B) | A | B | B |
| Q3(B) | B | B | B |
| Q4(B) | D | D | D |
| Q5(C) | A | A | A |

| Llava/vqa2 | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1(carnival ride) | The kids are riding in a small red car at a carnival | The kids are riding in a red and yellow toy car at a carnival | The kids are riding in a red and yellow toy car at a carnival |
| Q2(yes) | Yes, the boy is a good pitcher | Yes, the boy in the black hat is a good pitcher | Yes, the boy in the black hat is a good pitcher |
| Q3(wetsuit) | The person is wearing a wetsuit | The person is wearing a wetsuit | The person is wearing a wetsuit |
| Q4(4) | There are two sinks in this bathroom | There are two sinks in this bathroom | There are two sinks in this bathroom |
| Q5(soccer) | The girls are playing soccer | The girls are playing soccer | The girls are playing soccer |

# Quantization demo results [Llava]

| Llava/docvqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1('1128 SIXTEENTH ST., N. W., WASHINGTON, D. C. 20036', '1128 sixteenth st., N. W., washington, D. C. 20036') | The location address of NSDA is 1200 N. 12th St., Milwaukee, WI 53202 | The location address of NSDA is 1000 N. 12th St., Milwaukee, WI 53202 | The location address of NSDA is 1000 N. 16th St., Milwaukee, WI 53202 |
| Q2('975.00') | According to the budget request summary, the total amount of other expenses is $1,000,000 | The total amount of other expenses is $1,000 | The total amount of other expenses is $1,000 |
| Q3('TRRF Vice President', 'lee a. waller') | Robert H. Cottrell is presiding the TRRF General Session (Part 1) | Robert H. Hartburg is presiding the TRRF General Session (Part 1) | Robert H. Hartberg is presiding the TRRF General Session (Part 1) |
| Q4('2') | Y. C. Deveshwar has attended 10 nomination committee meetings | Y. C. Deveshwar has attended 1 nomination committee meeting | Y. C. Deveshwar has attended 1 nomination committee meeting |
| Q5('2') | S. Banerjee has attended 10 nomination committee meetings | S. Banerjee has attended 1 nomination committee meeting | S. Banerjee has attended 1 nomination committee meeting |

Georgia Tech.

# Quantization demo results [Phi3.5]

| Phi3.5/Scienceqa | 4 bit | 8 bit | 32 bit |
| --- | --- | --- | --- |
| Q1(B) | B | B | B |
| Q2(B) | B | B | B |
| Q3(B) | B | B | B |
| Q4(B) | B | B | B |
| Q5(C) | C | C | C |

| Phi3.5/vqa2 | 4 bit | 8 bit | 32 bit |
| --- | --- | --- | --- |
| Q1(carnival ride) | The kids are riding on a miniature train | The kids are riding on a miniature train | The kids are riding on a miniature train |
| Q2(yes) | The image shows a young baseball player in the middle of a pitch, wearing a black jersey with the number 5 and grey pants. Without additional context or information about his performance, it is not possible to determine if he is a good pitcher | The image shows a boy in a baseball uniform in the middle of a pitching motion on a baseball field. Without additional context or observing the boy's technique, it is not possible to determine if he is a good pitcher | The image shows a young individual in a baseball uniform in the middle of a pitching motion on a baseball field. Without observing multiple pitches or having expert knowledge in baseball, it is not possible to accurately assess the skill level of the pitcher. |
| Q3(wetsuit) | The person is wearing a black wetsuit | The person is wearing a wetsuit | The person is wearing a wetsuit |
| Q4(4) | There are two sinks in the bathroom | There are two sinks in the bathroom | There are two sinks in the bathroom |
| Q5(soccer) | The girls are playing soccer | The girls are playing soccer | The girls are playing soccer |

# Quantization demo results [Phi3.5]

| Phi3.5/docvqa | 4 bit | 8 bit | 32 bit |
|---|---|---|---|
| Q1('1128 SIXTEENTH ST., N. W., WASHINGTON, D. C. 20036', '1128 sixteenth st., N. W., washington, D. C. 20036') | 1185 BIRKENBURG ST., N.V., WASHINGTON, D.C. 20035 | The location address of NSDA is 1185 BIRKENSTUhl St., N.V., Washington, D.C. 20036 | 1185 BIRKENHEAD STREET, N.V., WASHINGTON, D.C. 20036 |
| Q2('975.00') | $9,755.00 | $975.00 | $975.00 |
| Q3('TRRF Vice President', 'lee a. waller') | Lee A. Waller | Lee A. Waller | Lee A. Waller |
| Q4('2') | 2 | 2 | 2 |
| Q5('2') | 2 | 2 | S. Banerjee has attended 2 nomination committee meetings |

Georgia Tech

# Quantization Results (Accuracy)

| benchmark | idefics2_32bit | idefics2_4bit | idefics2_8bit | llavanext_32bit | llavanext_4bit | llavanext_8bit | phi3_5_32bit | phi3_5_4bit | phi3_5_8bit | qwen2_32bit | qwen2_4bit | qwen2_8bit | video_llava_32bit | video_llava_4bit | video_llava_8bit |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| docvqa | 0.926 | 0.904 | 0.922 | 0.734 | 0.722 | 0.745 | 0.888 | 0.885 | 0.883 | 0.858 | 0.730 | 0.805 | 0.680 | 0.675 | 0.676 |
| scienceqa | 0.612 | 0.611 | 0.611 | 0.612 | 0.584 | 0.581 | 0.563 | 0.563 | 0.563 | 0.645 | 0.706 | 0.645 | 0.882 | 0.623 | 0.629 |
| vqa2 | 0.931 | 0.928 | 0.933 | 0.777 | 0.740 | 0.803 | 0.701 | 0.694 | 0.702 | 0.931 | 0.840 | 0.908 | 0.745 | 0.734 | 0.732 |
| flickr30k | 0.797 | 0.788 | 0.798 | 0.805 | 0.806 | 0.803 | 0.833 | 0.842 | 0.839 | 0.822 | 0.762 | 0.829 | 0.819 | 0.820 | 0.818 |

# Quantization Complete Results



DOCVQA Benchmark Results: Accuracy vs Memory Usage

SCIENCEQA Benchmark Results: Accuracy vs Memory Usage

VQA2 Benchmark Results: Accuracy vs Memory Usage

FLICKR30K Benchmark Results: Accuracy vs Memory Usage

# Quantization Results

| Model Name | Original Size | 4-Bit Quantized Size | 8-Bit Quantized Size |
|---|---|---|---|
| Qwen 2VL 2B | 6.17 GB | 1.12 GB **(-81%)** | 1.77 GB **(-71%)** |
| Llava 7B | 29.46 GB | 4.08 GB **(-86%)** | 7.63 GB **(-74%)** |
| LLava-Next | 28.35 GB | 3.92 GB **(-86%)** | 7.32 GB **(-74%)** |
| Idefics2 8B | 16.80 GB | 4.60 GB **(-72%)** | 8.67 GB **(-48%)** |
| Phi-3.5 Vision | 16.58 GB | 2.37 GB **(-85%)** | 4.34 GB **(-73%)** |

# Quantization Performance Analysis

- As expected, the accuracy either decreases or stays relatively similar as the 8 bit and 4 bit quantizations are applied
  - Exception is Qwen for ScienceQA, where the 4 bit quantized model performs the best relatively- since the benchmark is MCQ based, the BERT score analysis is not applicable here, as the output is either completely correct or incorrect
- Idefics is the most consistent performer across the different benchmarks
  - However, has the highest memory usage as well
- LlavaNext performs either equal to or better than Llava
  - This again is expected as LlavaNext is an improvement over the latter
- The quantized models are far superior in terms of the memory usage, across all the different models
- Qwen seems to the overall best, requiring minimal memory while also highly accurate- followed by Phi3.5 which sacrifices accuracy

Georgia Tech.

# Quantization Memory Analysis

- Qwen is the most compact model post-quantization
  - Best for low-resource environments
  - Even the original model size is comparable to other models' quantized versions
- Llava and LlavaNext benefit most from compression in absolute and percentage terms
  - Since LlavaNext is an optimized and improved version of Llava, the model sizes are somewhat similar
- Idefics is the worst performing in terms of the reduction in model size
  - Perhaps due to existing optimizations in its architecture

Georgia Tech

# Sample input for video inference

- We used the same video across the 5 models to get an idea about how the response varies over models and compressions ratios

# Comparison of outputs (4 bit quantized)

| Qwen | Phi3.5 | Idefics2 | LlavaNext | Llava |
|---|---|---|---|---|
| The video captures a serene and colorful scene of a person in a blue shirt and white pants, wearing a red and white striped hat, walking on a beach. The person is seen walking along the shoreline, with the ocean and beach visible in the background. The video showcases the beauty of the beach and the person's attire, creating a harmonious and picturesque scene. | In this video, a young girl is sitting on a bed and reading a book. She is wearing a green dress and glasses. The room appears to be a bedroom with a dresser and a window. The girl is focused on the book and seems to be enjoying her reading time. | A baby wears glasses as he plays | In the video, we see a young child sitting on a bed, wearing glasses and holding a book. The child appears to be engaged in reading or looking at the book, possibly interacting with someone off-camera. The room has a cozy, lived-in feel with various pieces of furniture and personal items, suggesting a home environment. The child's focus and the interaction with the book indicate a moment of learning or playtime. | In the video, a baby is sitting on the floor and playing with a Wii remote. (or a Wii controller). The baby is holding the remote in their hands and appears to be enjoying themselves. |
| • Incorrect description<br>• 4 bit loses all inference accuracy<br>• No correlation to original video | • Gives an accurate description of the subject and environment appearance<br>• Does not mention emotional undertones | • Very concise<br>• Does not elaborate on either the subject or the action, and does not correctly describe the activity as well | • Most descriptive output with correct description and adds emotional aspect (cozy, lived-in, moment of learning)<br>• "interacting with someone off-camera" is speculative and cannot be confirmed | • Semi-correct description of the events<br>• Unable to identify the smaller object such as the book |

# Comparison of outputs (8 bit quantized)

| Qwen | Phi3.5 | Idefics2 | LlavaNext | Llava |
|---|---|---|---|---|
| The video appears to show a person lying on a bed, possibly asleep or resting. The background includes a bed and some furniture, suggesting a bedroom setting. The person is wearing glasses and seems to be in a relaxed state. The overall atmosphere of the video is calm and serene | In this video, a young girl is sitting on a bed and reading a book. She is wearing a green shirt and pink pants. The room appears to be a bedroom with a dresser and a window. The girl is holding the book with both hands and seems to be engrossed in the story. She occasionally looks up and around the room, but mostly remains focused on the book | A baby is sitting on a bed with a book | In the video, we see a young child sitting on the floor, wearing glasses, and holding a book. The child appears to be engaged in reading or looking at the book, possibly trying to figure out how to read it. The child's focus is on the book, and they seem to be in a calm and relaxed state, possibly enjoying the activity of reading. The room has a cozy and homely atmosphere, with a couch and a chair visible in the background, suggesting a comfortable indoor setting. The child's posture and the way they are holding the book indicate that they are act | In the video, a baby is sitting on a bed and playing with a Wii remote. The baby is holding the remote in their hands and appears to be enjoying themselves. |
| • Semi-correct description- improvement from 4 bit output<br>• Does not highlight the undertaken action/activity- seems to be more of a frame analysis | • An accurate description of the subject and the environment<br>• Missing emotional undertones | • Concise semi-correct description<br>• Improvement from the 4 bit output as the smaller subject (book) is identified and described in the output | • Similar to the previous output<br>• Accurately describes the environment, subject, activity and emotional undertone<br>• Abruptly cuts off in the middle however | • Misidentifies the book as a Wii remote<br>• Describes the emotional state of the subject |

# Comparison of outputs (32 bit quantized)

| Qwen | Phi3.5 | Idefics2 | LlavaNext | Llava |
|------|--------|----------|-----------|-------|
| In the video, a young child is sitting on a bed wearing glasses. The child is holding a book and appears to be reading it. The background shows a bedroom with a bed, a pillow, and some furniture. The child seems to be enjoying the book and is focused on reading | In this video, a young girl is sitting on a bed and reading a book. She is wearing a green shirt and pink pants. The room appears to be a bedroom with a dresser and a window. The girl is focused on the book and seems to be enjoying her reading time | A baby is sitting on a bed with a book | In the video, we see a young child sitting on a bed, wearing glasses and holding a book. The child appears to be engaged in reading or looking at the book, possibly trying to figure out how to read it. The child's focus is on the book, and they seem to be in a comfortable and relaxed environment, possibly at home. The room has a cozy and lived-in feel, with various items like a bed, a chair, and a bedside table visible in the background. The child's expression suggests curiosity and interest in the book, and the scene captures a moment of learning | In the video, a baby is sitting on a bed and playing with a Wii remote |
| • Correctly describes the video<br>• Significant improvement over the previous quantized version- model is able to correctly identify and describe the video | • A shorter, albeit correct, description than the quantized versions | • The short descriptions do not improve with the model size<br>• Similar to the 8 bit quantized output | • Similar to the previous outputs<br>• Most accurate description out of all the given outputs | • The largest model also does not seem to recognize the book correctly<br>• Problem with the base model rather than the quantization |

Tech.

# Analysis

- Idefics gives really short and concise outputs
  - Outputs improve as the quantizations are removed
- Qwen has the highest performance degradation as the smaller quantizations are applied
  - The original model provides an accurate description of the video, the 8 bit quantized model loses some information but provides a correct description and the 4 bit quantized model gives a completely incorrect output
- LlavaNext provides the most consistent and accurate descriptions

Georgia Tech.

# Pruning

# Pruning Results

| Benchmark | Phi3.5 | Phi 3.5 Pruned | Fine-tuning time (s) | Inference Time (s) | Inference Time Pruned (s) |
|-----------|--------|----------------|----------------------|--------------------|---------------------------|
| DocVQA | 0.94 | 0.879 | 1150.09 | 108.92 | 138.51 |
| VQAv2 | 0.708 | 0.920 | 3087.37 | 256.88 | 106.40 |
| ScienceQA | 0.842 | 0.896 | 2991.36 | 80.00 | 80.00 |

| Model Type | Model Size |
|------------|------------|
| Phi 3.5 | 16.58GB |
| Phi 3.5 Pruned | 8.26GB |

Size of Finetuning Dataset: 2000
Train/Test: 4:1
Epoch: 3
Learning Rate: 4e-5

Georgia Tech

# Pruning Results for Cross Benchmark

|  | Finetune on ScienceQA | Finetune on VAQ2 | Finetune on DocVQA |
|---|---|---|---|
| ScienceQA | 0.896 | 0.842 | 0.536 |
| VQA2 | 0.72 | 0.92 | 0.74 |
| DocVQA | 0.63 | 0.79 | 0.879 |

Size of Finetuning Dataset: 1200
Train/Test: 4:1
Epoch: 3
Learning Rate: 4e-5

Georgia Tech

# Pruning Result Analysis

When doing prune on Phi 3.5, I pruned the decoder part of the model. The decoders are 32 repeating structures of layers, I delete half of them and then do fine-tuning on the dataset.

I kept the feature extractor submodel to keep the ability of extracting the features. By deleting half of the decoder part and fine-tuning, the ability of interpreting the features is mitigated.
Looking at pruning results of fine-tuning and evaluating on the same benchmark, we find the performance is good. The performance of VQA2 and ScienceQA is even better because of the fine-tuning. The performance of DocVQA gets worse but not by a lot.

When doing cross benchmark evaluation, the performace on VQA2 is more stable as it is not MCQ benchmark.

Considering the size of fine-tuning dataset, the overall performance of pruning is acceptable.

Georgia Tech

# Low Rank Factorization

# Why LoRA does not work for Model Compression?

- **LoRA (Low-Rank Adaptation)** is used to fine-tune large language models by adding low-rank matrices to the existing weights.

- LoRA adds trainable parameters without reducing the base model size
  - Hence, unsuitable for Model Size Reduction.

- **What to do?**
  - Use **Low-Rank Factorization**: the core concept of LoRA
  - Decompose weight matrices directly into smaller low-rank components.
  - Reduces number of model parameters, reducing the memory footprint of the model.

Georgia Tech®

# Why Factoring by Retained Variance is sensible?

- A fixed rank across layers may miss layer-specific importance.
    - We would end up removing more information from information rich (high variance) layers, leading to poor model performance.

- Retained variance preserves important information while optimizing compression.

- Last layers (QWen2-VL-2B-Instruct VLM) are not factorized to retain critical features for output generation.
    - Factorizing last layers leads to garbage characters in the output.

Georgia Tech®

# Low Rank Factorization Results and Comparisions

| Benchmark | Qwen2-VL-2B Base Model | | 80% Retained Variance | | 50% Retained Variance | | Pruning | | 8 Bit Quantized Model | |
|---|---|---|---|---|---|---|---|---|---|---|
| | BERT | Time | BERT | Time | BERT | Time | BERT | Time | BERT | Time |
| DocVQA | 0.769 | 54.22s | 0.681 | 70.5s | 0.622 | 94.37s | 0.724 | 53.16s | 0.805 | 227s |
| VQAv2 | 0.865 | 40.69s | 0.860 | 54.67s | 0.810 | 44.32s | 0.854 | 28.39s | 0.645 | 172s |
| Flickr30k | 0.865 | 40.71s | 0.845 | 51.77s | 0.700 | 51.04s | 0.854 | 29.34s | 0.908 | 187s |

| Model Type | Model Size |
|---|---|
| Base | 8.23GB |
| 80% Retained Variance | 6.00GB |
| 50% Retained Variance | 3.93GB |
| Pruned | 5.79GB |
| Quantized 8 bit | 1.65GB |

Replaced ScienceQA benchmark with Flickr30k as ScienceQA has A/B/C/D as the output. And finetuning is not working well on succinct outputs for Qwen2.

Also, Flickr30k is an image-captioning benchmark, which is more aligned with the task at hand.

Georgia Tech

# Finetuning Information

- Qwen2-VL-2B is originally trained on a diverse corpus of Image to Text dataset, Image-Text pairs, Visual QA, Video dialogues etc[1].

- After applying pruning and Low Rank Factorization, we observed that we're losing great amount of performance on the benchmarks. Hence, we finetune the compressed model on the benchmarks.

- We finetune over 1000 examples of each of the mentioned benchmark and then use 200 unseen examples for testing.

- For each benchmark, we finetune baseline, pruned model, and low-rank models, to do a fair comparison of performance.

[1] Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., Fan, Y., Dang, K., Du, M., Ren, X., Men, R., Liu, D., Zhou, C., Zhou, J., & Lin, J. (2024). Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. *arXiv preprint arXiv.2409.12191*.

Georgia Tech®

# Analyzing VQA-V2



Query: Where is the television?

Expected: on dresser

80% Var: above window

50% Var: no  (model in general answers "no" a lot; Garbage Output)

Pruned: above wall

Baseline: above window

Analysis: We can see that the 80% variance version retains the wrong answer of baseline, whereas the pruned model misinterprets the dresser as wall.

# Analyzing VQA-V2



Query: What color lights are here?

Expected: red

80% Var: yellow

50% Var: white

Pruned: red

Baseline: red

Analysis: Pruned model is able to retain the correct answer of the Baseline whereas low-rank models fumble. This indicates that low-rank models are not able to capture the intricate features of the image well.

# Analyzing Flickr30k



Image Captioning Task:

Expected: 'Two little girls playing with sand on the beach.'

80% Var: 'Two girls playing in the sand.'

50% Var: 'A man in a white shirt is standing on a large table.' (Garbage Output)

Pruned: 'Two girls are sitting on the ground.'

Baseline: 'Two young girls are playing with sand.'

Analysis: In this case, the 80% variance model retains the correct output, whereas the pruned model is also able to capture the major detail about two girls sitting on the ground. At 50% variance, we just see garbage output.

Georgia Tech.

# Analyzing Flickr30k



Image Captioning Task:

Expected: 'A group of gymnasts performing in the Basketball court.'

80% Var: 'A group of people in a gym playing a game.'

50% Var: 'A man in a white shirt is standing on a large table.' (Garbage Output)

Pruned: 'A group of people in a basketball field.'

Baseline: 'A group of people in white dance on a basketball court.'

Analysis: As the dataset we finetune on is limited, none of the models are correctly able to detect the activity as "gymnastics". But the pruned and baseline model give semantically correct output. Whereas 80% var model confuses the court with the gym. This further signifies that low-rank models' ability to capture visual features is impacted severely.

# Analyzing Flickr30k



Image Captioning Task:

Expected: 'A man in khakis pants is lying on concrete and looking into the sky.'

80% Var: 'A man is sleeping on a bench in front of a building.'

50% Var: 'A man in a white shirt is standing on a large table.' (Garbage Output)

Pruned: 'One person sitting on the steps of a building.'

Baseline: 'A man in a beige jacket is sitting on a ledge.'

Analysis:
80% variation captures the scene reasonably well but lacks detail, while pruned and baseline models provide semantically correct but generalized outputs, and 50% variance model shows severe degradation with nonsensical output.

# Analyzing Low-rank Video Demo

- 80% Retained Variance: 'A baby in a highchair is being held by a baby.

- 50% Retained Variance: 'A man in a red shirt is sitting on a bench.'

- Pruned: 'A young boy is playing with books.'

- Baseline: 'In the video, a young child is sitting on a bed wearing glasses. The child is holding a book and appears to be reading it. The background shows a bedroom with a bed, a pillow, and some furniture. The child seems to be enjoying the book and is focused on reading.'

- 8-bit: 'The video appears to show a person lying down on a bed. The person is wearing glasses and seems to be resting or sleeping. The background includes some furniture, such as a chair and a bed, and there is a visible logo or watermark in the bottom right corner of the image.'

- **Note**: Compressed model are succinct as they are finetuned on image-captioning dataset due to resource constraints and hence are not capturing various aspects of Video.

# Analysis of Demo Descriptions

- Baseline Model: The baseline model provides the most detailed and coherent description, capturing the child's actions, appearance, and surroundings accurately.

- Compressed Models:
  - As the compressed models are finetuned of Flickr30k, which contains one-liner caption for images, the video descriptions tend to be brief.
  - At 80% variance, the model is able to get some details out of the video such as there is a baby, and the baby is holding something. But the model is not able to frame a proper and coherent sentence.
  - At 50% variance, the model is outputting complete gibberish.
  - The pruned model gives a fairly accurate description but misses the details of the scene.
  - 8-bit model interprets the video incorrectly. And thinks that it is a person lying down. It also focuses on the watermark, which is accurate, but not very informative for the task at hand.

- We can see that the pruning retains the most performance out of all these compression techniques. Moreover, the output of the quantized model is promising. The model might perform better with Quantization if we have access to full training data and if we do a Quantization Aware Training (QAT).

Georgia Tech

# Analysis

- Pruning out-performs low rank factorization in almost all aspects (time, model-size, output quality).

- At 50% retained variance, qwen2-vl starts to become repetitive, and does not provide meaningful output.

- Factorization is <u>not leading to a significant improvement in the inference time</u>.
    - We're replacing nn.Linear layer in Transformers with 2 low-rank nn.Linear layers, wrapped in LowRankLinear layer.
    - Hypothesis: PyTorch is not able to optimize LowRankLinear to the same extent as nn.Linear. Also, having 2 internal linear layers would have higher overhead cost.

- Quantization does a lot better at improving the memory footprint but does not improve the inference speed.
    - Hypothesis: Quantization Overhead. Unoptimized hardware for the quantized parameters.

Georgia Tech.

# Interesting Piece of Code

```python
class CustomQwen2VL(Qwen2VL):
    def __init__(self, quantization_mode, model, tokenizer, processor, name=None):
        self.name = 'qwen2-custom'

        if quantization_mode is not None:
            print('WARNING: CustomQwen2VL ignores the quantization mode. Passed value: ' + str(quantization_mode))

    def get_model_name(self):
        return self.name


    @staticmethod
    def from_low_rank_path(metadata_path, pt_model_path, model_name=None):
        from low_rank import patch_model_using_metadata, read_metadata
        tokenizer = AutoTokenizer.from_pretrained(
            QWEN2_MODEL_NAME,
            trust_remote_code=True
        )
        processor = AutoProcessor.from_pretrained(QWEN2_MODEL_NAME)

        cache_dir = setup_cache_dir()
        config = Qwen2VLConfig.from_pretrained(
            QWEN2_MODEL_NAME,
            trust_remote_code=True
        )
        model = Qwen2VLForConditionalGeneration.from_pretrained(
            QWEN2_MODEL_NAME,
            config=config,
            device_map="auto",
            trust_remote_code=True,
            cache_dir=cache_dir
        )

        model = patch_model_using_metadata(model, read_metadata(metadata_path), pt_model_path)

        # Return an instance of CustomQwen2VL
        return CustomQwen2VL(
            quantization_mode=None,
            model=model,
            tokenizer=tokenizer,
            processor=processor,
            name=model_name
        )
```

```python
def save_metadata(path):
    with open(path, "w") as f:
        json.dump(metadata, f, indent=2)

def read_metadata(path):
    with open(path, "r") as f:
        metadata = json.load(f)
    return metadata

def clear_metadata():
    global metadata
    metadata = {}

def patch_model_using_metadata(model, metadata, pt_model_path=None):
    """
    Replace linear layers with low-rank layers using metadata.
    """
    for name, layer_info in metadata.items():
        # Locate the module to replace
        parent_module = model
        components = name.split('.')
        for comp in components[:-1]:
            parent_module = getattr(parent_module, comp)

        original_layer = getattr(parent_module, components[-1])
        if isinstance(original_layer, torch.nn.Linear):
            # Replace with a LowRankLinear equivalent using metadata
            in_features = layer_info["low_rank_shape"]["linear1"][1]
            out_features = layer_info["low_rank_shape"]["linear2"][0]
            rank = layer_info["low_rank_shape"]["linear1"][0]
            low_rank_layer = LowRankLinear(in_features, out_features, rank)
            setattr(parent_module, components[-1], low_rank_layer)

    if pt_model_path is not None:
        loc = "cuda" if torch.cuda.is_available() else "cpu"
        state_dict = torch.load(pt_model_path, map_location=loc)
        model.load_state_dict(state_dict)
        print('Loaded model weight on ' + loc)

    return model
```

As Qwen2VL uses specific config classes, it is challenging to save and reload the transformer model.

Solution: We save the patched weights and metadata and upon first load, we patch the model to use the compressed architecture.

```python
metadata_json = os.path.join(output_dir, "metadata.json")
save_metadata(metadata_json)

print("Saved metadata for low rank factorization")

pytorch_model_path = os.path.join(output_dir, "pytorch_model.pt")
torch.save(model.state_dict(), pytorch_model_path)
print("Saved low rank weights")
```
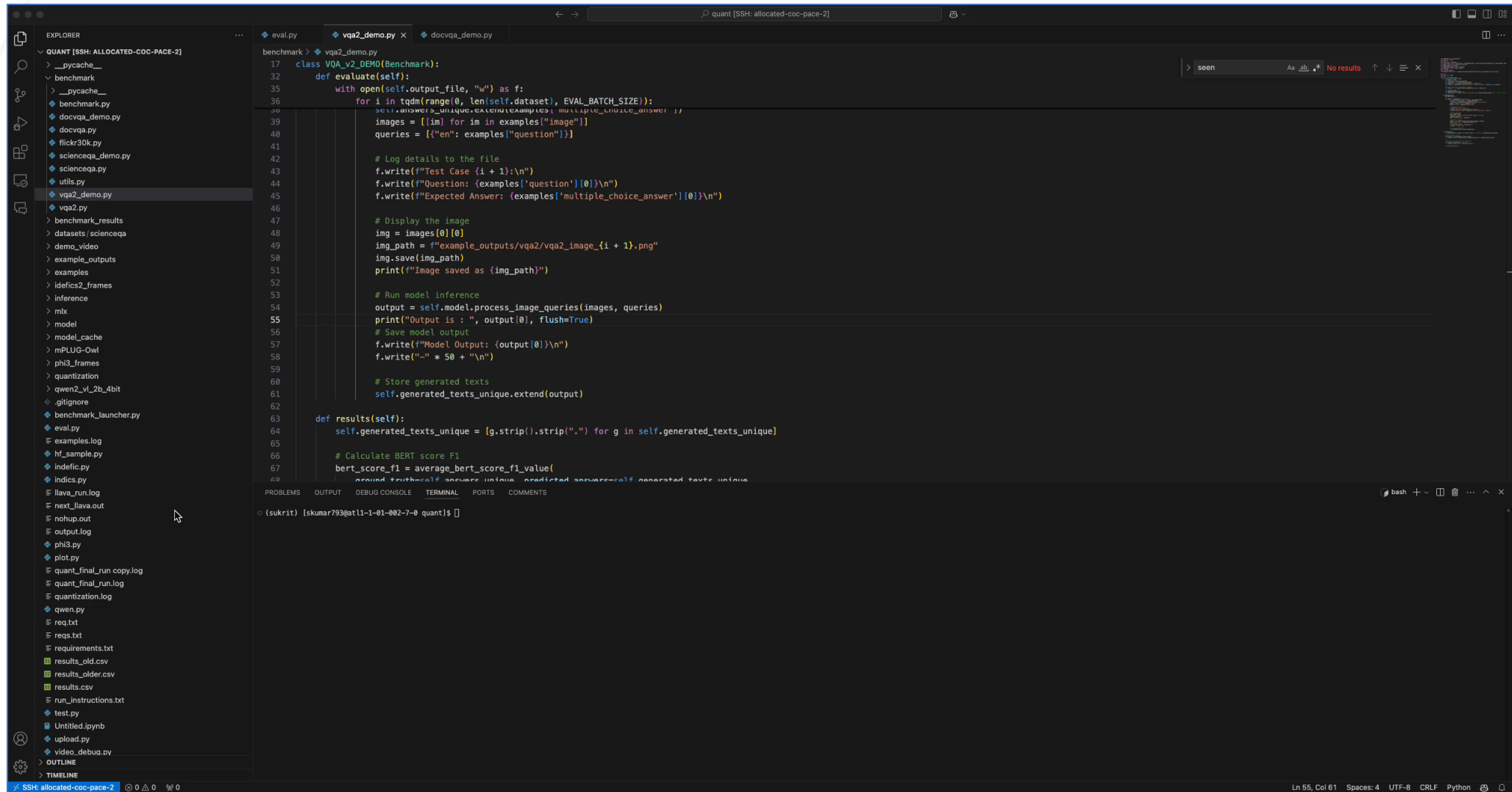
Georgia Tech

# Sample Outputs : The BAD

# Demo: Eval.py : [Video Link](#)

# Demo: Video Inference : [Video Link](Video Link)

# Edge specific optimizations

- These are the additional edge specific optimizations we did in-addition to quantization, pruning and low-rank weight decomposition
- We used Eager attention implementation instead of FlashAttentionv2 due to lack of support for it on edge-devices like laptops
- We additionally compiled the PyTorch models using torch.compile to allow for faster execution
- In addition to these optimization, we allowed for auto-offloading of weights to the CPU/Memory from the GPU.

Georgia Tech.

# Lessons, Learnt, Takeaways, Future Improvement

- We performed a comprehensive analysis of <u>3 different model compression methods : Quantization, Pruning and Low-rank factorization</u> for model layers across <u>4 different benchmark</u> datasets across <u>5 different models</u>.
- We provide a comprehensive plug and play system with easy extension, structured evaluation system and handle model specific things without messing up the abstraction
- We learnt that Quantization provides the best memory reduction with minimal performance loss. Pruning leads to slight inference speed gains, where as low-rank weight factorization shows a slight reduction with marginal degradation in output quality.
- We faced multiple technical challenges dealing with VLM specific issues, overhead of quantizing models on the fly, additional model specific issues like with Qwen and mPLug. Plus, we saw very limited hardware support for pruned and sparse weight matrix.
- Future improvements could be optimizing for inference time, using hardware specific abstraction to optimize inference speed, combining multiple techniques like Low-rank weights with quantization to further improve performance and model size.

# README and Code pointers

- Link to README.md:
  https://github.com/ksukrit/BDA_project/blob/master/README.md
- Pruning samples:
  - https://github.com/ksukrit/BDA_project/blob/master/prune/phi_prune.py
  - https://github.com/ksukrit/BDA_project/blob/master/prune/qwen2_prune.py
- Low Rank Model Patching:
  - https://github.com/ksukrit/BDA_project/blob/master/low_rank/low_rank.py#L95
- Plug & Play:
  - Benchmark example (all benchmarks follow the same interface):
    https://github.com/ksukrit/BDA_project/blob/master/benchmark/docvqa.py
  - Model example (all models follow the same interface):
    https://github.com/ksukrit/BDA_project/blob/master/model/llavanext.py

Georgia Tech.

# Open-Source Tools Used

- [HuggingFace model Hub](#)

- [HuggingFace Transformers](#)

- [HuggingFace Accelerate](#)

- [BitsandBytes](#)

- [FlashAttention2](#)

- [av](#)

# Bibliography/References

1. K. Alizadeh, I. Mirzadeh, D. Belenko, S. K. Khatamifard, M. Cho, C. C. D. Mundo, M. Rastegari, and M. Farajtabar. Llm in a flash: Efficient large language model inference with limited memory. In ACL, 2024.

2. A. Annavajjala, A. Khare, A. Agrawal, I. Fedorov, H. Latapie, M. Lee, and A. Tumanov. $D\epsilon ps$ : $Delayed\epsilon - shrinking f$ or $f$ asteronce $- f$ or $-$ alltraining, 2024.4

3. H. Bai, W. Lin, M. Zhang, A. Wang, and Y. Zhao. Multiple description video coding based on human visual system characteristics. IEEE Transactions on Circuits and Systems for Video Technology, 24(8):1390–1394, 2014.

4. K.-H. Le Minh, K.-H. Le, and Q. Le-Trung. Dlase: A light-weight framework supporting deep learning for edge devices. In 2020 4th International Conference on Recent Advances in Signal Processing, Telecommunications Computing (SigTelCom), pages 103–108, 2020.

5. M. Sahni, S. Varshini, A. Khare, and A. Tumanov. Compofa: Compound once-for-all networks for faster multi-platform deployment, 2021.

6. Y. Tang, J. Bi, S. Xu, L. Song, S. Liang, T. Wang, D. Zhang, J. An, J. Lin, R. Zhu, A. Vosoughi, C. Huang, Z. Zhang, P. Liu, M. Feng, F. Zheng, J. Zhang, P. Luo, J. Luo, and C. Xu. Video understanding with large language models: A survey, 2024.

7. P. Warden and D. Situnayake. Tinyml: Machine learning with tensorflow lite on arduino and ultra-low-power microcontrollers. O'Reilly Media, 2019.

8. X. Zhu, J. Li, Y. Liu, C. Ma, and W. Wang. A survey on model compression for large language models, 2024.

9. Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., Fan, Y., Dang, K., Du, M., Ren, X., Men, R., Liu, D., Zhou, C., Zhou, J., & Lin, J. (2024). Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. arXiv preprint arXiv:2409.12191.

Georgia Tech