

ECE 1512 Project B Part 2

Yingshun Lu 1006029049

Minghao Ma 1010800536

1. Summary Document on VILA[1]

1.1 Key Points

Focus: VILA investigates pre-training strategies for VLMs that better align the visual and textual modalities.

Objective: To impart vision capability to the LLMs while retaining functionality typical of text-only LLMs.

Architecture: Auto-regressive architecture, taking the visual tokens as equivalent to textual tokens makes the model flexible and unified for taking multi-modal inputs.

Applications: Excels in vision-language tasks like VQA, caption generation, and multi-image reasoning.

1.2 Technical Contributions

Pre-training Strategies:

That demonstrates fine-tuning in visual pre-training of LLMs as highly critical for deep embedding alignment and in-context learning.

Interleaving image-text datasets enjoy the benefits of better alignment and minimized degradation in the text-only capability compared to plain image-text pairs.

Data Blending:

It reintroduces text-only instruction data during supervised fine-tuning to regain degraded text-only capabilities and improve VLM task performance.

This paper proposes a mixture of interleaved datasets with image-text pairs to enhance the diversity and also the accuracy of downstream tasks.

Performance:

Consistently outperforms state-of-the-art models like LLaVA-1.5 on several benchmarks.

Demonstrates improved robust multiple-image reasoning and better world knowledge retention.

Efficiency:

Employs scalable approaches employing resolution adjustments and lightweight projection layers that result in improved performance with better cost trade-offs.

1.3 Areas for Improvement

Scaling:

Very limited pre-training dataset of 50M images compared to others' work in billion-scale datasets.

Further improvements may be possible by expanding the training dataset.

Instruction Dataset Quality:

Though successful, there is room for improving the diversity and quality of the prompts that constitute the instruction-tuning dataset.

Edge Deployment:

While it might be deployable with systems such as Jetson Orin, the full size of this model alone can already be a challenging prospect in environments starved for computational resources.

Generalization:

While VILA still seems to keep up with competitive text-only performance on this end, smaller models significantly degrade in performance, indicating room for improvement in preserving text-only abilities during pre-training.

2. Efficiency Optimization Report: Dynamic Token Pruning in Transformer Models

2.1. Introduction

Transformer models are very expensive due to heavy computations of self-attention and feed-forward layers at the token level. Reducing this while keeping the accuracy intact has become one of the most important factors that will help deploy efficient models into production. Thus, finding out bottlenecks is followed by optimization identification and evaluation to make up the core part of efficiency enhancement in the course of this report with the aid of Dynamic Token Pruning on the model in question.

2.2 Efficiency Bottleneck Identification

2.2.1 Profiling Results

We have profiled a baseline transformer model in order to get an idea about which are the main bottlenecks concerning efficiency. Results are summarised below:

Major Bottlenecks:

Self-Attention Layers: Dominant contributor of FLOPs and computational time.

Feedforward Layers (Linear): Responsible for secondary computational costs.

Baseline Model FLOPs and Time:

FLOPs: **12.908G**

CUDA Time: **12.490ms**

2.2.2 Conclusion

This inefficiency primarily comes from redundant computations over the input tokens that make minor contributions to the output of the whole model. In our attempt to improve it, dynamic token pruning focuses on better usage of tokens.

2.3 Optimization Method: Dynamic Token Pruning

2.3.1 Overview

Dynamic Token Pruning is based on adaptively removing irrelevant tokens computed by saliency from computation. This reduces the input sequence length dynamically during the inference process, saving considerably on FLOPs and time

2.3.2 Algorithm

1. **Compute Saliency Scores:** Each token's importance is estimated using its L2 norm.
2. **Generate Token Mask:** Saliency scores less than some threshold are pruned off the tokens.

3. **Layer-Wise Pruning:** Pass the remaining tokens to the next transformer block.
4. **Final Classification:** Pool the outputs of all active tokens.

2.3.3 Pseudo Code

```
class PrunedTransformerEncoder(nn.Module):
    def forward(self, src):
        keep_tokens = torch.ones(src.shape[:2], device=src.device).bool()
        for i, layer in enumerate(self.layers):
            saliency = self.token_pruning.compute_saliency(src)
            keep_tokens = keep_tokens & (saliency > self.token_pruning.saliency_threshold)
            src = layer(src, keep_tokens=keep_tokens)
        return src
```

2.3.4 Implementation Details

Saliency Score: Computed as the L2 norm of each token vector.

Token Mask: Dynamically updated across layers, preserving only the most relevant tokens.

2.4 Results and Evaluation

2.4.1 Accuracy

The optimized model maintains the same accuracy as the baseline:

Baseline Model Accuracy: **52.5%**

Pruned Model Accuracy: **52.5%**

2.4.2 FLOPs and Time Comparison

Model	FLOPs	CUDA Time	FLOPs Reduction	Time Reduction
Baseline	12.908G	12.490ms	-	-
Pruned	8.874G	9.968ms	31.2%	20.2%

The pruning method significantly reduces FLOPs and computational time while preserving accuracy.

2.4.3 Profiling Results

Baseline Model

- FLOPs: **12.908G**
- CUDA Memory Usage: **72.02MB**
- CUDA Time: **12.490ms**

Pruned Model

- FLOPs: **8.874G**
- CUDA Memory Usage: **72.02MB**
- CUDA Time: **9.968ms**

2.4.4 Efficiency Analysis

Dynamic Token Pruning successfully reduces computation while maintaining the accuracy of the model. Profiling indicates that FLOPs and execution time reductions are most prominent in the self-attention layers.

2.5 Discussion and Considerations

2.5.1 Strengths

Efficiency Gains:

Much smaller FLOPs and execution time, over 31.2% reduction in FLOPs.

Preserves the accuracy of a binary classification task.

Scalability:

The pruning is dynamic, layer-wise for the integration of large-scale models.

2.5.2 Limitations

Static Threshold: A static setting of the saliency threshold may not generalize so well for a wide range of datasets.

Evaluation on Toy Data: While the accuracy may be fine on simulated tasks, this may need further validation on real-world datasets.

2.6. Conclusion

This report now presents the possibility of Dynamic Token Pruning in optimizing transformer efficiency. The dynamic pruning of the length of a sequence using saliency-based pruning produced:

- Significant reductions both in FLOPs and in execution time.
- Maintenance of accuracy on a binary classification task.

Future work will be conducted on the validation of this approach with larger datasets, as well as exploring adaptive thresholds for saliency computation.

Reference:

[1] Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023.

Appendix:

Github links:

https://github.com/ADglory/ECE1512_2024F_ProjectB_PartB_Repo_YingshunLu_Minghao-Ma

Codes:

```
11 1. Introduce necessary dependencies
12 """
13
14 !pip install fvcore
15
16 import torch
17 import torch.nn as nn
18 import torch.nn.functional as F
19 from fvcore.nn import FlopCountAnalysis, flop_count_table
20 import torch.profiler as profiler
21 from torch.utils.data import DataLoader, TensorDataset
22 from copy import deepcopy
23
24 """2. Token Saliency computing module"""
25
26 class TokenSaliency(nn.Module):
27     """
28     Compute saliency scores for visual tokens based on their contribution.
29     """
30     def __init__(self, method="norm"):
31         super(TokenSaliency, self).__init__()
32         self.method = method
33
34     def forward(self, tokens):
35         """
36         Args:
37             tokens: Tensor of shape (B, N, D), where
38                     B = Batch size,
39                     N = Number of tokens,
40                     D = Dimension of each token.
41         Returns:
42             saliency_scores: Tensor of shape (B, N), saliency scores for each token.
43         """
44         if self.method == "norm":
45             saliency_scores = tokens.norm(dim=-1) # Use L2 norm
46         else:
47             raise ValueError(f"Unsupported method: {self.method}")
48
49         return saliency_scores
50
51 """3. Adaptive Token pruning module"""
52
53 class AdaptiveTokenPruning(nn.Module):
54     def __init__(self, saliency_threshold=0.5):
55         super(AdaptiveTokenPruning, self).__init__()
56         self.saliency_threshold = saliency_threshold
57
58     def forward(self, x):
59         """
60         Compute token saliency and generate a pruning mask.
61         """
62         saliency_scores = self.compute_saliency(x)
63         keep_tokens = saliency_scores > self.saliency_threshold
64         return keep_tokens
65
```

```

66  ✓ def compute_saliency(self, x):
67      """
68      Compute saliency scores (e.g., L2 norm across embedding dimensions).
69      """
70      saliency_scores = x.norm(p=2, dim=-1) # Shape: (batch_size, seq_len)
71      return saliency_scores
72
73      """4. Pruned Transformer Encoder"""
74
75  ✓ class PrunedTransformerEncoder(nn.Module):
76      """
77      Transformer encoder with token pruning capability.
78      """
79      def __init__(self, encoder_layer, num_layers, saliency_threshold=0.5):
80          super().__init__()
81          self.layers = nn.ModuleList([deepcopy(encoder_layer) for _ in range(num_layers)])
82          self.token_pruning = AdaptiveTokenPruning(saliency_threshold=saliency_threshold)
83
84  ✓ def forward(self, src):
85      """
86      Forward pass with token pruning.
87      Args:
88          src: Input tensor of shape (batch_size, seq_len, d_model).
89      Returns:
90          Output tensor after pruning.
91      """
92      batch_size, seq_len, d_model = src.shape
93      keep_tokens = torch.ones((batch_size, seq_len), device=src.device).bool() # Initialize with all True
94
95      for i, layer in enumerate(self.layers):
96          # Calculate saliency scores
97          saliency = self.token_pruning.compute_saliency(src)
98
99          # Update keep_tokens
100         new_keep_tokens = (saliency > self.token_pruning.saliency_threshold)
101         keep_tokens = keep_tokens & new_keep_tokens # Retain the accumulated crop state
102
103         # Dynamically crop the input tensor
104         pruned_src = []
105         pruned_keep_tokens = []
106
107         for batch_idx in range(batch_size):
108             active_token_indices = keep_tokens[batch_idx].nonzero(as_tuple=True)[0]
109             pruned_src.append(src[batch_idx, active_token_indices])
110             pruned_keep_tokens.append(keep_tokens[batch_idx, active_token_indices])
111
112         # Update src and keep_tokens with the trimmed tensor
113         src = torch.nn.utils.rnn.pad_sequence(pruned_src, batch_first=True)
114         keep_tokens = torch.nn.utils.rnn.pad_sequence(pruned_keep_tokens, batch_first=True)
115
116         # Print debugging information
117         print(f"Layer {i}: Active tokens per batch = {[len(t) for t in pruned_src]}")
118
119         # Pass the clipped tensor to the next layer
120         src = layer(src)

```



```

121
122         return src
123
124     """5. FLOPs evaluation tool"""
125
126     from fvcore.nn import FlopCountAnalysis, flop_count_table
127
128     def calculate_dynamic_flops(model, x, keep_tokens):
129         """
130         Calculate FLOPs dynamically based on active tokens.
131         Args:
132             model: The pruned Transformer model.
133             x: Input tensor of shape (batch_size, seq_len, d_model).
134             keep_tokens: Boolean tensor indicating active tokens for the pruned model.
135         """
136         # Get the maximum number of active tokens
137         active_tokens = keep_tokens.sum(dim=1).max().item()
138         x = x[:, :active_tokens, :] # Crop to active Token
139         flops = FlopCountAnalysis(model, x)
140         print(flop_count_table(flops))
141
142     """6. Memory usage evaluation tool"""
143
144     def profile_memory_and_time_safe(model, input_tensor):
145         """
146         Profile memory and time for the given model and input.
147         Args:
148             model: PyTorch model to profile.
149             input_tensor: Tensor input to pass through the model.
150         """
151         try:
152             with torch.profiler.profile(
153                 activities=[
154                     torch.profiler.ProfilerActivity.CPU,
155                     torch.profiler.ProfilerActivity.CUDA,
156                 ],
157                 record_shapes=True,
158                 profile_memory=True,
159                 with_stack=False, # Disable stack tracing to reduce possible conflicts
160             ) as prof:
161                 model(input_tensor) # Perform model forward propagation
162                 print(prof.key_averages().table(sort_by="cuda_time_total", row_limit=10))
163             except RuntimeError as e:
164                 print(f"Profiler failed: {e}")
165
166     """7. Model accuracy evaluation"""

```

```

167
168 ✓ def evaluate_model_accuracy(model, train_data, train_labels, test_data, test_labels):
169     """
170     Train and evaluate model accuracy on a toy dataset.
171     Args:
172         model: PyTorch model to evaluate.
173         train_data, train_labels, test_data, test_labels: Dataset tensors.
174     """
175     # Make sure the shape of the label is 1D
176     train_labels = train_labels.squeeze()
177     test_labels = test_labels.squeeze()
178
179     model.train()
180
181     # Dataset and DataLoader
182     train_dataset = TensorDataset(train_data, train_labels)
183     test_dataset = TensorDataset(test_data, test_labels)
184     train_loader = DataLoader(train_dataset, batch_size=16, shuffle=True)
185     test_loader = DataLoader(test_dataset, batch_size=16)
186
187     # Optimizer and Loss
188     optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
189     loss_fn = nn.CrossEntropyLoss()
190
191     # Training loop
192     for epoch in range(5):
193         for inputs, labels in train_loader:
194             inputs, labels = inputs.cuda(), labels.cuda()
195             optimizer.zero_grad()
196             outputs = model(inputs)
197             outputs = outputs.mean(dim=1)
198             loss = loss_fn(outputs, labels)
199             loss.backward()
200             optimizer.step()
201
202     # Evaluation loop
203     model.eval()
204     correct, total = 0, 0
205     with torch.no_grad():
206         for inputs, labels in test_loader:
207             inputs, labels = inputs.cuda(), labels.cuda()
208             outputs = model(inputs)
209             outputs = outputs.mean(dim=1)
210             _, predicted = torch.max(outputs, 1)
211             total += labels.size(0)
212             correct += (predicted == labels).sum().item()

```

```

213
214     accuracy = 100 * correct / total
215     print(f"Accuracy: {accuracy:.2f}%")
216
217     from fvcore.nn import FlopCountAnalysis, flop_count_table
218
219     def evaluate_pruned_model(baseline_model, pruned_model, test_data):
220         """
221         Compare baseline and pruned models in terms of FLOPs and active token efficiency.
222         Args:
223             baseline_model: The baseline Transformer model.
224             pruned_model: The pruned Transformer model.
225             test_data: Sample input tensor for efficiency evaluation.
226         """
227         print("=== Baseline Model Efficiency ===")
228         flops_baseline = FlopCountAnalysis(baseline_model, test_data)
229         print(flop_count_table(flops_baseline))
230
231         print("\n=== Pruned Model Efficiency ===")
232         # Assuming the PrunedTransformerEncoder dynamically prunes tokens
233         with torch.no_grad():
234             pruned_outputs = pruned_model[0](test_data) # Get the intermediate result of PrunedTransformer
235             active_tokens = pruned_outputs.shape[1] # The number of valid tokens remaining
236             flops_pruned = FlopCountAnalysis(pruned_model, test_data[:, :active_tokens, :])
237             print(flop_count_table(flops_pruned))
238
239     """8. Prepare the data set"""
240
241     def prepare_data():
242         """
243         Prepare simulated toy dataset for training and testing.
244         Returns:
245             train_data, train_labels, test_data, test_labels
246         """
247         train_data = torch.rand(1000, 128, 512).cuda() # 1000 samples, 128 tokens, 512 dimensions
248         train_labels = torch.randint(0, 2, (1000,), dtype=torch.long).cuda() # Make sure it's a 1D long integral tensor
249         test_labels = torch.randint(0, 2, (200,), dtype=torch.long).cuda()
250         test_data = torch.rand(200, 128, 512).cuda() # 200 samples for testing
251         return train_data, train_labels, test_data, test_labels
252
253     train_data, train_labels, test_data, test_labels = prepare_data()
254     print(train_data.shape, train_labels.shape)
255
256     """9. Define the model"""
257
258     # Keep the SimpleClassifierHead class
259     class SimpleClassifierHead(nn.Module):
260         """
261         A simple classification head for transformer output.
262         """
263         def __init__(self, input_dim, num_classes):
264             super(SimpleClassifierHead, self).__init__()
265             self.fc = nn.Linear(input_dim, num_classes)
266
267         def forward(self, x):
268             return self.fc(x)

```

```

269
270
271 # TransformerEncoderLayerWithPruning class
272 ✓ class TransformerEncoderLayerWithPruning(nn.TransformerEncoderLayer):
273     """
274     A customized TransformerEncoderLayer that supports dynamic token skipping.
275     """
276     def __init__(self, *args, **kwargs):
277         super().__init__(*args, **kwargs)
278
279 ✓ def forward(self, src, src_mask=None, src_key_padding_mask=None, keep_tokens=None):
280     """
281     Args:
282         src: Input tensor of shape (batch_size, seq_len, d_model).
283         keep_tokens: Boolean tensor of shape (batch_size, seq_len).
284     """
285     if keep_tokens is not None:
286         # Dynamically crop the tensor shape, keeping only tokens marked True
287         batch_size, seq_len, d_model = src.shape
288         active_indices = keep_tokens.nonzero(as_tuple=True) # Gets the index of active tokens
289         max_active_tokens = keep_tokens.sum(dim=1).max().item() # Maximum number of active tokens
290         pruned_src = torch.zeros(batch_size, max_active_tokens, d_model, device=src.device)
291
292         for batch_idx in range(batch_size):
293             active_token_indices = keep_tokens[batch_idx].nonzero(as_tuple=True)[0]
294             pruned_src[batch_idx, :len(active_token_indices)] = src[batch_idx, active_token_indices]
295
296         src = pruned_src # Update to the clipped tensor
297
298     # A forward method that passes the trimmed tensor to the parent class
299     return super().forward(src, src_mask, src_key_padding_mask)
300
301
302
303
304
305 # create_models function
306 ✓ def create_models():
307     """
308     Create baseline and pruned Transformer models, each with a classification head.
309     Returns:
310         baseline_model, pruned_model
311     """
312     num_classes = 2 # dichotomy
313
314     # Baseline model
315     baseline_encoder = nn.TransformerEncoderLayer(d_model=512, nhead=8)
316     baseline_transformer = nn.TransformerEncoder(baseline_encoder, num_layers=2).cuda()
317     baseline_model = nn.Sequential(
318         baseline_transformer,
319         SimpleClassifierHead(input_dim=512, num_classes=num_classes).cuda()
320     )

```

```

321
322     # Pruned model
323     pruned_encoder = TransformerEncoderLayerWithPruning(d_model=512, nhead=8)
324     pruned_transformer = PrunedTransformerEncoder(pruned_encoder, num_layers=2, saliency_threshold=13.0).cuda()
325     pruned_model = nn.Sequential(
326         pruned_transformer,
327         SimpleClassifierHead(input_dim=512, num_classes=num_classes).cuda()
328     )
329
330     return baseline_model, pruned_model
331
332 baseline_model, pruned_model = create_models()
333 print(baseline_model)
334 print(pruned_model)
335
336 input_tensor = torch.rand(16, 128, 512).cuda()
337
338 pruned_outputs = pruned_model[0](input_tensor)
339 print(f"Output shape after pruning: {pruned_outputs.shape}")
340
341 saliency_scores = pruned_model[0].token_pruning.compute_saliency(input_tensor)
342 print(f"Saliency scores range: {saliency_scores.min().item()} - {saliency_scores.max().item()}")
343
344 keep_tokens = pruned_model[0].token_pruning.compute_saliency(input_tensor) > pruned_model[0].token_pruning.saliency_threshold
345 print(f"Keep tokens mask (sample batch): {keep_tokens[0].cpu().numpy()}")
346
347 """10. Evaluate model accuracy"""
348
349 ✓ def compare_models_accuracy(baseline_model, pruned_model, train_data, train_labels, test_data, test_labels):
350     """
351     Compare accuracy of baseline and pruned models.
352     """
353     print("\n=== Baseline Model Accuracy ===")
354     evaluate_model_accuracy(baseline_model, train_data, train_labels, test_data, test_labels)
355
356     print("\n=== Pruned Model Accuracy ===")
357     evaluate_model_accuracy(pruned_model, train_data, train_labels, test_data, test_labels)
358
359     compare_models_accuracy(baseline_model, pruned_model, train_data, train_labels, test_data, test_labels)
360
361 """11.FLOPs versus memory performance"""
362
363 from torch.profiler import profile, ProfilerActivity
364 from fvcnn.nn import FlopCountAnalysis, flop_count_table
365
366 ✓ def calculate_dynamic_flops_and_profile(pruned_model, input_tensor):
367     """
368     Calculate dynamic FLOPs and memory usage for the pruned model.
369     Args:
370         pruned_model: Model with dynamic token pruning.
371         input_tensor: Input tensor.
372     """
373     # Dynamic computing FLOPs
374     print("\n=== Pruned Model ===")
375     flops_pruned = FlopCountAnalysis(pruned_model, input_tensor)
376     print(flop_count_table(flops_pruned))
377
378     # Dynamic profile memory and time
379     with profile(activities=[ProfilerActivity.CPU, ProfilerActivity.CUDA]) as prof:
380         _ = pruned_model(input_tensor)
381     print(prof.key_averages().table(sort_by="cuda_time_total"))
382

```

```

383  ✓ def compare_efficiency(baseline_model, pruned_model):
384      """
385      Compare FLOPs and memory usage for baseline and pruned models.
386      Args:
387          baseline_model, pruned_model: Models to compare.
388      """
389      input_tensor = torch.rand(16, 128, 512).cuda() # Simulated input: batch size=16, tokens=128, dim=512
390
391      # FLOPs and performance evaluation of Baseline Model
392      print("\n=== Baseline Model ===")
393      flops_baseline = FlopCountAnalysis(baseline_model, input_tensor)
394      print(flop_count_table(flops_baseline))
395      profile_memory_and_time_safe(baseline_model, input_tensor)
396
397      # Pruned Model dynamic FLOPs and performance evaluation
398      calculate_dynamic_flops_and_profile(pruned_model, input_tensor)
399
400      compare_efficiency(baseline_model, pruned_model)

```