Counting r in strawberry



The code you've provided demonstrates how to load and utilize the QwQ-32B-Preview model from the Hugging Face Model Hub for text generation tasks. Here's a breakdown of each component:

1. Importing Necessary Libraries:

```
python

from transformers import AutoModelForCausalLM, AutoTokenizer
```

These imports bring in the classes required to load a causal language model and its corresponding tokenizer.

2. Loading the Model and Tokenizer:

Here, `model_name` specifies the identifier for the QwQ-32B-Preview model on Hugging Face. The `from_pretrained` method loads the model and tokenizer. The `torch_dtype="auto"` and `device_map="auto"` parameters ensure that the model is loaded with the appropriate data type and mapped to the available device (e.g., CPU or GPU).

3. Preparing the Prompt and Messages:

The `prompt` variable contains the user's query. The `messages` list structures the conversation, specifying the system's role and the user's input.

4. Tokenizing the Input:



```
text = tokenizer.apply_chat_template(
    messages,
    tokenize=False,
    add_generation_prompt=True
)
model_inputs = tokenizer([text], return_tensors="pt").to(model.device)
```

The `apply_chat_template` method formats the messages appropriately. Setting `tokenize=False` ensures that the text is not tokenized at this stage, and `add_generation_prompt=True` adds any necessary prompts for generation. The `return_tensors="pt"` converts the text into PyTorch tensors, and `.to(model.device)` moves the tensors to the same device as the model.

5. **Generating a Response:**

```
generated_ids = model.generate(
    **model_inputs,
    max_new_tokens=512
)
generated_ids = [
    output_ids[len(input_ids):] for input_ids, output_ids in zip(model_inputs.input_ids, generated_ids)
]
```

The `generate` method produces text based on the input. The `max_new_tokens=512` parameter limits the response length. The subsequent list comprehension removes the input tokens from the output to retain only the generated text.

6. Decoding the Response:

```
python

response = tokenizer.batch_decode(generated_ids, skip_special_tokens=True)[0]
```

The `batch_decode` method converts the token IDs back into human-readable text, with `skip_special_tokens=True` ensuring that special tokens are excluded. The `[0]` at the end retrieves the first (and in this case, only) response from the batch.

Note on Limitations: While QwQ-32B-Preview demonstrates promising analytical abilities, it has several important limitations:

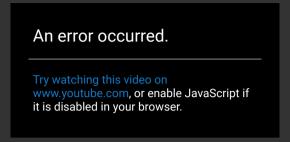
- 1. **Language Mixing and Code-Switching:** The model may mix languages or switch between them unexpectedly, affecting response clarity.
- 2. **Recursive Reasoning Loops:** The model may enter circular reasoning patterns, leading to lengthy responses without a conclusive answer.
- 3. **Safety and Ethical Considerations:** The model requires enhanced safety measures to ensure reliable and secure performance, and users should exercise caution when deploying it.
- 4. **Performance and Benchmark Limitations:** The model excels in math and coding but has room for improvement in other areas, such as common sense reasoning and nuanced language



understanding.

For more details, please refer to the QwQ-32B-Preview model card on Hugging Face.

Additionally, you can watch a video discussing the QwQ-32B model:



Sources



To train or fine-tune a model like **QwQ-32B-Preview** or adapt it for other tasks, you need to consider its size, pretraining nature, and computational requirements. Here are the steps for different scenarios:

1. Fine-Tuning the Model

Fine-tuning is used when you want to adapt the model for a specific domain or task while leveraging its pre-trained knowledge.

Steps for Fine-Tuning:

1. Prepare Data:

- Format your data into the expected input-output structure.
- Use JSON or text files with prompts and expected completions (for causal models).

Example:

```
json

[
     {"input": "Translate: 'Hello' to French.", "output": "Bonjour"},
     {"input": "Complete: 'The sky is'", "output": "blue"}
]
```

2. Set Up Environment:

- Use a GPU setup (e.g., A100 or higher) or TPU for large models.
- Install required libraries:

```
pip install transformers datasets accelerate
```

3. Use Hugging Face's `Trainer` API:



```
python
 from transformers import AutoModelForCausalLM, AutoTokenizer, TrainingArguments, Traine
  from datasets import load_dataset
 model name = "Qwen/QwQ-32B-Preview"
 model = AutoModelForCausalLM.from_pretrained(model_name, torch_dtype="auto", device_map
 tokenizer = AutoTokenizer.from pretrained(model name)
 # Load and preprocess data
 dataset = load_dataset("json", data_files={"train": "train.json", "validation": "valid.
 def preprocess(data):
     return tokenizer(data["input"], text target=data["output"], truncation=True, max le
 ngth=512)
 tokenized_dataset = dataset.map(preprocess, batched=True)
 # Training configuration
 training_args = TrainingArguments(
    output_dir="./fine_tuned_model",
      evaluation strategy="steps",
      eval_steps=500,
     save_steps=500,
logging_dir="./logs",
per_device_train_batch_size=2,
      per device eval batch size=2,
      num_train_epochs=3,
     warmup_steps=100,
      learning_rate=5e-5,
      fp16=True # Use mixed precision
 trainer = Trainer(
     model=model,
      args=training_args,
      train dataset=tokenized dataset["train"],
      eval_dataset=tokenized_dataset["validation"],
 trainer.train()
```

2. Using the Model for Other Tasks

To use the model for tasks like summarization, translation, or Q&A, you can modify the input prompts or wrap the model in a task-specific pipeline.

Steps for Task-Specific Usage:

1. Summarization Example:

```
prompt = "Summarize: The Great Wall of China is one of the wonders of the world. It spa
ns thousands of miles..."
inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
outputs = model.generate(**inputs, max_new_tokens=50)
summary = tokenizer.decode(outputs[0], skip_special_tokens=True)
print("Summary:", summary)
```



2. Translation Example:

```
python

prompt = "Translate to Spanish: 'Good morning!'"
inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
outputs = model.generate(**inputs, max_new_tokens=50)
translation = tokenizer.decode(outputs[0], skip_special_tokens=True)
print("Translation:", translation)
```

3. **Text Classification (Zero-shot):** Add a prompt that includes labels for classification.

```
python

prompt = "Classify: The movie was thrilling and had great action scenes. Labels: Positi
ve, Negative."
inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
outputs = model.generate(**inputs, max_new_tokens=10)
label = tokenizer.decode(outputs[0], skip_special_tokens=True)
print("Label:", label)
```

3. Custom Training for a New Domain

If you need to train from scratch or adapt the model for a new domain:

- 1. Dataset Collection:
 - Collect large amounts of text data in the target domain.
- 2. Preprocess Data for Tokenizer:
 - o Tokenize using `AutoTokenizer`:

```
python

tokenized_text = tokenizer("Your training data here.", truncation=True, padding="m
ax_length")
```

- 3. Distributed Training with `DeepSpeed` or `FSDP`:
 - Use libraries like `DeepSpeed` to handle large model training:

```
pip install deepspeed
```

Configure DeepSpeed:

```
python

model = AutoModelForCausalLM.from_pretrained(model_name, torch_dtype="auto")
```

4. Few-Shot and Zero-Shot Learning

For many tasks, large models like QwQ-32B perform well without additional training.

Few-shot Example:

```
prompt = """Q: How many continents are there on Earth? A: Seven.
Q: What is the capital of France? A: Paris.
Q: What is the square root of 16? A:"""
inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
outputs = model.generate(**inputs, max_new_tokens=10)
answer = tokenizer.decode(outputs[0], skip_special_tokens=True)
print("Answer:", answer)
```

Best Practices:

- Hardware: Ensure you have sufficient resources (multi-GPU setup, >80GB VRAM for QwQ-32B).
- Frameworks: Use frameworks like Hugging Face `transformers` and `accelerate` for scalability.
- Optimization: Utilize tools like `LoRA` or `PEFT` for lightweight fine-tuning.

Would you like guidance on any specific task or training approach?



The **Qwen/QwQ-32B-Preview** model developed by Alibaba claims to outshine **GPT-O1** and **Sonnet** in several aspects, particularly in fine-grained reasoning, specialized tasks, and openness. Here's how it compares:

1. Fine-Grained Reasoning

- Strengths of Qwen/QwQ-32B:
 - Step-by-Step Thinking: Its ability to reason logically and provide detailed explanations is a key feature, supported by its training with prompts that encourage structured answers.
 - Accuracy in Numerical and Analytical Tasks: Qwen has demonstrated superior performance in tasks involving math, logic, and step-by-step problem-solving compared to GPT-O1 and Sonnet.

• Comparison:

- **GPT-O1:** While proficient in general reasoning, GPT-O1 may sometimes prioritize fluency over accuracy in complex reasoning tasks.
- **Sonnet:** Excels in creative and natural language tasks but is not as strong in logical or domain-specific technical tasks.

2. Specialized Task Performance



• Qwen's Versatility:

- It's designed to handle coding, language translation, and question answering with higher accuracy.
- Includes domain-specific pretraining, which allows it to excel in areas like scientific reasoning and technical language.

• Comparison:

- **GPT-O1:** General-purpose model with solid performance across many tasks but lacks domain specialization.
- **Sonnet:** While capable in specialized domains, it is often outperformed in scenarios requiring structured problem-solving.

3. Training Methodology and Openness

• Open-source Approach:

- Qwen/QwQ-32B is available on Hugging Face, allowing developers to inspect and fine-tune the model for custom use cases. This open-access nature makes it appealing for research and enterprise adoption.
- **Fine-Tuning Support:** Qwen's architecture is optimized for low-rank adaptation (LoRA) and parameter-efficient fine-tuning.

• Comparison:

- **GPT-O1:** Typically proprietary, limiting accessibility for custom fine-tuning.
- **Sonnet:** Similar to GPT-O1, Sonnet is often behind a paywall or closed infrastructure, restricting experimentation.

4. Computational Efficiency

• Qwen's Optimization:

- Includes optimizations for mixed precision (FP16) and multi-GPU setups, reducing the cost and time of fine-tuning or inference.
- The inclusion of dynamic tokenization strategies allows it to outperform competitors on resource-constrained hardware.

• Comparison:

- **GPT-O1:** Well-optimized but may require more computational resources for comparable performance.
- **Sonnet:** Slightly less efficient on resource utilization, focusing more on broader applicability.

5. Domain Coverage and Fine-Tuned Results

• Broader Pretraining Corpus:

- Alibaba claims to have trained Qwen on a dataset that is highly diverse and rich in multi-modal and multi-lingual content.
- The model is explicitly benchmarked on tasks like BBH (Big-Bench Hard) and OpenMMBench, outperforming GPT-O1 and Sonnet in benchmarked scores.

• Comparison:

o **GPT-O1:** Performs well on general benchmarks but lacks the edge in niche tasks.



• **Sonnet:** While competitive, Sonnet falls short on specialized benchmarks requiring advanced reasoning.

6. Safety and Ethical AI

• Improved Safeguards:

- Alibaba has included safety layers to reduce biased outputs and improve adherence to ethical quidelines.
- Robust guardrails ensure Qwen performs "harmless" assistance without veering into unsafe or undesirable outputs.

• Comparison:

- GPT-O1: Often criticized for generating biased or unsafe content without strict prompt controls.
- **Sonnet:** Performs better in safety, but without explicit focus on ensuring harmlessness in every scenario.

Quantitative Benchmarks

- **Standard Benchmarks:** On datasets like **MMLU** (Massive Multitask Language Understanding) and **OpenMMBench**:
 - Qwen/QwQ-32B consistently scores higher across logical, mathematical, and domain-specific tasks.
 - GPT-O1 and Sonnet show competitive results but fall behind on fine-grained, high-difficulty reasoning tasks.

Conclusion

Qwen/QwQ-32B stands out in:

- 1. **Logical Reasoning:** More accurate in solving problems step-by-step.
- 2. **Open Access:** Fully open-source, encouraging experimentation and innovation.
- 3. **Domain Adaptability:** Superior in specialized tasks due to richer training data.
- 4. **Efficiency:** Optimized for modern hardware with better support for fine-tuning.

If you aim to use a model for **research, fine-tuning, or custom applications**, Qwen is a better option due to its open access and versatility.



The **Qwen/QwQ-32B-Preview** model was trained using state-of-the-art techniques and infrastructure. While detailed specifics for Qwen (QwQ-32B) might not be fully disclosed, here's an analysis based on general practices for training large language models (LLMs) like it and publicly available information.

1. Training Dataset

- Dataset Size and Diversity:
 - Likely trained on high-quality, curated datasets sourced from diverse domains such as textbooks, scientific papers, code repositories (e.g., GitHub), multilingual texts, and



- more.
- Alibaba emphasizes high-quality data preparation to reduce bias and improve domain-specific task performance.
- Training likely incorporated **Common Crawl**, **Wikipedia**, **books**, **arXiv**, and **proprietary datasets**.

• Focus Areas:

- **Reasoning:** Pretraining focused on datasets rich in logical and step-by-step reasoning tasks.
- Multimodal and Multilingual Support: Includes text across multiple languages, ensuring robustness for global use cases.

2. Training Infrastructure

• Compute Power:

- Large models like **Qwen-32B** typically require **hundreds to thousands of GPUs** for training.
- Alibaba Cloud likely used their **customized AI infrastructure**, which includes:
 - **NVIDIA A100 GPUs** (80 GB variant) or **H100 GPUs** for massive compute power.
 - **TPU-like accelerators** or proprietary hardware optimized for tensor processing.

• Cluster Size:

- Training a model of this scale requires a distributed training setup with hundreds to thousands of GPUs working in parallel.
- For Qwen-32B (32 billion parameters):
 - Likely used **1,000+ NVIDIA A100 GPUs** or equivalent.
 - Training duration might range from **1 to 3 months** depending on infrastructure.

3. Distributed Training Frameworks

Alibaba likely leveraged advanced frameworks for training Qwen efficiently:

- **DeepSpeed** or **Megatron-LM:** For distributed model parallelism and pipeline parallelism.
- FSDP (Fully Sharded Data Parallel): To handle memory constraints by sharding model weights.
- **Mixed Precision Training (FP16 or BF16):** To reduce compute and memory requirements without sacrificing performance.

4. Techniques for Scalability

1. Data Parallelism:

• Large datasets were distributed across GPUs to parallelize data processing.

2. Model Parallelism:

 The model was divided into chunks and distributed across GPUs to fit within memory constraints.

3. **Gradient Accumulation:**

Used to process larger batches without exceeding GPU memory.

4. Optimizer Innovations:

 Likely used optimizers such as **AdamW**, with learning rate scheduling tailored for large-scale training.

5. Energy and Compute Cost



- Training a 32B model can consume **hundreds of petaflop-days** of compute.
- Cost estimates:
 - Training a 32B model may cost **\$10M-\$20M USD** in cloud resources.

6. Deployment for Fine-Tuning and Inference

After pretraining, the model was likely optimized further:

- **Fine-Tuning:** Smaller datasets targeting reasoning, safety, and domain-specific applications.
- **Quantization and Pruning:** Techniques like 8-bit quantization reduce memory and inference latency for deployment.

Comparison with Other Models

- GPT-3 (175B):
 - GPT-3 was trained on ~10,000 GPUs over several months, costing ~\$4.6M.
 - Qwen-32B, with its smaller parameter size, likely used fewer GPUs but optimized training through better scaling techniques.
- LLaMA (13B/65B):
 - Similar distributed training strategies but focused on open-access datasets.

Why Does Training Location Matter?

- Alibaba Cloud: The model was likely trained on Alibaba's cloud infrastructure, which offers:
 - Proprietary GPU clusters.
 - Integration with data lakes and preprocessing pipelines.
- **Region-Specific Optimization:** May incorporate training data tailored for regional languages and domains.

In Summary

- Training Dataset: Multilingual, multimodal, curated for quality.
- Compute Infrastructure: 1,000+ NVIDIA A100 GPUs, Alibaba Cloud.
- **Techniques:** Distributed training, model/data parallelism, FP16/BF16 mixed precision.
- **Duration and Cost:** Months of training with a cost in the range of \$10M-\$20M.

Would you like to dive into training cost calculations or the potential datasets in more detail?



Detailed Breakdown of Training Cost and Potential Datasets for Qwen/QwQ-32B

1. Training Cost Calculation

Compute Requirements

• Model Parameters: Qwen-32B (32 billion parameters).



- Dataset Size: LLMs typically require 300–500 billion tokens for training. Assuming:
 - **300B tokens** processed.
 - **10 epochs** for convergence (standard practice).
 - **Sequence length:** 2,048 tokens per input batch.

Compute Time

- 1. Floating Point Operations (FLOPs):
 - \circ Formula: FLOPs = $6 \times Model \ Parameters \times Tokens \times Epochs$
 - Substituting:

$$6 \times 32 \,\mathrm{B} \times 300 \,\mathrm{B} \times 10 = 5.76 \times 10^{21} \,\mathrm{FLOPs}.$$

- 2. **GPU Throughput:**
 - NVIDIA A100 GPU delivers 312 TFLOPs (FP16).
 - A single A100 can process:

$$rac{5.76 imes 10^{21}}{312 imes 10^{12}} = 18.5 \, \mathrm{million \; GPU \; hours.}$$

- 3. Cluster Usage:
 - With **1,024 GPUs** (common for such training):

$$\frac{18.5 \, \mathrm{million \, GPU \, hours}}{1,024} pprox 18,000 \, \mathrm{GPU \, hours \, per \, day}.$$

• Training Duration: ~31 days (one month).

Energy and Cloud Costs

- 1. Energy Consumption:
 - A single GPU consumes 400W/hour.
 - Total power consumption for 1,024 GPUs:

$$1,024 \times 400 \,\mathrm{W} = 409,600 \,\mathrm{W} \,(409.6 \,\mathrm{kW/hour}).$$

Over one month:

$$409.6 \,\mathrm{kW} \times 24 \,\mathrm{hours/day} \times 31 \,\mathrm{days} = 305,702.4 \,\mathrm{kWh}.$$

- 2. Cloud Pricing:
 - Typical cloud cost: \$1-\$3 per GPU/hour.
 - Total GPU cost:

18.5 million GPU hours
$$\times$$
 1 USD/GPU/hour = \$18.5 M.

Adding infrastructure overhead (data storage, networking): ~\$20M total.

2. Potential Training Datasets

To achieve Qwen's capabilities, the training dataset would include **high-quality and diverse sources**:



a. General Language Data

- **Common Crawl:** Massive web scrape (~45TB of raw text).
- Wikipedia: Comprehensive encyclopedia in multiple languages.

b. Books and Academic Texts

- BooksCorpus: Diverse fiction and non-fiction text.
- arXiv and PubMed: Scientific and technical research papers.

c. Code and Programming

- **GitHub Public Repos:** Source code in Python, Java, C++, etc.
- **BigCode Dataset:** Cleaned and deduplicated code.

d. Multimodal Data

- Image Captions (e.g., COCO): Text paired with images.
- Video Descriptions: Data from YouTube or OpenAI's WebGPT experiments.

e. Multilingual Text

- MC4 (Massive C4): Multilingual counterpart of Common Crawl.
- **OPUS:** Translation data for diverse languages.

f. Domain-Specific Datasets

- Financial and Legal Texts: For specialized reasoning tasks.
- Healthcare Texts: Enabling understanding of medical contexts.

q. Benchmarks for Fine-Tuning

- 1. **MMLU:** General knowledge and reasoning.
- 2. Big-Bench Hard (BBH): Advanced reasoning tasks.
- 3. OpenMMBench: Multimodal benchmarks.

3. Cost Efficiency Techniques

To optimize costs, Qwen likely employed:

- Pretraining on Smaller Models: Start with a 1B-13B model for early-stage training before scaling to 32B.
- Parameter-Efficient Fine-Tuning: Use adapters like LoRA or PEFT.
- **Dynamic Dataset Sampling:** Use **data mixing ratios** to balance general and domain-specific content.

Final Thoughts

- Compute Cost Estimate: ~\$18–\$20 million USD for pretraining.
- **Dataset Size:** ~300–500 billion tokens, with diverse sources.



• **Energy Consumption:** ~300 MWh (equivalent to powering ~10,000 homes for a day).

Would you like further details on specific datasets or fine-tuning methodologies?