Featherless AI Description </i>
⟨/> View API Code Dismiss Send Gemma is a series of best-in-class open models and draws inspiration and technological lineage </>
View Code Snippets Open Playground from the Gemini family of models. They are text-to-text, decoder-only large language models with open weights. Gemma models are well-suited for a variety of text generation tasks, including Model tree for google/gemma-2-2b-jpn-it ⑤ question answering, summarization, and reasoning. Base model google/gemma-2-2b Gemma-2-JPN is a Gemma 2 2B model fine-tuned on Japanese text. It supports the Japanese google/gemma-2-2b-it **G** Finetuned language with the same level of performance of English only queries on Gemma 2. **Finetuned** (<u>697</u>) Adapters <u>8 models</u> Usage **Finetunes** 27 models Merges <u>2 models</u> Below we share some code snippets on how to get quickly started with running the model. First, Quantizations 22 models install the Transformers library with: **∷** Spaces using google/gemma-2-2b-jpn-it 20 pip install -U transformers Sarath0x8f/Document-QA-bot | 🟃 bhaskartripathi/LLM_Quantization FallnAI/Quantize-HF-Models > seawolf2357/LLM_Quantization Then, copy the snippet from the section that is relevant for your usecase. Running with the pipeline API openfree/LLM_Quantization
KBaba7/Quant import torch

from transformers import pipeline **□ Collections including** google/gemma-2-2b-jpn-it pipe = pipeline("text-generation", Gemma 2 JPN Release Collection model="google/gemma-2-2b-jpn-it", A Gemma 2 2B model fine-tuned on Japa... • 3 items • Updated Jul 10 • \triangle 28 model_kwargs={"torch_dtype": torch.bfloat16}, device="cuda", # replace with "mps" to run on a Mac device Google's Gemma models family Collection 328 items • Updated 10 days ago • △ 500 messages = [{"role": "user", "content": "マシーンラーニングについての詩を書いてください。"},

outputs = pipe(messages, return_full_text=False, max_new_tokens=256) assistant_response = outputs[0]["generated_text"].strip() print(assistant_response) ► Example output

English:\n\n{assistant_response}" messages = [{"role": "user", "content": translation_input_text}, outputs = pipe(messages, return_full_text=False, max_new_tokens=1024) translated_response = outputs[0]["generated_text"].strip() print(translated_response) ► Example output Running the model on a single / multi GPU # pip install accelerate from transformers import AutoTokenizer, AutoModelForCausalLM import torch

inputs['input_ids'].shape[1]:], skip_special_tokens=True)[0] print(generated_text.strip()) Running the model on a GPU using different precisions The native weights of this model were exported in bfloat16 precision. You can also use float32 if you skip the dtype, but no precision increase will occur (model weights will just be upcasted to float32). See examples below. • Upcasting to torch.float32 # pip install accelerate from transformers import AutoTokenizer, AutoModelForCausalLM

It can also be used for translation, as follows:

translation_input_text = f"Translate the following poem from Japanese to

tokenizer = AutoTokenizer.from_pretrained("google/gemma-2-2b-jpn-it")

inputs = tokenizer.apply_chat_template(messages, return_tensors="pt",

tokenizer = AutoTokenizer.from_pretrained("google/gemma-2-2b-jpn-it")

inputs = tokenizer.apply_chat_template(messages, return_tensors="pt",

add_generation_prompt=True, return_dict=True).to(model.device)

inputs['input_ids'].shape[1]:], skip_special_tokens=True)[0]

outputs = model.generate(**inputs, max_new_tokens=256)

generated_text = tokenizer.batch_decode(outputs[:,

Data used for model training and how the data was processed.

representation, and to address mathematical queries.

8 trillion tokens. Here are the key components:

model = AutoModelForCausalLM.from_pretrained(

"google/gemma-2-2b-jpn-it",

device_map="auto",

print(generated_text.strip())

messages = [

Inputs and outputs

Model Data

Training Dataset

from training sets.

Implementation Information

Details about the model internals.

Hardware

domain:

Software

models like these ones.

defaulting to English.

Benchmark

Ethics and Safety

Evaluation Approach

Preference vs GPT-3.5 -0.25 ± 0.05

Ethics and safety evaluation approach and results.

personally identifiable information exposure.

These models have certain limitations that users should be aware of.

poems, scripts, code, marketing copy, and email drafts.

grammar correction or providing writing practice.

summaries or answering questions about specific topics.

and nuclear (CBRN) risks.

part of model training and development.

Content Creation and Communication

advancement of the field.

Limitations

Training Data

responses.

effectively.

Factual Accuracy

Common Sense

Bias and Fairness

■ System theme

Misinformation and Misuse

Generative AI Toolkit.

statements.

sarcasm, or figurative language.

Usage and Limitations

Intended Usage

Language correctness | 86.47%

add_generation_prompt=True, return_dict=True).to(model.device)

outputs = model.generate(**inputs, max_new_tokens=256)

generated_text = tokenizer.batch_decode(outputs[:,

{"role": "user", "content": "マシーンラーニングについての詩を書いてください。"},

model = AutoModelForCausalLM.from_pretrained(

"google/gemma-2-2b-jpn-it",

torch_dtype=torch.bfloat16,

device_map="auto",

messages = [

• Input: Text string, such as a question, a prompt, or a document to be summarized. • Output: Generated Japanese-language text in response to the input, such as an answer to a question, or a summary of a document.

These models were trained on a dataset of text data that includes a wide variety of sources, totaling

• Web Documents: A diverse collection of web text ensures the model is exposed to a broad range

languages, which improves its ability to generate code or understand code-related questions.

• Mathematics: Training on mathematical text helps the model learn logical reasoning, symbolic

• Instruction data set: large-scale and high-quality Japanese and multilingual instruction data.

The combination of these diverse data sources is crucial for training a powerful language model that

• Code: Exposing the model to code helps it to learn the syntax and patterns of programming

of linguistic styles, topics, and vocabulary. Primarily English-language content.

{"role": "user", "content": "マシーンラーニングについての詩を書いてください。"},

can handle a wide variety of different tasks and text formats. **Data Preprocessing** Here are the key data cleaning and filtering methods applied to the training data: • CSAM Filtering: Rigorous CSAM (Child Sexual Abuse Material) filtering was applied at multiple stages in the data preparation process to ensure the exclusion of harmful and illegal content.

• Sensitive Data Filtering: As part of making Gemma pre-trained models safe and reliable, we

• Additional methods: Filtering based on content quality and safety in line with <u>our policies</u>.

used automated techniques to filter out certain personal information and other sensitive data

• Performance: TPUs are specifically designed to handle the massive computations involved in training LLMs. They can speed up training considerably compared to CPUs. • Memory: TPUs often come with large amounts of high-bandwidth memory, allowing for the handling of large models and batch sizes during training. This can lead to better model quality. • Scalability: TPU Pods (large clusters of TPUs) provide a scalable solution for handling the growing complexity of large foundation models. You can distribute training across multiple TPU

Cost-effectiveness: In many scenarios, TPUs can provide a more cost-effective solution for

JAX allows researchers to take advantage of the latest generation of hardware, including TPUs, for

ML Pathways is Google's latest effort to build artificially intelligent systems capable of generalizing

across multiple tasks. This is specially suitable for <u>foundation models</u>, including large language

Together, JAX and ML Pathways are used as described in the <u>paper about the Gemini family of</u>

These advantages are aligned with <u>Google's commitments to operate sustainably</u>.

training large models compared to CPU-based infrastructure, especially when considering the

devices for faster and more efficient processing.

time and resources saved due to faster training.

Training was done using <u>JAX</u> and <u>ML Pathways</u>.

faster and more efficient training of large models.

Gemma was trained using the latest generation of <u>Tensor Processing Unit (TPU)</u> hardware (TPUv5p).

Training large language models requires significant computational power. TPUs, designed

specifically for matrix operations common in machine learning, offer several advantages in this

<u>models</u>; "the 'single controller' programming model of Jax and Pathways allows a single Python process to orchestrate the entire training run, dramatically simplifying the development workflow." **Evaluation** To assess the quality of this model, we collected a diverse set of Japanese prompts and evaluated performance using an LLM-as-a-judge approach against GPT-3.5. The rating system is based on a 7scale assessments, which are MuchBetterThan, BetterThan, SlightlyBetterThan, AboutTheSame, SlightlyWorse, WorseThan, MuchWorseThan associated with the numerical scores 1.5, 1.0, 0.5, 0,

-0.5, -1.0, -1.5 respectively. We also tracked the ability of the model to answer in the correct

Gemma-2-IT Gemma-2-IT-JPN

 0.03 ± 0.04

98.24%

language: for a Japanese prompt, the model should typically answer in Japanese rather than

categories relevant to ethics and safety, including: Text-to-Text Content Safety: Human evaluation on prompts covering safety policies including child sexual abuse and exploitation, harassment, violence and gore, and hate speech.

• Memorization: Automated evaluation of memorization of training data, including the risk of

• Large-scale harm: Tests for "dangerous capabilities," such as chemical, biological, radiological,

Open Large Language Models (LLMs) have a wide range of applications across various industries and

provide contextual information about the possible use-cases that the model creators considered as

• Text Generation: These models can be used to generate creative text formats such as

domains. The following list of potential uses is not comprehensive. The purpose of this list is to

• Text-to-Text Representational Harms: Benchmark against relevant academic datasets.

Our evaluation methods include structured evaluations and internal red-teaming testing of relevant

content policies. Red-teaming was conducted by a number of different teams, each with different

goals and human evaluation metrics. These models were evaluated against a number of different

• Chatbots and Conversational AI: Power conversational interfaces for customer service, virtual assistants, or interactive applications. • Text Summarization: Generate concise summaries of a text corpus, research papers, or reports. Research and Education • Natural Language Processing (NLP) Research: These models can serve as a foundation for researchers to experiment with NLP techniques, develop algorithms, and contribute to the

• Language Learning Tools: Support interactive language learning experiences, aiding in

Knowledge Exploration: Assist researchers in exploring large bodies of text by generating

• The quality and diversity of the training data significantly influence the model's

capabilities. Biases or gaps in the training data can lead to limitations in the model's

• The scope of the training dataset determines the subject areas the model can handle

 Context and Task Complexity • LLMs are better at tasks that can be framed with clear prompts and instructions. Openended or highly complex tasks might be challenging. A model's performance can be influenced by the amount of context provided (longer context generally leads to better outputs, up to a certain point). Language Ambiguity and Nuance • Natural language is inherently complex. LLMs might struggle to grasp subtle nuances,

• LLMs generate responses based on information they learned from their training datasets,

but they are not knowledge bases. They may generate incorrect or outdated factual

• LLMs rely on statistical patterns in language. They might lack the ability to apply common sense reasoning in certain situations. **Ethical Considerations and Risks** The development of large language models (LLMs) raises several ethical concerns. In creating an open model, we have carefully considered the following:

• LLMs trained on large-scale, real-world text data can reflect socio-cultural biases

pre-processing described and posterior evaluations reported in this card.

• LLMs can be misused to generate text that is false, misleading, or harmful.

• Guidelines are provided for responsible use with the model, see the <u>Responsible</u>

embedded in the training material. These models underwent careful scrutiny, input data

• Transparency and Accountability: • This model card summarizes details on the models' architecture, capabilities, limitations, and evaluation processes. • A responsibly developed open model offers the opportunity to share innovation by making LLM technology accessible to developers and researchers across the AI ecosystem. Risks identified and mitigations:

Perpetuation of biases: It's encouraged to perform continuous monitoring (using evaluation

metrics, human review) and the exploration of de-biasing techniques during model training,

fine-tuning, and other use cases. • Generation of harmful content: Mechanisms and guidelines for content safety are essential. Developers are encouraged to exercise caution and implement appropriate content safety safeguards based on their specific product policies and application use cases. • Misuse for malicious purposes: Technical limitations and developer and end-user education can help mitigate against malicious applications of LLMs. Educational resources and reporting mechanisms for users to flag misuse are provided. Prohibited uses of Gemma models are outlined in the **Gemma Prohibited Use Policy**.

Identifiable Information). Developers are encouraged to adhere to privacy regulations with privacy-preserving techniques. **Benefits**

• Privacy violations: Models were trained on data filtered for removal of PII (Personally

TOS

At the time of release, this family of models provides high-performance open large language model implementations designed from the ground up for Responsible AI development compared to similarly sized models.

Privacy

About

Jobs

Models

Datasets

Spaces

Pricing

Docs