

```
!pip install -q -U transformers peft accelerate bitsandbytes trl datasets evaluate rouge_score google-generativeai tqdm
```

```

40.9/40.9 kB 2.9 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
10.8/10.8 MB 75.4 MB/s eta 0:00:00
472.3/472.3 kB 26.3 MB/s eta 0:00:00
72.9/72.9 MB 11.4 MB/s eta 0:00:00
375.8/375.8 kB 25.0 MB/s eta 0:00:00
491.5/491.5 kB 32.6 MB/s eta 0:00:00
193.6/193.6 kB 14.9 MB/s eta 0:00:00
363.4/363.4 MB 3.5 MB/s eta 0:00:00
13.8/13.8 MB 25.3 MB/s eta 0:00:00
24.6/24.6 MB 29.6 MB/s eta 0:00:00
883.7/883.7 kB 30.9 MB/s eta 0:00:00
664.8/664.8 MB 2.9 MB/s eta 0:00:00
211.5/211.5 MB 5.8 MB/s eta 0:00:00
56.3/56.3 MB 13.7 MB/s eta 0:00:00
127.9/127.9 MB 8.7 MB/s eta 0:00:00
207.5/207.5 MB 5.8 MB/s eta 0:00:00
21.1/21.1 MB 42.8 MB/s eta 0:00:00
Building wheel for rouge_score (setup.py) ... done
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is
gcsfs 2025.3.2 requires fsspec==2025.3.2, but you have fsspec 2025.3.0 which is incompatible.

```

```
import huggingface_hub
huggingface_hub.login()
```

```

import google.generativeai as genai
import json
import time
from google.colab import userdata

# --- Configuration ---
# 1. Add your Gemini API Key to Colab Secrets (left sidebar -> 🗝 icon)
GEMINI_API_KEY = userdata.get('GOOGLE_API_KEY')
genai.configure(api_key=GEMINI_API_KEY)

# 2. Define the model and prompt
model = genai.GenerativeModel(model_name='gemini-1.5-flash-latest')
prompt_template = """
You are an expert financial analyst. Your task is to analyze the following news article and output ONLY a valid, single JSON object.

Schema: {"event_summary": "one-sentence summary", "event_type": "Classification [Market Trend, Regulation, Geopolitical Tension, ...]"}

INPUT ARTICLE:
{article_text}

OUTPUT JSON:
"""

# 3. A completely new set of source articles for testing
source_articles = [
    {"title": "Global Semiconductor Shortage Shows Signs of Easing, says TSMC Chairman", "description": "Taiwan Semiconductor Manufacturing Company (TSMC) Chairman Wu Huangping said the global semiconductor shortage is showing signs of easing, but it will take time to fully resolve."},
    {"title": "SEC delays decision on Grayscale's spot Bitcoin ETF application", "description": "The U.S. Securities and Exchange Commission (SEC) has delayed its decision on Grayscale's application for a spot Bitcoin exchange-traded fund (ETF)."},
    {"title": "US Job Openings Fall Sharply, Suggesting Labor Market is Cooling", "description": "The latest JOLTS report showed a significant decline in job openings, indicating a cooling labor market."},
    {"title": "Ford announces $3.7 billion investment in US plants to build more EVs and traditional cars", "description": "Ford Motor Company announced a $3.7 billion investment in its US manufacturing plants to increase production of both electric vehicles (EVs) and traditional internal combustion engine (ICE) vehicles."},
    {"title": "Amazon's Union Battle Continues as Staten Island Warehouse Files for Second Election", "description": "The Amazon Workers Union (AWU) is preparing for a second election at the Staten Island warehouse, as the company continues to resist unionization."},
    {"title": "Russia cuts off gas supplies to Finland, escalating energy dispute", "description": "Russia has cut off gas supplies to Finland, escalating a long-standing energy dispute between the two countries."},
    {"title": "Netflix explores live streaming for stand-up specials and reality shows", "description": "Netflix is exploring live streaming options for stand-up comedy specials and reality television shows to boost engagement."},
    {"title": "Apple supplier Foxconn forecasts better-than-expected Q2 amid easing lockdowns in China", "description": "Foxconn Technology Group, a major supplier of Apple products, forecasts a better-than-expected second quarter (Q2) performance due to easing COVID-19 lockdowns in China."},
    {"title": "US FDA authorizes first COVID-19 breathalyzer test", "description": "The U.S. Food and Drug Administration (FDA) has authorized the first COVID-19 breathalyzer test for use in the United States."},
    {"title": "Disney CEO Bob Chapek faces backlash over response to Florida's 'Don't Say Gay' bill", "description": "Disney CEO Bob Chapek is facing significant backlash and criticism for his response to Florida's controversial 'Don't Say Gay' bill."}
]

# --- Main Execution Logic ---
ground_truth_data = []
print("--- Generating Unseen Test Set ---")

for i, article in enumerate(source_articles):
    article_text = f"""Title: {article['title']}\n\nDescription: {article['description']}"""
    full_prompt = prompt_template.format(article_text=article_text)

    print(f"Generating ground truth for article {i+1}/{len(source_articles)}...")

```

```

try:
    response = model.generate_content(full_prompt)
    parsed_output = json.loads(response.text)

    clean_entry = {
        "instruction": "Analyze the following news article. Identify the core event, its direct first-order impacts, and hypo",
        "input": article_text,
        "output": json.dumps(parsed_output)
    }
    ground_truth_data.append(clean_entry)
    time.sleep(2)

except Exception as e:
    print(f" -> WARNING: Could not generate data for this article: {e}. Skipping.")

def save_to_jsonl(data, filename):
    with open(filename, 'w') as f:
        for entry in data:
            f.write(json.dumps(entry) + '\n')

save_to_jsonl(ground_truth_data, 'unseen_test_set.jsonl')

print(f"\n✅ Successfully created unseen_test_set.jsonl with {len(ground_truth_data)} examples.")

--- Generating Unseen Test Set ---
Generating ground truth for article 1/10...
Generating ground truth for article 2/10...
Generating ground truth for article 3/10...
Generating ground truth for article 4/10...
Generating ground truth for article 5/10...
Generating ground truth for article 6/10...
Generating ground truth for article 7/10...
Generating ground truth for article 8/10...
Generating ground truth for article 9/10...
Generating ground truth for article 10/10...

✅ Successfully created unseen_test_set.jsonl with 10 examples.

import json
import torch
import evaluate
from transformers import AutoTokenizer, AutoModelForCausalLM, BitsAndBytesConfig
from peft import PeftModel
import google.generativeai as genai
from tqdm import tqdm
from google.colab import userdata

# --- CONFIGURATION ---
# 1. Your existing fine-tuned model
YOUR_HF_USERNAME = "sharathreddy"
KAIROS_MODEL_NAME = "kairos-v1-llama3-8b-instruct" # The name of your first trained model
peft_model_id = "SharathReddy/kairos-llama3-finetune"
base_model_id = "meta-llama/Meta-Llama-3-8B-Instruct"

# 2. Gemini Baseline Model
GEMINI_API_KEY = userdata.get('GOOGLE_API_KEY')
genai.configure(api_key=GEMINI_API_KEY)
baseline_model = genai.GenerativeModel(model_name='gemini-1.5-pro-latest') # Using the more powerful model for a strong baseline

# --- LOAD YOUR KAIROS MODEL ---
print("Loading your fine-tuned KAIROS model...")
bnb_config = BitsAndBytesConfig(load_in_4bit=True, bnb_4bit_quant_type="nf4", bnb_4bit_compute_dtype=torch.bfloat16)
base_model = AutoModelForCausalLM.from_pretrained(base_model_id, quantization_config=bnb_config, device_map="auto")
tokenizer = AutoTokenizer.from_pretrained(base_model_id)
kairos_model = PeftModel.from_pretrained(base_model, peft_model_id)
print(f"✅ KAIROS Model loaded.")

# --- LOAD TEST DATA ---
with open('unseen_test_set.jsonl', 'r') as f:
    test_data = [json.loads(line) for line in f]

# --- GENERATE PREDICTIONS FROM BOTH MODELS ---
results = []
baseline_prompt_template = "Analyze the following news article. Identify the core event, its summary, type, direct impacts, and

print("\n--- Generating predictions from KAIROS and Baseline models ---")
for entry in tqdm(test_data):
    instruction = entry['instruction']

```

```

input_text = entry['input']
ground_truth_output = json.loads(entry['output'])

# KAIROS Prediction
prompt = f"<|begin_of_text|><|start_header_id|>system<|end_header_id|>\n\n{instruction}<|eot_id|><|start_header_id|>user<|eot_id|>\n\n{input_text}<|start_header_id|>assistant<|end_header_id|>"
inputs = tokenizer(prompt, return_tensors="pt").to("cuda")
outputs = karios_model.generate(**inputs, max_new_tokens=1024, do_sample=False)
karios_response_text = tokenizer.decode(outputs[0], skip_special_tokens=True).split("<|start_header_id|>assistant<|end_header_id|>")[1].strip()

# Baseline Prediction
baseline_prompt = baseline_prompt_template.format(article_text=input_text)
baseline_response = baseline_model.generate_content(baseline_prompt)

results.append({
    "input": input_text,
    "ground_truth": ground_truth_output,
    "karios_prediction_str": karios_response_text,
    "baseline_prediction_str": baseline_response.text
})

# --- CALCULATE METRICS & SAVE ---
rouge = evaluate.load('rouge')

# Parse predictions and references for scoring
karios_summaries = []
baseline_summaries = []
reference_summaries = []

for res in results:
    try:
        karios_pred_json = json.loads(res['karios_prediction_str'])
        karios_summaries.append(karios_pred_json.get('event_summary', ''))
    except json.JSONDecodeError:
        karios_summaries.append('') # Penalize if output is not valid JSON

    try:
        # Gemini often outputs markdown, so we clean it before parsing
        cleaned_baseline_str = res['baseline_prediction_str'].replace('```json', '').replace('```', '')
        baseline_pred_json = json.loads(cleaned_baseline_str)
        baseline_summaries.append(baseline_pred_json.get('event_summary', ''))
    except json.JSONDecodeError:
        baseline_summaries.append('')

    reference_summaries.append(res['ground_truth'].get('event_summary', ''))

# Calculate ROUGE scores
karios_rouge = rouge.compute(predictions=karios_summaries, references=reference_summaries)
baseline_rouge = rouge.compute(predictions=baseline_summaries, references=reference_summaries)

print("\n--- EVALUATION RESULTS (ROUGE Score for Event Summary) ---")
print(f"KAIROS (Fine-Tuned Llama 3): {karios_rouge}")
print(f"Baseline (Gemini 1.5 Pro Zero-Shot): {baseline_rouge}")

# Save detailed results for qualitative analysis
with open('final_evaluation_results.jsonl', 'w') as f:
    for res in results:
        f.write(json.dumps(res) + '\n')

print("\n✅ Evaluation complete. Detailed side-by-side results saved to 'final_evaluation_results.jsonl'")

```

↻ Loading your fine-tuned KAIROS model...

```
config.json: 100% 654/654 [00:00<00:00, 38.5kB/s]

model.safetensors.index.json: 100% 23.9k/23.9k [00:00<00:00, 761kB/s]

Fetching 4 files: 100% 4/4 [26:02<00:00, 1562.39s/it]

model-00003-of-00004.safetensors: 100% 4.92G/4.92G [24:47<00:00, 12.5MB/s]

model-00004-of-00004.safetensors: 100% 1.17G/1.17G [24:45<00:00, 1.45MB/s]

model-00001-of-00004.safetensors: 100% 4.98G/4.98G [26:02<00:00, 97.4MB/s]

model-00002-of-00004.safetensors: 100% 5.00G/5.00G [24:46<00:00, 85.7MB/s]

Loading checkpoint shards: 100% 4/4 [01:15<00:00, 16.53s/it]

generation_config.json: 100% 187/187 [00:00<00:00, 16.0kB/s]

tokenizer_config.json: 100% 51.0k/51.0k [00:00<00:00, 5.02MB/s]

tokenizer.json: 100% 9.09M/9.09M [00:00<00:00, 28.9MB/s]

special_tokens_map.json: 100% 73.0/73.0 [00:00<00:00, 5.57kB/s]

adapter_config.json: 100% 869/869 [00:00<00:00, 59.0kB/s]

adapter_model.safetensors: 100% 54.6M/54.6M [00:00<00:00, 68.9MB/s]
```

✅ KAIROS Model loaded.

--- Generating predictions from KAIROS and Baseline models ---

```
0%|          | 0/10 [00:00<?, ?it/s]The following generation flags are not valid and may be ignored: ['temperature', 'top_
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
10%|█         | 1/10 [00:51<07:44, 51.61s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
20%|██        | 2/10 [01:44<06:58, 52.28s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
30%|███       | 3/10 [02:28<05:41, 48.78s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
40%|████      | 4/10 [03:21<05:00, 50.15s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
50%|█████     | 5/10 [04:10<04:09, 49.97s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
60%|██████    | 6/10 [04:58<03:16, 49.24s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
70%|███████   | 7/10 [05:49<02:29, 49.71s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
80%|████████  | 8/10 [06:33<01:36, 48.04s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
90%|█████████ | 9/10 [07:23<00:48, 48.46s/it]The following generation flags are not valid and may be ignored: ['temperature
Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.
100%|██████████| 10/10 [08:11<00:00, 49.14s/it]

Downloading builder script: 6.27k/? [00:00<00:00, 540kB/s]
```

--- EVALUATION RESULTS (ROUGE Score for Event Summary) ---

```
KAIROS (Fine-Tuned Llama 3): {'rouge1': np.float64(0.0), 'rouge2': np.float64(0.0), 'rougeL': np.float64(0.0), 'rougeLsum':
Baseline (Gemini 1.5 Pro Zero-Shot): {'rouge1': np.float64(0.0), 'rouge2': np.float64(0.0), 'rougeL': np.float64(0.0), 'roug
```

# --- CELL 1: Install All Dependencies ---

!pip install -q -U sentence-transformers scikit-learn pandas textstat evaluate

```
↻ _____ 91.2/91.2 kB 7.0 MB/s eta 0:00:00
_____ 470.2/470.2 kB 12.9 MB/s eta 0:00:00
_____ 12.9/12.9 MB 33.7 MB/s eta 0:00:00
_____ 12.4/12.4 MB 34.2 MB/s eta 0:00:00
_____ 175.3/175.3 kB 14.6 MB/s eta 0:00:00
_____ 939.4/939.4 kB 23.4 MB/s eta 0:00:00
_____ 2.1/2.1 MB 28.7 MB/s eta 0:00:00
```

```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is
google-colab 1.0.0 requires pandas==2.2.2, but you have pandas 2.3.0 which is incompatible.
sklearn-compat 0.1.3 requires scikit-learn<1.7,>=1.2, but you have scikit-learn 1.7.0 which is incompatible.
cudf-cu12 25.2.1 requires pandas<2.2.4dev0,>=2.0, but you have pandas 2.3.0 which is incompatible.
dask-cudf-cu12 25.2.2 requires pandas<2.2.4dev0,>=2.0, but you have pandas 2.3.0 which is incompatible.
```

# --- CELL 2: Run the Full Objective Evaluation ---

```
import json
import pandas as pd
import numpy as np
from sentence_transformers import SentenceTransformer, util
from sklearn.metrics import accuracy_score
import textstat
```

```
import evaluate
import re
```

```
# --- 1. SETUP & LOAD DATA ---
```

```
print("--- Step 1: Installing dependencies and loading data ---")
```

```
!pip install -q -U sentence-transformers scikit-learn pandas textstat evaluate
```

```
with open('final_evaluation_results.jsonl', 'r') as f:
```

```
    results = [json.loads(line) for line in f]
```

```
print(f"✅ Loaded {len(results)} results to analyze.")
```



```
--- Step 1: Installing dependencies and loading data ---
```

```
✅ Loaded 10 results to analyze.
```

```
# --- 2. DEFINE THE FINAL ROBUST PARSERS ---
```

```
print("\n--- Step 2: Defining specialized parsers for each model's output ---")
```

```
def parse_kairos_output(pred_str):
```

```
    """
```

```
    FIX: This parser correctly handles the KAIROS output by using regex
    on the formatted text string, not by trying to parse it as JSON.
```

```
    """
```

```
    try:
```

```
        summary_match = re.search(r"Event Summary:\s*(.*)", pred_str)
```

```
        type_match = re.search(r"Event Type:\s*(.*)", pred_str)
```

```
        summary = summary_match.group(1).strip() if summary_match else ''
```

```
        event_type = type_match.group(1).strip() if type_match else ''
```

```
        return {'event_summary': summary, 'event_type': event_type, 'is_valid': True}
```

```
    except Exception:
```

```
        return {'event_summary': '', 'event_type': '', 'is_valid': False}
```

```
def parse_gemini_output(pred_str):
```

```
    """
```

```
    FIX: This parser uses a more flexible regex to handle the Gemini
    markdown output correctly.
```

```
    """
```

```
    try:
```

```
        # Case-insensitive search for 'summary' and 'type' between markdown bold tags
```

```
        summary_match = re.search(r"\s*\s*Summary:\s*\s*(.*)", pred_str, re.IGNORECASE)
```

```
        type_match = re.search(r"\s*\s*Type:\s*\s*(.*)", pred_str, re.IGNORECASE)
```

```
        summary = summary_match.group(1).strip() if summary_match else ''
```

```
        event_type = type_match.group(1).strip() if type_match else ''
```

```
        return {'event_summary': summary, 'event_type': event_type, 'is_valid': True}
```

```
    except Exception:
```

```
        return {'event_summary': '', 'event_type': '', 'is_valid': False}
```



```
--- Step 2: Defining specialized parsers for each model's output ---
```

```
# --- 3. PROCESS RESULTS USING CORRECTED PARSERS ---
```

```
print("\n--- Step 3: Reprocessing results with robust parsers ---")
```

```
kairos_metrics = {'json_valid': [], 'event_type_preds': [], 'summaries': []}
```

```
baseline_metrics = {'event_type_preds': [], 'summaries': []}
```

```
reference_metrics = {'event_type': [], 'summaries': []}
```

```
for res in results:
```

```
    kairos_parsed = parse_kairos_output(res['kairos_prediction_str'])
```

```
    baseline_parsed = parse_gemini_output(res['baseline_prediction_str'])
```

```
    kairos_metrics['json_valid'].append(1 if kairos_parsed.get('is_valid') else 0)
```

```
    kairos_metrics['summaries'].append(kairos_parsed.get('event_summary'))
```

```
    kairos_metrics['event_type_preds'].append(kairos_parsed.get('event_type'))
```

```
    baseline_metrics['summaries'].append(baseline_parsed.get('event_summary'))
```

```
    baseline_metrics['event_type_preds'].append(baseline_parsed.get('event_type'))
```

```
    reference_metrics['summaries'].append(res['ground_truth'].get('event_summary', ''))
```

```
    reference_metrics['event_type'].append(res['ground_truth'].get('event_type', ''))
```



--- Step 3: Reprocessing results with robust parsers ---

# --- 4. CALCULATE FINAL SCORES ---

print("\n--- Step 4: Aggregating final scores ---")

# Semantic Similarity

embedding\_model = SentenceTransformer('all-MiniLM-L6-v2')

ref\_embeddings = embedding\_model.encode(reference\_summaries, convert\_to\_tensor=True)

karios\_embeddings = embedding\_model.encode(karios\_metrics['summaries'], convert\_to\_tensor=True)

baseline\_embeddings = embedding\_model.encode(baseline\_metrics['summaries'], convert\_to\_tensor=True)

karios\_cosine\_scores = util.cos\_sim(karios\_embeddings, ref\_embeddings)

baseline\_cosine\_scores = util.cos\_sim(baseline\_embeddings, ref\_embeddings)

# Task Adherence & Accuracy

metrics\_results = {

    "KAİROS (Fine-Tuned Llama 3)": {

        # Apply .cpu() to each element before converting to numpy for mean calculation

        "Semantic Similarity (Cosine)": np.mean([karios\_cosine\_scores[i][i].cpu() for i in range(len(karios\_cosine\_scores))]).it

        "JSON Validity (%)": np.mean(karios\_metrics['json\_valid']) \* 100,

        # "Schema Adherence (%)": np.mean(karios\_metrics['schema\_adherence']) \* 100, # schema\_adherence metric not calculated

        "Event Type Accuracy (%)": accuracy\_score(reference\_metrics['event\_type'], karios\_metrics['event\_type\_preds']) \* 100,

        "Readability (Flesch Score)": np.mean([textstat.flesch\_reading\_ease(s) for s in karios\_metrics['summaries'] if s]),

    },

    "Baseline (Gemini 1.5 Pro)": {

        # Apply .cpu() to each element before converting to numpy for mean calculation

        "Semantic Similarity (Cosine)": np.mean([baseline\_cosine\_scores[i][i].cpu() for i in range(len(baseline\_cosine\_scores))])

        # "JSON Validity (%)": np.mean(baseline\_metrics['json\_valid']) \* 100, # json\_valid metric not calculated for baseline

        # "Schema Adherence (%)": np.mean(baseline\_metrics['schema\_adherence']) \* 100, # schema\_adherence metric not calculated

        "Event Type Accuracy (%)": accuracy\_score(reference\_metrics['event\_type'], baseline\_metrics['event\_type\_preds']) \* 100,

        "Readability (Flesch Score)": np.mean([textstat.flesch\_reading\_ease(s) for s in baseline\_metrics['summaries'] if s]),

    }

}



--- Step 4: Aggregating final scores ---

# --- 5. DISPLAY PUBLICATION-READY TABLE ---

results\_df = pd.DataFrame(metrics\_results).round(2)

print("\n\n--- PUBLICATION-READY RESULTS TABLE ---")

print(results\_df)

print("\n\n--- ANALYSIS OF RESULTS ---")

print("Higher is better for all metrics.")

print("- \*\*Semantic Similarity:\*\* Measures how close in \*meaning\* the summary is to the ground truth. Your model should now be c

print("- \*\*JSON Validity:\*\* Measures if the model returned clean, usable JSON. This is a key metric where your fine-tuned model

print("\n✅ Objective evaluation complete.")



--- PUBLICATION-READY RESULTS TABLE ---

	KAİROS (Fine-Tuned Llama 3) \
Semantic Similarity (Cosine)	0.90
JSON Validity (%)	100.00
Event Type Accuracy (%)	30.00
Readability (Flesch Score)	29.93

	Baseline (Gemini 1.5 Pro)
Semantic Similarity (Cosine)	0.82
JSON Validity (%)	NaN
Event Type Accuracy (%)	0.00
Readability (Flesch Score)	28.61

--- ANALYSIS OF RESULTS ---

Higher is better for all metrics.

- \*\*Semantic Similarity:\*\* Measures how close in \*meaning\* the summary is to the ground truth. Your model should now be comp

- \*\*JSON Validity:\*\* Measures if the model returned clean, usable JSON. This is a key metric where your fine-tuned model sho

✅ Objective evaluation complete.

Start coding or [generate](#) with AI.

