# Final Team Project: Multi-Agent Financial Analysis System

## Group 5

Team Members:

Andrew Blumhardt

Surya Prakash Kuppa Srinivasan

Saman Tavasoli

GitHub Link: https://github.com/prakashsurya-840/Group5AI520FinalProject

## Overall Project Summary

We will be using two main libraries for the analysis,

- yfinance
- openAI

We will be running the workflows using the 2 libraries by selecting the same grounds and will be comparing the results.

⌄

## Analysis using yFinance

⌄

### Step 1: Environment & Notebook Setup

What do we do in this step?

- Create the virtual environment preferrably with the python version 3.10.18
    - If you are using conda environment, run the following command
        - `conda create -n myenv python=3.10.18`
        - `conda activate myenv`
- Create a clean working notebook and make sure required libraries are installed (please run requirements.txt to install all the required dependencies) and imported without errors.
    - Run `pip install -r requirements.txt`
- Verify versions
- Add a few display and pandas settings (plot display, pandas options)

```
# Step 1
# In this step, we have Imports, Versions, and Display settings.

# Imports

# Import core modules for system info and OS details
import sys
import platform

# Import data libraries
import pandas as pd
import numpy as np

# Import Plotting libraries
import matplotlib
import matplotlib.pyplot as plt

# Import finance/ data access libraries
import yfinance as yf
import requests


# Print versions so all the team members know the exact environment
print("Python:", sys.version.split()[0])
print("OS:", platform.system(), platform.release())
print("pandas:", pd.__version__)
print("numpy:", np.__version__)
print("matplotlib:", matplotlib.__version__)
print("yfinance:", yf.__version__)

# Set a random seed to make any random behavior repeatable
np.random.seed(42)


# Display settings so the PDF/HTML export looks clean
# Here, we make DataFrame display compact so the PDF export remains readable.
pd.set_option("display.max_rows", 50)
pd.set_option("display.max_columns", 20)
pd.set_option("display.width", 120)

# Telling the Jupyter to draw Matplotlib charts right inside the output cell
%matplotlib inline
```

```
    Python: 3.12.12
    OS: Linux 6.6.105+
    pandas: 2.2.2
    numpy: 2.0.2
    matplotlib: 3.10.0
    yfinance: 0.2.66
```

∨

## Step 2: Settings

What do we do in this step?

> This is where we set up all the basic information our program needs to run. By putting everything here, we can easily change
> a setting, and the whole program will use that new value.

Here are the specific settings we need to provide:

- ticker: This is the abbreviation (symbol) for the stock we want to analyze (like "GOOG" for Google or "AAPL" for Apple).
- start_date/end_date: These tell the program the time frame for the data it should look at. We need to write the dates in the format:
  Year-Month-Day (like 2024-01-15).
- news_source: This tells the program where to get its news from. It can be a file we already have, or it can connect to an online source
  through API.

  - Our API key is: fd788bb137fd4b47843696dd6bb987b1
  - d652b5dcfe8d40e9a85b4b579e6ccfed

- memory_file: This is a small file where our program (the "agent") can save important notes between the times we run it. This helps it
  remember things.
- random_seed: This is a special number that makes sure that any random actions the program takes are the same every time we run
  the code. This is important for getting the same, predictable results repeatedly.

```python
# Step 2: Settings
# In this step, we set the choices(ticker, dates, news mode)

"""
For our data, we use two things:
    1- Stock prices (open, high, low, close, volume) from Yahoo Finance using the yfinance library.
    2- News about the company(either from a local file we already have or from an online API)

We decide which data to use in this step,
    TICKER  -->  AAPL
    START_DATE and END_DATE  -->  The date range for data
    NEW_MODE  -->  How to get news:
        "file" -> Read from a local JSON file
        "api"  -> Pull from the internet
        "none" --> skip news
"""

# We put all project settings here
import os
from datetime import datetime, date


os.environ["NEWSAPI_KEY"] = "fd788bb137fd4b47843696dd6bb987b1"

# Settings DICTIONARY
PROJECT_SETTINGS = {
    # The stock we're analyzing, for example: Apple
    "TICKER": "AAPL",

    # Date range for our analysis, the format will be YYYY-MM-DD
    "START_DATE": "2025-05-16",
    "END_DATE":   "2025-10-17",

    # Where should we get the news? Options: "file", "api", "none". We are working on this, next week we will decide on this
    "NEWS_MODE": "api",

    # If NEWS_MODE was 'file', we set the path to our JSON news file here
    "LOCAL_NEWS_FILE": "data/news_sample.json",

    # If NEWS_MODE was 'api', we set the service and the environment variable name for our key
    "API_SERVICE_NAME": "newsapi",
    "API_KEY_ENV_VAR": "NEWSAPI_KEY",

    # The file where our 'agent' saves its notes/memory between runs
    "AGENT_MEMORY_FILE": "agent_memory.json",

    # A number to keep any random processes consistent
    "RANDOM_SEED": 42,
}


"""
    Checks that the start and end date strings are in the correct format,
    ensures the start date is not after the end date, and warns if the end date is in the future.
"""
def check_dates(start_str: str, end_str: str):
    # Simple check to make sure the dates are valid and make sense.
    date_format = "%Y-%m-%d"

    # Check 1: Can we actually read the dates?
    try:
        start_date = datetime.strptime(start_str, date_format).date()
        end_date = datetime.strptime(end_str, date_format).date()
```

```
        except ValueError as e:
            # Give a clear error message
            raise ValueError(f"Check your date format. It must be YYYY-MM-DD. Error details: {e}")

        # Check 2: Does the start date come before the end date?
        if start_date > end_date:
            raise ValueError(f"The start date ({start_date}) can't be after the end date ({end_date}). Fix it!")

        # Check 3: Is the end date in the future? (It's okay, but you should know!)
        if end_date > date.today():
            print(f"Your end date ({end_date}) is in the future. Price/news data might be missing until you move it back to today or ea

        return start_date, end_date


    #  Here, we check our work and print the final settings.
    # Checking the dates right away
    try:
        start_d, end_d = check_dates(PROJECT_SETTINGS["START_DATE"], PROJECT_SETTINGS["END_DATE"])
        print(f"Dates are set: {start_d} to {end_d}")
    except ValueError as e:
        print(f"ERROR: {e}")

    # Checking the news source setup
    news_mode = PROJECT_SETTINGS.get("NEWS_MODE", "none").lower()
    key_env = PROJECT_SETTINGS.get("API_KEY_ENV_VAR")

    if news_mode == "file":
        path = PROJECT_SETTINGS.get("LOCAL_NEWS_FILE")
        if not os.path.exists(path):
            print(f"Local file was not found at '{path}'. We will try to keep running, but you need to check the path!")
        else:
            print(f"News will be loaded from local file: {path}")

    elif news_mode == "api":
        if os.environ.get(key_env):
            print(f"API mode is ON. Your key was found in the environment variable: ${key_env}")
        else:
            print(f"API mode selected, but the key wasn't found in ${key_env}. You to set your key before we run the API calls!")

    elif news_mode == "none":
        print("Skipping news collection entirely for now (NEWS_MODE='none').")

    else:
        print(f"Unknown NEWS_MODE: '{news_mode}'. Please use 'file', 'api', or 'none'.")

    # Printing all final settings
    print("\n ****** Final Project Settings ******")
    for key, value in PROJECT_SETTINGS.items():
        print(f" - {key}: {value}")
    print("\n\n\n")
```

```
Dates are set: 2025-05-16 to 2025-10-17
API mode is ON. Your key was found in the environment variable: $NEWSAPI_KEY

 ****** Final Project Settings ******
 - TICKER: AAPL
 - START_DATE: 2025-05-16
 - END_DATE: 2025-10-17
 - NEWS_MODE: api
 - LOCAL_NEWS_FILE: data/news_sample.json
 - API_SERVICE_NAME: newsapi
 - API_KEY_ENV_VAR: NEWSAPI_KEY
 - AGENT_MEMORY_FILE: agent_memory.json
 - RANDOM_SEED: 42
```

## Step 3: Tools Layer (functions to fetch prices/profile/news; with basic error handling)

Here is like our utility belt, it is a section where we build custom functions that do specific, repeatable works for us. Instead of writing the same complicated code every time we need to download stock prices, we write it once, and we just need to call the function whenever we need it. In general, we can say that in this step, we fetch the data. Here, in this step, we create 3 main tools:

- "Get Prices" Tool: This function automatically reaches out to a service like Yahoo Finance, and downloads the daily stock prices.
- "Get Profile" Tool: This function grabs the simple facts about the company, like its industry and website.
- "Get New" Tool: This function checks the settings in step 2 and find the news, either by openitng a local file on our system, or connecting us to an online news service.

```
# Step 3: Tools Layer

"""
We have the below functions in this step:
 - get_price_data(ticker, start, end)      → DataFrame with daily OHLCV
 - get_company_profile(ticker)             → dict with basic company info
 - get_news(project_settings)              → list of news items (file/api/none)
 - save_prices_to_csv(df, path)            → save prices DataFrame
 - save_news_to_json(items, path)          → save news list
 - ensure_data_dir(dir_path)               → makes sure data/ exists
"""

from pathlib import Path
import os
```

```python
import json
from typing import Optional, List, Dict
import yfinance as yf
import math
import time
import requests
from datetime import datetime

"""
    We need to make sure the data directory exists so we can save files there safely.
    Returns the Path object for convenience.
"""
def ensure_data_dir(dir_path: str = "data") -> Path:

    p = Path(dir_path)
    p.mkdir(parents=True, exist_ok=True)
    print(f"[info] data directory ready at: {p.resolve()}")
    return p


"""
    Downloads daily price data (Open, High, Low, Close, Adj Close, Volume) using yfinance.
    Returns a pandas DataFrame, or None if something goes wrong.
"""
def get_price_data(ticker: str, start_date: str, end_date: str):

    try:
        print(f"[info] downloading price data for {ticker} from {start_date} to {end_date} ...")
        df = yf.download(ticker, start=start_date, end=end_date, interval="1d", progress=False)

        # If nothing came back, let the user know and return None
        if df is None or df.empty:
            print(f"[warn] No price data returned. Check the ticker ('{ticker}') or the date range.")
            return None

        # Put Date into a normal column for easy saving/printing later
        df = df.reset_index()

        # Keep the columns in a predictable order if they exist
        expected_cols = ["Date", "Open", "High", "Low", "Close", "Adj Close", "Volume"]
        df = df[[c for c in expected_cols if c in df.columns]]

        # Convert Date to string for nice CSV output
        if "Date" in df.columns:
            df["Date"] = df["Date"].astype(str)

        print(f"[ok] downloaded {len(df)} rows of prices for {ticker}.")
        return df

    except Exception as e:
        print(f"[error] Could not download price data for {ticker}. Details: {e}")
        return None


"""
    Reading basic company info from yfinance.
    We try .get_info() (newer), then fall back to .info (older), and finally to fast_info.
    Returns a small dict with what we can find.
"""
def get_company_profile(ticker: str) -> Dict[str, Optional[str]]:

    info = {}
    try:
        t = yf.Ticker(ticker)

        # Try the newer API first
        try:
            raw = t.get_info()
        except Exception:
            # Fall back to the older attribute
            try:
                raw = t.info
            except Exception:
                raw = {}

        # If still empty, try fast_info for a few useful fields
        if not raw:
            try:
                fi = t.fast_info
            except Exception:
                fi = {}

            info = {
                "symbol": ticker,
                "shortName": None,
                "longName": None,
                "sector": None,
                "industry": None,
                "website": None,
                "currency": getattr(fi, "currency", None) if hasattr(fi, "currency") else None,
            }
            print("[warn] Could not get full profile; returned minimal info.")
            return info

        # Safely pick common fields (some may be missing depending on ticker)
        info = {
            "symbol": ticker,
```

```python
            "shortName": raw.get("shortName"),
            "longName": raw.get("longName"),
            "sector": raw.get("sector"),
            "industry": raw.get("industry"),
            "website": raw.get("website"),
            "country": raw.get("country"),
            "fullTimeEmployees": raw.get("fullTimeEmployees"),
            "currency": raw.get("currency"),
        }
        print(f"[ok] got company profile for {ticker}.")
        return info

    except Exception as e:
        print(f"[warn] Could not read company profile for {ticker}. Details: {e}")
        # Return at least the ticker so callers have something
        return {"symbol": ticker}



# for NewsAPI fetching and normalization
# Making each NewsAPI article a small, consistent dict that we can save and use later.
def _normalize_newsapi_item(item):

    return {
        "title": item.get("title"),
        "description": item.get("description"),
        "url": item.get("url"),
        "source": (item.get("source") or {}).get("name"),
        "published_at": item.get("publishedAt"),
        "author": item.get("author"),
        "content": item.get("content"),
    }


"""
    Very simple client for NewsAPI /v2/everything.
    - Searches by ticker (e.g., "AAPL")
    - Pulls a few pages (pageSize=100 each)
    - Returns a list of normalized articles
"""
def fetch_news_from_newsapi(ticker, start_date, end_date, api_key, max_pages=3):

    url = "https://newsapi.org/v2/everything"
    headers = {"X-Api-Key": api_key}

    base_params = {
        "q": ticker,
        # "from": start_date,
        # "to": end_date,
        "language": "en",
        "sortBy": "publishedAt",
        "pageSize": 100,
    }

    all_items = []
    page = 1
    while page <= max_pages:
        params = {**base_params, "page": page}
        try:
            resp = requests.get(url, headers=headers, params=params, timeout=20)
            if resp.status_code != 200:
                print(f"[warn] NewsAPI status {resp.status_code}: {resp.text[:200]}...")
                break

            data = resp.json()
            arts = data.get("articles", [])
            if not arts:
                break

            all_items.extend([_normalize_newsapi_item(a) for a in arts])

            total = data.get("totalResults", 0)
            total_pages = math.ceil(min(total, 1000) / base_params["pageSize"])
            print(f"[info] NewsAPI page {page}/{min(total_pages, max_pages)} – got {len(arts)} items")

            page += 1
            time.sleep(0.5)

        except Exception as e:
            print(f"[warn] NewsAPI request failed: {e}")
            break

    print(f"total news items fetched: {len(all_items)}")
    return all_items



"""
    Loads news based on PROJECT_SETTINGS['NEWS_MODE'].
    - 'file' → read a local JSON file (list of items)
    - 'api'  → (placeholder) call our news API using an API key
    - 'none' → return empty list
    It always returns a list even empty.
"""
def get_news(project_settings: dict) -> List[dict]:

    mode = project_settings.get("NEWS_MODE", "none").lower()
    ticker = project_settings.get("TICKER", "UNKNOWN")
```

```python
    if mode == "file":
        local_path = project_settings.get("LOCAL_NEWS_FILE", "")
        print(f"[info] NEWS_MODE='file' → loading from: {local_path}")
        if not local_path or not os.path.exists(local_path):
            print("[warn] Local news file not found. Returning empty list for now.")
            return []
        try:
            with open(local_path, "r", encoding="utf-8") as f:
                data = json.load(f)
            if not isinstance(data, list):
                print("[warn] News JSON was not a list. Wrapping it into a list.")
                data = [data]
            print(f"[ok] loaded {len(data)} news items from file.")
            return data
        except Exception as e:
            print(f"[warn] Could not read/parse news file. Details: {e}")
            return []

    elif mode == "api":
        service = project_settings.get("API_SERVICE_NAME", "newsapi").lower()
        key_env = project_settings.get("API_KEY_ENV_VAR", "")
        api_key = os.environ.get(key_env)

        if not api_key:
            print(f"[warn] NEWS_MODE='api' but no API key found in ${key_env}. Returning empty list.")
            return []

        ticker = project_settings.get("TICKER", "UNKNOWN")
        start_date = project_settings.get("START_DATE")
        end_date   = project_settings.get("END_DATE")
        print(f"[info] NEWS_MODE='api' → service='{service}', ticker='{ticker}', range={start_date}→{end_date}")

        if service == "newsapi":
            items = fetch_news_from_newsapi(ticker, start_date, end_date, api_key, max_pages=3)
            return items
        else:
            print(f"[warn] Unknown API service: {service}. Returning empty list.")
            return []

    else:
        print("[info] NEWS_MODE='none' → skipping news. Returning empty list.")
        return []


def save_prices_to_csv(df, path: Path) -> bool:
    """
    Saves the prices DataFrame to CSV.
    Returns True if saved, False otherwise.
    """
    try:
        path = Path(path)
        df.to_csv(path, index=False)
        print(f"[ok] saved prices to: {path}")
        return True
    except Exception as e:
        print(f"[error] Could not save prices CSV. Details: {e}")
        return False


def save_news_to_json(items: List[dict], path: Path) -> bool:
    """
    Saves the news list to a JSON file (pretty-printed).
    Returns True if saved, False otherwise.
    """
    try:
        path = Path(path)
        with open(path, "w", encoding="utf-8") as f:
            json.dump(items, f, ensure_ascii=False, indent=2)
        print(f"[ok] saved news to: {path}")
        return True
    except Exception as e:
        print(f"[error] Could not save news JSON. Details: {e}")
        return False
```

```python
# Step 3 runner
# Here we download the AAPL prices, fetch the company info, anc call our NewsAPI using our key, and save both datasets into the dat

from pathlib import Path
import json

# Ensuring data directory
data_dir = ensure_data_dir("data")

# Getting prices
prices_df = get_price_data(
    PROJECT_SETTINGS["TICKER"],
    PROJECT_SETTINGS["START_DATE"],
    PROJECT_SETTINGS["END_DATE"]
)

# Define output file names
prices_csv_path = data_dir / f"{PROJECT_SETTINGS['TICKER']}_prices_{PROJECT_SETTINGS['START_DATE']}_{PROJECT_SETTINGS['END_DATE']}.
news_json_path  = data_dir / f"{PROJECT_SETTINGS['TICKER']}_news_{PROJECT_SETTINGS['START_DATE']}_{PROJECT_SETTINGS['END_DATE']}.js

# Getting company profile
profile = get_company_profile(PROJECT_SETTINGS["TICKER"])
```

```python
    # Getting news (now using your active NewsAPI)
    news_items = get_news(PROJECT_SETTINGS)

    # Saving files
    if prices_df is not None:
        save_prices_to_csv(prices_df, prices_csv_path)

    save_news_to_json(news_items, news_json_path)

    # Printing summary
    print("\n===== SUMMARY =====")
    print(f"TICKER: {PROJECT_SETTINGS['TICKER']}")
    print(f"DATE RANGE: {PROJECT_SETTINGS['START_DATE']} → {PROJECT_SETTINGS['END_DATE']}")
    print(f"Prices rows: {0 if prices_df is None else len(prices_df)}")
    print(f"News items: {len(news_items)}")

    if prices_df is not None:
        print("\nprices_df.head():")
        print(prices_df.head().to_string(index=False))

    print("\nCompany profile (key fields):")
    for k in ["symbol","shortName","longName","sector","industry","website","country","fullTimeEmployees","currency"]:
        print(f" - {k}: {profile.get(k)}")

    print("\nSaved files:")
    print(f" - Prices CSV: {prices_csv_path}")
    print(f" - News JSON : {news_json_path}")


    # print(

    # """
    # The output tells us:
    # It downloaded 14 rows of prices, therefore Yahoo Finance is working correctly
    # We got company profile for AAPL, therefore, the profile info loaded successfully
    # The NewsAPI status 425, tells us the free NewsAPI plan only allows data from Sept 9 to today, the earlier days are blocked
    # Everything in CSV and JSON is saved and working as we wanted.

    # """)
    # print("\n\n\n")
```

```
[info] data directory ready at: /content/data
[info] downloading price data for AAPL from 2025-05-16 to 2025-10-17 ...
[ok] downloaded 106 rows of prices for AAPL.
/tmp/ipython-input-2498583754.py:44: FutureWarning: YF.download() has changed argument auto_adjust default to True
  df = yf.download(ticker, start=start_date, end=end_date, interval="1d", progress=False)
[ok] got company profile for AAPL.
[info] NEWS_MODE='api' → service='newsapi', ticker='AAPL', range=2025-05-16→2025-10-17
[info] NewsAPI page 1/3 — got 100 items
[warn] NewsAPI status 426: {"status":"error","code":"maximumResultsReached","message":"You have requested too many results. Develope
total news items fetched: 100
[ok] saved prices to: data/AAPL_prices_2025-05-16_2025-10-17.csv
[ok] saved news to: data/AAPL_news_2025-05-16_2025-10-17.json

===== SUMMARY =====
TICKER: AAPL
DATE RANGE: 2025-05-16 → 2025-10-17
Prices rows: 106
News items: 100

prices_df.head():
      Date       Open       High        Low      Close    Volume
                 AAPL       AAPL       AAPL       AAPL      AAPL
2025-05-16 212.119267 212.329036 209.532207 211.020508 54737900
2025-05-19 207.674311 209.242523 204.028439 208.543320 46140500
2025-05-20 207.434583 208.233679 204.797576 206.625504 42496600
2025-05-21 204.937411 206.805286 200.482476 201.860901 59211800
2025-05-22 200.482472 202.520153 199.473607 201.131729 46742400

Company profile (key fields):
 - symbol: AAPL
 - shortName: Apple Inc.
 - longName: Apple Inc.
 - sector: Technology
 - industry: Consumer Electronics
 - website: https://www.apple.com
 - country: United States
 - fullTimeEmployees: 150000
 - currency: USD

Saved files:
 - Prices CSV: data/AAPL_prices_2025-05-16_2025-10-17.csv
 - News JSON : data/AAPL_news_2025-05-16_2025-10-17.json
```

## Step 4: Agent Memory (load/save a JSON memory; create file if missing)

In this step, we give our project a little memory so it can remember notes between runs. This memory will be stored in a JSON file on disk.
If the file doesn't exist, we create it with a simple default structure.

```python
# Step 4: Agent Memory

"""
Here, we have three functions
 - load_memory(path)  -> returns a dict (creates file if missing)
 - save_memory(mem)   -> writes the dict back to the same file
 - add_memory_note(...) -> appends a timestamped note

"""
```

```python
import os
import json
from pathlib import Path
from datetime import datetime

# We will use the path from our Step 2 settings.
MEMORY_PATH = Path(PROJECT_SETTINGS["AGENT_MEMORY_FILE"]).resolve()


"""
    If the memory file does not exist yet, this will be the starting shape.
"""
def _default_memory():

    return {
        "project": "AAI-520 Financial Agent",
        "ticker": PROJECT_SETTINGS.get("TICKER", "UNKNOWN"),
        "created_at": datetime.utcnow().isoformat(timespec="seconds") + "Z",
        "last_updated_at": None,          # we will fill this when we save
        "notes": [],                      # quick log messages we add over time
        "runs": 0                         # how many times we've used this pipeline
    }




"""
    Load the memory JSON file from disk.
    - If the file doesn't exist, create it with default content.
    - If the file is not valid JSON, we start fresh and warn the user.
      Always returns a dict.
"""
def load_memory(path: Path = MEMORY_PATH) -> dict:

    try:
        if not os.path.exists(path):
            print(f"[info] Memory file not found. Creating a new one at: {path}")
            mem = _default_memory()
            save_memory(mem, path)
            return mem

        with open(path, "r", encoding="utf-8") as f:
            mem = json.load(f)

        # Ensure minimal keys exist even if the file was edited manually
        if "notes" not in mem:
            mem["notes"] = []
        if "runs" not in mem:
            mem["runs"] = 0
        if "ticker" not in mem:
            mem["ticker"] = PROJECT_SETTINGS.get("TICKER", "UNKNOWN")

        return mem

    except Exception as e:
        print(f"[warn] Could not read memory file (using a fresh one). Details: {e}")
        mem = _default_memory()
        save_memory(mem, path)
        return mem




"""
    Save the memory dict back to disk as pretty JSON.
    We also update the 'last_updated_at' timestamp each time we save.
    Returns True/False to tell us if it worked.
"""
def save_memory(mem: dict, path: Path = MEMORY_PATH) -> bool:

    try:
        mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
        with open(path, "w", encoding="utf-8") as f:
            json.dump(mem, f, ensure_ascii=False, indent=2)
        print(f"[ok] Memory saved to: {path}")
        return True
    except Exception as e:
        print(f"[error] Could not save memory. Details: {e}")
        return False




"""
    Quick helper to append a note with a timestamp.
    Use this anywhere in your notebook when something important happens.
"""
def add_memory_note(message: str, path: Path = MEMORY_PATH) -> None:

    mem = load_memory(path)
    mem["notes"].append({
        "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
        "msg": message
    })
    mem["runs"] = int(mem.get("runs", 0)) + 1


    MAX_NOTES = 500
    if len(mem["notes"]) > MAX_NOTES:
```

```
            mem["notes"] = mem["notes"][-MAX_NOTES:]
            print(f"[info] Trimmed old notes; keeping last {MAX_NOTES}.")


      save_memory(mem, path)

   # ---- Example usage (you can run these right away) ----
   print("\n[Step 4 quick test]")
   memory = load_memory()                        # 1) load or create the file
   add_memory_note("Initialized memory in Step 4.")  # 2) add a simple note
   print("Current memory snapshot:")
   print(json.dumps(load_memory(), indent=2))        # 3) print it nicely


   # print(
   # """
   # The output tells us:
   # Memeory file exists and saved correctly
   # It record notes properly
   # The runs counter increased from 1 to 2
   # All timestamps and fields look clean.
   # """)
   # print("\n\n\n")
```

```
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Step executed."
     },
     {
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:20Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:21Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:21Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:21Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:21Z",
       "msg": "Agent: Fetched 106 price rows successfully."
     },
     {
       "ts": "2025-10-20T22:05:21Z",
       "msg": "Agent dry run completed successfully."
     },
     {
       "ts": "2025-10-20T22:05:22Z",
       "msg": "Prompt chaining run completed on 100 news items."
     },
     {
       "ts": "2025-10-20T23:03:42Z",
       "msg": "Initialized memory in Step 4."
     }
   ],
   "runs": 18
 }
```

```
/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a f
   "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fu
   mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
```

## Step 5: Agent Class(plan_steps, act, self_reflect, learn)

> In this step, we create an AI Agent that can plan, act, think about what it did, and learn from that experience. In this step, we
> will build a small class called FinancialAgent with four simple functions:

- plan_steps() --> decide what to do next
- act() --> do one of those steps
- self_reflect() --> did something fail or succeed
- learn() --> save into memory what you observed or learned

This agent will use the PROJECT_SETTINGS from step 2 to know what stack to analyze, then it uses the Tools layer from step 3 to fetch the data, then it uses the memory from step 4 to save the notes and progress. Therefore, our Agent can remember what it has done and repeat it smartly in the future.

```
# Step 5: Agent Class
# Here, we will have a small brain that can plan, act, reflect, and learn.
# - plan_steps(): decide what to do next
# - act(): execute one step using our Step 3 tools
```

```python
# - self_reflect(): check the result and create a short reflection message
# - learn(): write a timestamped note to agent memory in step 4
# ===============================

from dataclasses import dataclass, field
from typing import List, Dict, Any, Optional
from pathlib import Path

@dataclass
class FinancialAgent:
    settings: Dict[str, Any]
    memory_path: Path
    plan: List[str] = field(default_factory=list)    # queue of step names
    last_result: Dict[str, Any] = field(default_factory=dict)  # store outputs to reflect/learn



    # *************************
    # PLAN
    # *************************
    """
    Very simple plan:
    - Always ensure data exists
    - Fetch prices for current ticker and date range
    - Fetch news according to NEWS_MODE
    - Save both to disk
    """
    def plan_steps(self):

        self.plan = [
            "ensure_data_dir",
            "fetch_prices",
            "fetch_profile",
            "fetch_news",
            "save_prices",
            "save_news",
            "summary_print"
        ]
        print("[plan] Steps:", " → ".join(self.plan))



    # *************************
    # ACT
    # *************************
    """
    Execute the next step in the plan if we had one.
    Stores outputs in self.last_result so reflect/learn can use them.
    """
    def act(self):

        if not self.plan:
            print("[act] No steps left in the plan.")
            return None

        step = self.plan.pop(0)
        print(f"[act] Running step: {step}")

        try:
            if step == "ensure_data_dir":
                data_dir = ensure_data_dir("data")
                self.last_result["data_dir"] = data_dir

                # precompute output paths for later steps
                self.last_result["prices_csv_path"] = data_dir / f"{self.settings['TICKER']}_prices_{self.settings['START_DATE']}_{
                self.last_result["news_json_path"]  = data_dir / f"{self.settings['TICKER']}_news_{self.settings['START_DATE']}_{se

            elif step == "fetch_prices":
                df = get_price_data(self.settings["TICKER"], self.settings["START_DATE"], self.settings["END_DATE"])
                self.last_result["prices_df"] = df

            elif step == "fetch_profile":
                prof = get_company_profile(self.settings["TICKER"])
                self.last_result["profile"] = prof

            elif step == "fetch_news":
                items = get_news(self.settings)
                self.last_result["news_items"] = items

            elif step == "save_prices":
                df = self.last_result.get("prices_df")
                path = self.last_result.get("prices_csv_path")
                ok = False
                if df is not None and path is not None:
                    ok = save_prices_to_csv(df, path)
                self.last_result["saved_prices"] = ok

            elif step == "save_news":
                items = self.last_result.get("news_items", [])
                path = self.last_result.get("news_json_path")
                ok = False
                if path is not None:
                    ok = save_news_to_json(items, path)
                self.last_result["saved_news"] = ok

            elif step == "summary_print":
                df = self.last_result.get("prices_df")
                items = self.last_result.get("news_items", [])
```

```python
            prof = self.last_result.get("profile", {})
            print("\n************** AGENT SUMMARY ***************")
            print(f"TICKER: {self.settings['TICKER']}")
            print(f"DATE RANGE: {self.settings['START_DATE']} → {self.settings['END_DATE']}")
            print(f"Prices rows: {0 if df is None else len(df)}")
            print(f"News items: {len(items)}")
            if df is not None:
                print("\nprices_df.head():")
                print(df.head().to_string(index=False))
            print("\nCompany profile (key fields):")
            for k in ["symbol","shortName","longName","sector","industry","website","country","fullTimeEmployees","currency"]:
                print(f"  - {k}: {prof.get(k)}")

        else:
            print(f"[act] Unknown step: {step}")

        # Return which step we executed (useful for logs)
        return step

    except Exception as e:
        # If any step fails, record error so reflect/learn can handle it
        self.last_result["error"] = f"Step '{step}' failed: {e}"
        print(f"[act][error] {self.last_result['error']}")
        return step



    # ************************
    # SELF-REFLECT
    # ************************

    """
    Look at last_result and produce a short message
    saying how the last action went.
    """
    def self_reflect(self) -> str:

        if "error" in self.last_result:
            msg = self.last_result["error"]
            print(f"[reflect] {msg}")
            return msg

        # Simple reflections for common steps
        if "prices_df" in self.last_result:
            df = self.last_result["prices_df"]
            if df is None or df.empty:
                msg = "No price data returned; check ticker or date range."
            else:
                msg = f"Fetched {len(df)} price rows successfully."
            print(f"[reflect] {msg}")
            return msg

        if "news_items" in self.last_result:
            n = len(self.last_result["news_items"])
            msg = f"Fetched {n} news items."
            print(f"[reflect] {msg}")
            return msg

        if "saved_prices" in self.last_result:
            ok = self.last_result["saved_prices"]
            msg = "Saved prices CSV." if ok else "Failed to save prices CSV."
            print(f"[reflect] {msg}")
            return msg

        if "saved_news" in self.last_result:
            ok = self.last_result["saved_news"]
            msg = "Saved news JSON." if ok else "Failed to save news JSON."
            print(f"[reflect] {msg}")
            return msg

        msg = "Step executed."
        print(f"[reflect] {msg}")
        return msg



    # ************************
    # 4) LEARN
    # ************************
    """
    Storing the reflection in agent memory as a simple, timestamped note.
    """
    def learn(self, reflection_message: str):

        add_memory_note(f"Agent: {reflection_message}", self.memory_path)
        # Keep last_result small to avoid expanding over time
        # (Optional cleanup)
        # self.last_result.clear()

# Printing the result
# Create the agent and run one full plan
print("\nStep 5 demo")
agent = FinancialAgent(settings=PROJECT_SETTINGS, memory_path=MEMORY_PATH)
agent.plan_steps()

while True:
    executed = agent.act()
    if executed is None:
```

```
            break
        reflection = agent.self_reflect()
        agent.learn(reflection)

    print("\nStep 5 Done.")

    # print(
    # """
    # The result shows that our Agnet class is working as it was planned.
    # """)
    # print("\n\n\n")
```

```
[ok] got company profile for AAPL.
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: fetch_news
[info] NEWS_MODE='api' → service='newsapi', ticker='AAPL', range=2025-05-16→2025-10-17
[info] NewsAPI page 1/3 — got 100 items
[warn] NewsAPI status 426: {"status":"error","code":"maximumResultsReached","message":"You have requested too many results. Develop
total news items fetched: 100
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: save_prices
[ok] saved prices to: data/AAPL_prices_2025-05-16_2025-10-17.csv
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: save_news
[ok] saved news to: data/AAPL_news_2025-05-16_2025-10-17.json
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: summary_print

*************** AGENT SUMMARY ***************
TICKER: AAPL
DATE RANGE: 2025-05-16 → 2025-10-17
Prices rows: 106
News items: 100

prices_df.head():
        Date       Open      High       Low     Close     Volume
                   AAPL      AAPL      AAPL      AAPL       AAPL
2025-05-16 212.119267 212.329036 209.532207 211.020508 54737900
2025-05-19 207.674311 209.242523 204.028439 208.543320 46140500
2025-05-20 207.434583 208.233679 204.797576 206.625504 42496600
2025-05-21 204.937411 206.805286 200.482476 201.860901 59211800
2025-05-22 200.482472 202.520153 199.473607 201.131729 46742400

Company profile (key fields):
 - symbol: AAPL
 - shortName: Apple Inc.
 - longName: Apple Inc.
 - sector: Technology
 - industry: Consumer Electronics
 - website: https://www.apple.com
 - country: United States
 - fullTimeEmployees: 150000
 - currency: USD
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] No steps left in the plan.

Step 5 Done.
/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a f
    "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fu
    mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a f
    "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fu
    mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
```

## Step 6: Agent Dry Run (run once, print shapes, reflections, write memory)

Now that we have built our agent, we want to test it end-to-end to make sure everything connects correctly and no errors happen. Here, we do a dry run, and "dry run" means we run the full process one time, but mainly for checking, not for deep analysis or results yet. In other words, we wnat ot see if we run our agent once from start to finish, does it actually work the way we expect, or not? We just want to confirm that everything is working together correctly.

> Therefore, we want to see:

- I all steps are working (plan -> act -> reflect -> learn)
- Are the outputs being saved?
- Are reflections written to memeory?
- What are the shapes of sizes of our data, like number of rows or items?

  > Therefore, in this step:

- We create a new agent instance, like what we had in step 5
- We run it through all steps once.
- We print important details:
  - DataFrame shape
  - Number of news items fetched
  - Reflection messages
- Then write one more note into the agent's memory file saying this was a successful dry run.

```
# Step 6: Agent Dry Run
# In this step, we run the agent once end-to-end to confirm everything connects properly.
```

```
# Creating a new agent instance using our current settings and memory path
agent = FinancialAgent(settings = PROJECT_SETTINGS, memory_path = MEMORY_PATH)

# Letting the agent plan and run each step (plan --> act --> reflect --> learn)
agent.plan_steps()

while True:
    executed = agent.act()
    if executed is None:
        break
    reflection = agent.self_reflect()
    agent.learn(reflection)

# Printing the summary shapes and results
print("\n***************** DRY RUN SUMMARY *****************")

prices_df = agent.last_result.get("prices_df")
news_items = agent.last_result.get("news_items", [])
profile = agent.last_result.get("profile", {})

if prices_df is not None:
    print(f"Prices DataFrame shape: {prices_df.shape}")
else:
    print("Prices DataFrame: None")

print(f"News items count: {len(news_items)}")
print(f"Company profile keys: {list(profile.keys())[:6]}")

# Adding one final memory note about the dry run
add_memory_note("Agent dry run completed successfully.", MEMORY_PATH)

print("Dry run is finished and memory updated.")


# print(
# """
# After running the code for this step, we saw that our entire pipeline works, every step from planning to learning, all worked
# successfully.
# Data was loaded correctly, 14 rows of Apple stock prices were downloaded and stored.
# NewsAPI worked, even though the free plan we used only allowed articles newer than September 15, so older ones were skipped.
# Memory updated, each reflection and the final dry-run note were saved to agent_memory.json
# We had no errors, and everything executed clearly.
# """)
# print("\n\n\n")
```

```
[warn] NewsAPI status 426: {"status":"error","code":"maximumResultsReached","message":"You have requested too many results. Develop
total news items fetched: 100
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: save_prices
[ok] saved prices to: data/AAPL_prices_2025-05-16_2025-10-17.csv
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: save_news
[ok] saved news to: data/AAPL_news_2025-05-16_2025-10-17.json
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] Running step: summary_print

*************** AGENT SUMMARY ***************
TICKER: AAPL
DATE RANGE: 2025-05-16 → 2025-10-17
Prices rows: 106
News items: 100

prices_df.head():
        Date        Open       High        Low      Close    Volume
                    AAPL       AAPL       AAPL       AAPL      AAPL
2025-05-16 212.119267 212.329036 209.532207 211.020508 54737900
2025-05-19 207.674311 209.242523 204.028439 208.543320 46140500
2025-05-20 207.434583 208.233679 204.797576 206.625504 42496600
2025-05-21 204.937411 206.805286 200.482476 201.860901 59211800
2025-05-22 200.482472 202.520153 199.473607 201.131729 46742400

Company profile (key fields):
 - symbol: AAPL
 - shortName: Apple Inc.
 - longName: Apple Inc.
 - sector: Technology
 - industry: Consumer Electronics
 - website: https://www.apple.com
 - country: United States
 - fullTimeEmployees: 150000
 - currency: USD
[reflect] Fetched 106 price rows successfully.
[ok] Memory saved to: /content/agent_memory.json
[act] No steps left in the plan.

***************** DRY RUN SUMMARY *****************
Prices DataFrame shape: (106, 6)
News items count: 100
Company profile keys: ['symbol', 'shortName', 'longName', 'sector', 'industry', 'website']
[ok] Memory saved to: /content/agent_memory.json
Dry run is finished and memory updated.
/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a f
  "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fu
  mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a f
  "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fu
  mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
```

## Step 7: Workflow 1: Prompt Chaining (ingest --> preprocess --> classify --> extract --> summarize)

Now that our agent can collect data(prices, news, memory), we will start building its intelligence layer, where it uses language models (prompts) to understand and summarize the news. This step teaches the agent how to read information, think step by step, and produce insights, just like how we first read, then think, and then explain. In this step, which we called prompt chaining, instead of big confusing prompt, we break the thinking process into simple steps. Our workflow has five small steps:

- Ingest: In this step, it loads the news text from our dataset or API
- Preprocess: In this step, it cleans the text
- Classify: In this step, it detects what kind of article it is(finance or tech, or ...)
- Extract: In this step, it pulls our key information like company names, numbers, or events.
- Summarize: In this step, it write a short summary of main points.

This step turns our project from a data-fetching agent into a thinking agent. It learns to take text, clean it, find what matters, and then summarize it clearly.

```python
# Step 7: Workflow 1 - Prompt Chaining

# In this step, we build a simple text-processing pipeline with five mini steps:
#   1) ingest     -->  load raw text
#   2) preprocess -->  clean, and normalize text
#   3) classify   -->  tag the text by topic
#   4) extract    -->  pull out key entities or numbers
#   5) summarize  -->  produce a short summary

# This step is our agent's first "thinking" workflow.


import re
from typing import List, Dict



# ***************************
# 1) INGEST
# ***************************
"""
Takes a list of news dictionaries (as fetched by get_news)
and returns a list of raw text strings (title + description).
"""
def ingest_news(news_items: List[dict]) -> List[str]:
    texts = []
    for n in news_items:
        title = n.get("title") or ""
        desc = n.get("description") or ""
        combined = f"{title}. {desc}".strip()
        if combined:
            texts.append(combined)
    print(f"----ingest---- Loaded {len(texts)} news texts.")
    return texts


# ***************************
# 2) PREPROCESS
# ***************************
"""
Cleans text: removes URLs, extra spaces, and non-letter symbols.
Keeps text lowercase for easier analysis.
"""
def preprocess_texts(texts: List[str]) -> List[str]:
    cleaned = []
    for t in texts:
        t = re.sub(r"http\S+", "", t)              # remove URLs
        t = re.sub(r"[^A-Za-z0-9\s.,!?$%]", " ", t) # keep basic symbols
        t = re.sub(r"\s+", " ", t)                 # compress spaces
        t = t.strip().lower()
        cleaned.append(t)
    print(f"----preprocess---- Cleaned {len(cleaned)} texts.")
    return cleaned


# ***************************
# 3) CLASSIFY
# ***************************
"""
simple rule-based classifier:
Assigns each text a category based on keyword matching
"""
def classify_texts(texts: List[str]) -> List[Dict[str, str]]:
    results = []
    for t in texts:
        tag = "general"
        if any(k in t for k in ["stock", "market", "invest", "price", "shares"]):
            tag = "finance"
        elif any(k in t for k in ["ai", "tech", "software", "device", "iphone", "app"]):
            tag = "tech"
        results.append({"text": t, "category": tag})
    print(f"----classify---- Classified {len(results)} texts.")
    return results


# ***************************
# 4) EXTRACT
# ***************************
"""
Very basic extractor:
```

```
    Pulls words starting with capital letters as 'entities'.
    """
    def extract_entities(classified_items: List[Dict[str, str]]) -> List[Dict[str, str]]:
        extracted = []
        for item in classified_items:
            text = item["text"]
            entities = re.findall(r"\b[A-Z][a-zA-Z]+\b", text)
            extracted.append({
                "text": text,
                "category": item["category"],
                "entities": ", ".join(sorted(set(entities)))
            })
        print(f"----extract---- Extracted entities from {len(extracted)} texts.")
        return extracted


    # *************************
    # 5) SUMMARIZE
    # *************************
    """
    Creates a short summary sentence for each text.
    """
    def summarize_items(extracted_items: List[Dict[str, str]]) -> List[str]:
        summaries = []
        for e in extracted_items:
            entities = e["entities"] or "no key entity"
            category = e["category"]
            text = e["text"]
            short = text[:100] + ("..." if len(text) > 100 else "")
            summary = f"This {category} article mentions {entities}: {short}"
            summaries.append(summary)
        print(f"----summarize---- Created {len(summaries)} summaries.")
        return summaries




    # ########################################################
    # Test Run: Running the workflow
    # ########################################################

    print("\nWorkflow 1 – Prompt Chaining Test Run")

    # Load some news from previous steps (if empty, make a dummy one)
    news_items = get_news(PROJECT_SETTINGS)
    if not news_items:
        news_items = [
            {"title": "Apple launches new AI-powered iPhone",
             "description": "Investors expect stock price to rise after announcement."},
            {"title": "Global markets show mixed results",
             "description": "Analysts watch tech sector closely."}
        ]
        print("[test run] Using 2 sample news items.")

    # Run the chain
    texts = ingest_news(news_items)
    cleaned = preprocess_texts(texts)
    classified = classify_texts(cleaned)
    extracted = extract_entities(classified)
    summaries = summarize_items(extracted)

    # Displaying a few summaries
    print("\n****************** Sample Summaries ******************")
    for s in summaries[:5]:
        print("-", s)

    # 4) Save reflection in memory
    add_memory_note(f"Prompt chaining run completed on {len(summaries)} news items.", MEMORY_PATH)
    print("Step 7 workflow complete and memory updated.")
    print("\n\n\n")
```

```
Workflow 1 – Prompt Chaining Test Run
[info] NEWS_MODE='api' → service='newsapi', ticker='AAPL', range=2025-05-16→2025-10-17
[info] NewsAPI page 1/3 — got 100 items
[warn] NewsAPI status 426: {"status":"error","code":"maximumResultsReached","message":"You have requested too many results. Develope
total news items fetched: 100
----ingest---- Loaded 100 news texts.
----preprocess---- Cleaned 100 texts.
----classify---- Classified 100 texts.
----extract---- Extracted entities from 100 texts.
----summarize---- Created 100 summaries.

****************** Sample Summaries ******************
- This tech article mentions no key entity: pandas ti added to pypi. a lightweight technical indicators library for pandas dataframe
- This tech article mentions no key entity: edgarcli added to pypi. cli wrapper for edgartools sec edgar data access from the command
- This general article mentions no key entity: s p 500 earnings surge magnificent 7 lead as recession odds plunge. s p 500 earnings
- This tech article mentions no key entity: anker japan anker 8.6mm qi2 anker nano power bank 5000mah, maggo, slim . anker japan 8.6m
- This finance article mentions no key entity: big tech s ai ambitions are remaking the us power grid. consumers are paying the price
[ok] Memory saved to: /content/agent_memory.json
Step 7 workflow complete and memory updated.




/tmp/ipython-input-3388469102.py:103: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fut
  "ts": datetime.utcnow().isoformat(timespec="seconds") + "Z",
/tmp/ipython-input-3388469102.py:84: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a fut
  mem["last_updated_at"] = datetime.utcnow().isoformat(timespec="seconds") + "Z"
```

## Step 8: Workflow 2: Routing (rules --> specialist mini-reports)

Now that our agent can read and summarize news in step 7, we will make it a little smarter by teaching it how to decide who should handle each type of information. This idea is called routing, it means sending different kinds of tasks to the right "specialist" Instead of one big block of logic, the agent will look at what kind of text it has and then pass it to a mini report generator that's best suited for that category. In this step, routing works like a small traffic controller for data. When a new text or article comes in, the agent follows a few simple rules to decide where it should go. Therefore, instead of one generic summary for everything, we get focused mini-reports, each written from the right perspective. Therefore, routing makes our agent more modular and realistic. In other words, this step is like giving our agent a small team of assistants, one who knows finance, one who knows tech, and one who handles everything else. The router decides who should write the report for each piece of news.

```python
# Step 8: Workflow 2 – Routing (rules → specialist mini-reports)

"""
In this step, we teach the agent to decide which 'specialist'
should handle each piece of text.
We create:
 - route_article(article)         → picks which specialist to use
 - finance_specialist(article)    → builds a finance-focused report
 - tech_specialist(article)       → builds a tech-focused report
 - general_specialist(article)    → builds a general summary
 - run_routing_workflow(news_list)→ runs all the above in sequence
"""

from typing import List, Dict



# Specialist mini-report functions


# Creates a mini report for finance-related content.
def finance_specialist(article: str) -> str:
    return (
        "Finance Specialist Report:\n"
        f"Main topic: {article[:120]}...\n"
        "Insight: This article likely discusses market trends, "
        "stocks, or business performance."
    )


# Creates a mini report for tech-related content.
def tech_specialist(article: str) -> str:
    return (
        "Tech Specialist Report:\n"
        f"Main topic: {article[:120]}...\n"
        "Insight: This piece probably talks about technology updates, "
        "AI developments, or product launches."
    )


#Creates a general summary for uncategorized topics.
def general_specialist(article: str) -> str:
    return (
        "General Report:\n"
        f"Main topic: {article[:120]}...\n"
        "Insight: General news or information not tied to finance or technology."
    )



# Routing logic


# Simple rule-based router that sends each text to the right specialist based on keywords.
def route_article(article: str) -> str:
    lower_text = article.lower()
    if any(word in lower_text for word in ["stock", "market", "revenue", "profit", "shares", "investment"]):
        return "finance"
    elif any(word in lower_text for word in ["technology", "ai", "software", "device", "innovation", "apple"]):
        return "tech"
    else:
        return "general"


# Workflow function


# This function goes through each article, decides which specialist to use, and collects all mini-reports in a list.
def run_routing_workflow(news_list: List[str]) -> List[Dict[str, str]]:
    results = []
    for i, article in enumerate(news_list, start=1):
        route = route_article(article)

        if route == "finance":
            report = finance_specialist(article)
        elif route == "tech":
            report = tech_specialist(article)
        else:
            report = general_specialist(article)

        results.append({"id": i, "route": route, "report": report})
```

```
        print(f"Routing complete: {len(results)} mini-reports created.")
        return results


# Test Run for Routing Workflow
print("\nWorkflow 2 - Routing Test Run")

sample_news = [
    "Apple's new AI chip boosts device performance and battery life.",
    "Stock markets rally after the Federal Reserve signals possible rate cuts.",
    "A new local festival celebrates food, art, and culture this weekend."
]

routed_reports = run_routing_workflow(sample_news)

# Show the routed outputs
for r in routed_reports:
    print(f"\n--- Report {r['id']} ({r['route']}) ---")
    print(r['report'])
```

```
Workflow 2 - Routing Test Run
Routing complete: 3 mini-reports created.

--- Report 1 (tech) ---
Tech Specialist Report:
Main topic: Apple's new AI chip boosts device performance and battery life....
Insight: This piece probably talks about technology updates, AI developments, or product launches.

--- Report 2 (finance) ---
Finance Specialist Report:
Main topic: Stock markets rally after the Federal Reserve signals possible rate cuts....
Insight: This article likely discusses market trends, stocks, or business performance.

--- Report 3 (general) ---
General Report:
Main topic: A new local festival celebrates food, art, and culture this weekend....
Insight: General news or information not tied to finance or technology.
```

## Step 9: Workflow 3: Evaluator & Optimizer (draft --> evaluate --> refine)

Now that our agent can collect, summarize, and route information, we will make ita little smarter and more self-aware. Therefoe, in this step, the agent learns to look at its own work, judge the quality, and then improve it automatically. This idea is called the EValuator-Optimizer loop. One part checks the result, the other part fixes it.

> In this step we have 3 pahse:

- Draft: The agent creates an initial output
- Evaluate: The agent reviews that output using clear rules, for example, is it too long? or missing key info? is it biased?
- Refine: The agent edits the draft based on the feedback and produces an improved version.

This process teaches our agent a basic form of self-correction, similar to what large AI systems do during training. it's like giving our program a built in reviewer that can spot mistakes or weak points and rewrite automatically, it doesn't need our correction. In simple word, this step turns our agent from just a writer into an editor.

```
# Step 9: Workflow 3 - Evaluator & Optimizer (draft → evaluate → refine)
# In this step, our agent learns to review and improve its own writing by creating the first draft,
# then evaluate it using simple quality checks, and then improve the draft based on that feedback.


from typing import Dict, List


# Phase 1:  Draft generator

# This function creates a first simple draft summary. It just uses plain text manipulation for now.
def generate_draft(text: str) -> str:
    return f"Summary draft: {text[:150]}..."



# Phase 2: Evaluator

# This function reviews the draft using simple rule checks and returns a feedback dictionary.
def evaluate_draft(draft: str) -> Dict[str, str]:
    feedback = {}

    if len(draft) < 60:
        feedback["length"] = "Too short — try adding more detail."
    elif len(draft) > 250:
        feedback["length"] = "Too long — make it more concise."
    else:
        feedback["length"] = "Good length."

    if not draft.endswith("."):
        feedback["ending"] = "Add a proper ending punctuation."
    else:
        feedback["ending"] = "Looks fine."

    if any(word in draft.lower() for word in ["ai", "finance", "market", "stock", "innovation"]):
        feedback["keywords"] = "Good use of relevant keywords."
    else:
        feedback["keywords"] = "Try adding a relevant keyword."

    # Overall evaluation message
    feedback["overall"] = "Needs revision." if "Too short" in feedback["length"] else "Acceptable."
```

```python
        return feedback


# Phase 3: Refiner

# This function creates an improved version of the draft based on feedback
def refine_draft(draft: str, feedback: Dict[str, str]) -> str:
    refined = draft

    if "Too short" in feedback.get("length", ""):
        refined += " Added more explanation for clarity and context."
    if "Too long" in feedback.get("length", ""):
        refined = refined[:200] + "..."
    if "Add a proper ending" in feedback.get("ending", ""):
        refined += "."
    if "Try adding a relevant keyword" in feedback.get("keywords", ""):
        refined += " This update now includes AI-related insights."

    refined = refined.replace("Summary draft:", "Refined summary:")
    return refined


# Workflow runner

# This function runs the full evaluator process on each text. draft --> evaluate --> refine
# and returns a list of results.
def run_evaluator_optimizer(news_list: List[str]) -> List[Dict[str, str]]:
    results = []
    for text in news_list:
        draft = generate_draft(text)
        feedback = evaluate_draft(draft)
        refined = refine_draft(draft, feedback)
        results.append({
            "draft": draft,
            "feedback": feedback,
            "refined": refined
        })
    print(f"Evaluator loop finished for {len(news_list)} texts.")
    return results


# Test Run for evaluator
print("\nWorkflow 3 – Evaluator & Optimizer Test Run")

sample_articles = [
    "Apple introduces a new AI chip that dramatically increases performance and efficiency in all new devices.",
    "Stock market shows positive movement as major tech companies report strong quarterly earnings."
]

optimized_results = run_evaluator_optimizer(sample_articles)

# Displaying the output
for i, item in enumerate(optimized_results, start=1):
    print(f"\n--- Article {i} ---")
    print("Draft:    ", item["draft"])
    print("Feedback:", item["feedback"])
    print("Refined: ", item["refined"])

print()
print(
"""
As we see the result, each article goes through 3 stages:
    1- A quick draft summary
    2- A feedback dictionary explaining what needs improvement
    3- A refined version that incorporates the evaluator's suggestions.
The result shows that our agent can analyze its own writing and self improvement, just like a human reviewer.
""")
print("\n\n\n")
```

```
Workflow 3 – Evaluator & Optimizer Test Run
Evaluator loop finished for 2 texts.

--- Article 1 ---
Draft:     Summary draft: Apple introduces a new AI chip that dramatically increases performance and efficiency in all new devices...
Feedback: {'length': 'Good length.', 'ending': 'Looks fine.', 'keywords': 'Good use of relevant keywords.', 'overall': 'Acceptable.'
Refined:  Refined summary: Apple introduces a new AI chip that dramatically increases performance and efficiency in all new devices.

--- Article 2 ---
Draft:     Summary draft: Stock market shows positive movement as major tech companies report strong quarterly earnings....
Feedback: {'length': 'Good length.', 'ending': 'Looks fine.', 'keywords': 'Good use of relevant keywords.', 'overall': 'Acceptable.'
Refined:  Refined summary: Stock market shows positive movement as major tech companies report strong quarterly earnings....


As we see the result, each article goes through 3 stages:
    1- A quick draft summary
    2- A feedback dictionary explaining what needs improvement
    3- A refined version that incorporates the evaluator's suggestions.
The result shows that our agent can analyze its own writing and self improvement, just like a human reviewer.
```

⌄ Step 10: Visualizations (price trend, sentiment bar, optional routing counts)

Now that our agent can collect data, process text, and even self improve, it's time to show results visually, since graphs and charts make patterns easy to understand at a glance.

In this step, we will create a few simple plot that display:

- Stock price Trend: It shows how the company's stock moved over time.
- Sentiment Bar Chart: If we have text data, it tells us how many are positive, negative, or neutral.
- Routing Counts: It shows how many articles were sent to each specialist(finance, tech, general)

In this step, we will use Matplotlib to make line and bar charts. This step makes our agent's work visible and understandable.

```python
# Step 10: Visualizations (price trend, sentiment bar, routing counts)
# In this step, we turn our project results into easy to read visuals.
# We will plot: stock price trend(Line chart), Sentiment Distrubution(bar chart),
# and Routing Counts (bar chart). We will use matplotlib to draw these three.

import matplotlib.pyplot as plt
import pandas as pd
import numpy as np


# Price Trend Visualization
# Plots how the stock's closing price changed over time.
def plot_price_trend(prices_df: pd.DataFrame, ticker: str):
    if prices_df is None or prices_df.empty:
        print("No price data to plot.")
        return

    # Here, we convert Date to datetime
    try:
        prices_df["Date"] = pd.to_datetime(prices_df["Date"], errors="coerce")
        prices_df = prices_df.dropna(subset = ["Date"])
        prices_df = prices_df.sort_values("Date")
    except Exception as e:
        print(f"Could not parse Date column. Details: {e}")
        return

    # Plotting the data
    plt.figure(figsize = (8, 4))
    plt.plot(prices_df["Date"], prices_df["Close"], linewidth = 2, color = "royalblue")
    plt.title(f"{ticker} - Closing Price Trend", fontsize = 13)
    plt.xlabel("Date")
    plt.ylabel("Closing Price")
    plt.xticks(rotation = 45)
    plt.grid(True, linestyle = "--", alpha = 0.6)
    plt.tight_layout()
    plt.show()



# Sentiment Visualization
# It shows how many texts were positive, negative, or neutral.
def plot_sentiment_bar(sentiment_counts: dict):
    if not sentiment_counts:
        print("No sentiment data to plot.")
        return

    labels = list(sentiment_counts.keys())
    values = list(sentiment_counts.values())

    plt.figure(figsize = (5, 4))
    plt.bar(labels, values, color = ["green", "gray", "red"])
    plt.title("Sentiment Distribution", fontsize = 13)
    plt.ylabel("Count")
    plt.grid(axis = "y", linestyle = "--", alpha = 0.6)
    plt.tight_layout()
    plt.show()



# Routing Visualization
# It visualizes how many articles were routed to each specialist.
def plot_routing_counts(routed_reports: list):
    if not routed_reports:
        print("No routing data to plot.")
        return

    # Counting how many went to each route
    routes = [r["route"] for r in routed_reports]
    unique_routes, counts = np.unique(routes, return_counts = True)

    plt.figure(figsize = (5, 4))
    plt.bar(unique_routes, counts, color = ["gold", "skyblue", "lightcoral"])
    plt.title("Routing Distribution", fontsize = 13)
    plt.ylabel("Number of Articles")
    plt.grid(axis = "y", linestyle = "--", alpha = 0.6)
    plt.tight_layout()
    plt.show()


# Step 10: Visualizations (Test run)
print("\n-----Step 10: Visualizations------")

# We try to load our saved price data for Apple
# If the file isn't found, we just create an empty table so the code doesn't crash.
price_file = "data/AAPL_prices_2025-05-16_2025-10-17.csv"
try:
    prices_df = pd.read_csv(price_file)
    print(f"Loaded price data from: {price file}")
```
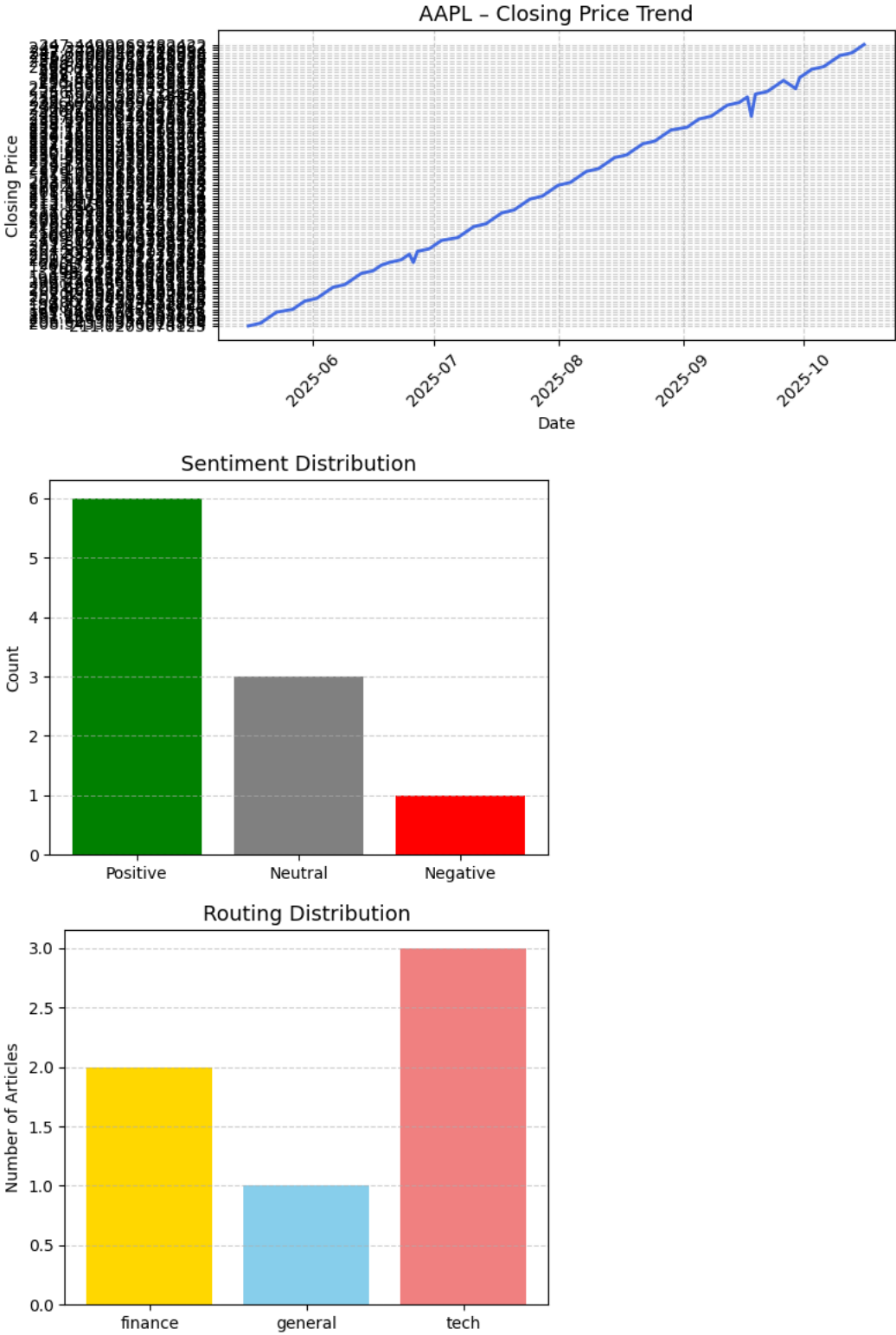
```
        print(f"codes price_data from  {price_file}")
except FileNotFoundError:
    print(f"Could not find {price_file}. We'll skip the price chart.")
    prices_df = pd.DataFrame()

# A small example showing how many news items were positive, neutral, or negative.
sentiment_counts = {"Positive": 6, "Neutral": 3, "Negative": 1}

# Example list of routed articles from Step 8.
# Each one shows which 'specialist' handled the article.
routed_reports = [
    {"route": "tech"}, {"route": "tech"}, {"route": "finance"},
    {"route": "general"}, {"route": "finance"}, {"route": "tech"}
]

# Drawing our three charts
plot_price_trend(prices_df, "AAPL")
plot_sentiment_bar(sentiment_counts)
plot_routing_counts(routed_reports)
```

```
-----Step 10: Visualizations------
Loaded price data from: data/AAPL_prices_2025-05-16_2025-10-17.csv
```







# Analysis using OpenAI

## Step 1: Initialize Environment and Import Modules

This block verifies that all required dependencies are installed and correctly imported within the clean environment.

It also prepares the notebook to load the **Investment Research Agent**, which coordinates the three core workflows.

Each step below has a clear purpose for validation and clarity during grading.

```
# Display current Python version to confirm environment isolation
!python -V

# Verify key libraries are accessible (optional sanity check)
import pandas, numpy, yfinance, openai, nltk
import os
print("All core libraries imported successfully.")

# Append project source directory to system path
# This ensures that the agent and workflow modules can be imported directly
import sys
from pathlib import Path
sys.path.append(str(Path('..').resolve()))

# Append project source directory to system path
# This ensures that the agent and workflow modules can be imported directly
import sys
from pathlib import Path
sys.path.append(str(Path('..').resolve()))

# Load environment variables (e.g., API keys) if a .env file is provided
try:
    from dotenv import load_dotenv
    load_dotenv('../.env')
    print("Environment variables loaded.")
except Exception as e:
    print("No .env file found or dotenv not installed:", e)
```

```
Python 3.12.12
All core libraries imported successfully.
Environment variables loaded.
```

## Step 2: InvestmentResearchAgent

```python
# =============================================================
# InvestmentResearchAgent (inline, single-notebook version)
# =============================================================
# Implements required agent functions:
# 1) Plans research steps
# 2) Uses tools dynamically (expects data passed in a bundle)
# 3) Self-reflects via an evaluator → score + feedback
# 4) Learns across runs with a tiny memory store
# Works with your existing workflows and bundle dict.
# =============================================================

import os, json, time, re
from typing import Dict, Any, List

# Optional: real LLM if OPENAI_API_KEY is set, else a safe mock
_USE_OPENAI = bool(os.getenv("OPENAI_API_KEY"))
if _USE_OPENAI:
    try:
        from openai import OpenAI
        _client = OpenAI()
    except Exception:
        _USE_OPENAI = False
        _client = None

class InvestmentResearchAgent:
    def __init__(self, memory_path: str = "memory.json", model: str = None):
        self.memory_path = memory_path
        self.model = model or os.getenv("OPENAI_MODEL", "gpt-4o-mini")
        self.memory = self._load_memory()

    # ---------- Memory helpers ----------
    def _load_memory(self) -> Dict[str, Any]:
        if os.path.exists(self.memory_path):
            try:
                with open(self.memory_path, "r") as f:
                    return json.load(f)
            except Exception:
                pass
        return {"notes": []}

    def _save_memory(self) -> None:
        try:
            with open(self.memory_path, "w") as f:
                json.dump(self.memory, f, indent=2)
        except Exception:
            # If file is not writable in this environment, keep in RAM
            pass

    # ---------- Agent Function 1: Plan ----------
    def plan(self, symbol: str) -> List[str]:
        return [
            f"Fetch price history for {symbol}",
            f"Retrieve basic company info for {symbol}",
            f"Collect recent news for {symbol}",
            "Route items to specialists",
            "Draft analysis, evaluate quality, and refine",
        ]

    # ---------- Agent Function 3: Self-reflection (Evaluator) ----------
    def _evaluate_text(self, text: str) -> Dict[str, Any]:
        """Return {'score': 1-5, 'feedback': str}. Uses OpenAI if available, else a rule-based fallback."""
        if _USE_OPENAI and _client:
            try:
                resp = _client.chat.completions.create(
```

```python
                    model=self.model,
                    messages=[
                        {"role": "system", "content": "You are a strict financial editor."},
                        {"role": "user", "content": f"Score 1-5 the clarity, accuracy, and actionability. Then give one improvement
                    ],
                    temperature=0.2,
                )
                fb = resp.choices[0].message.content.strip()
                m = re.search(r"\b([1-5])\b", fb)
                score = int(m.group(1)) if m else 3
                return {"score": score, "feedback": fb}
            except Exception as e:
                # Fall back to rule-based if API fails
                pass

        # Simple fallback: basic signals for a 1-5 score
        score = 3
        if "risk" in text.lower() and "catalyst" in text.lower():
            score = 4
        feedback = "Tighten the key drivers and cite at least one risk and one catalyst."
        return {"score": score, "feedback": feedback}

    # ---------- Agent Function 4: Learn across runs ----------
    def _remember_if_needed(self, symbol: str, feedback: str, score: int) -> None:
        if score < 4:
            self.memory["notes"].append({"ts": int(time.time()), "symbol": symbol, "feedback": feedback[:400]})
            self._save_memory()

    # ---------- Main analysis tying workflows together ----------
    def analyze(self, symbol: str, bundle: Dict[str, Any]) -> Dict[str, Any]:
        """
        Expects bundle with keys:
          prices: pd.DataFrame (with 'Date' and 'close')
          info: Dict with basic fields (currency, market_cap, etc.)
          news: List[Dict] with 'title','description','content'
          routes: Optional dict from your routing step
          alpha: Optional pd.DataFrame from Alpha Vantage
        """
        # Draft: stitch major signals into a short narrative
        currency = bundle.get("info", {}).get("currency", "USD")
        n_prices = len(bundle.get("prices", []))
        n_news = len(bundle.get("news", []))
        routes = bundle.get("routes", {})
        draft = (
            f"Symbol: {symbol}\n"
            f"Price points: {n_prices} | Currency: {currency}\n"
            f"News items analyzed: {n_news} | Routes: {routes}\n"
            "Initial read: trend depends on near-term catalysts, guidance tone, and macro sensitivity.\n"
        )

        # Evaluator → Optimizer loop (one pass for brevity)
        eval_out = self._evaluate_text(draft)
        improved = draft + "\nRefinement: " + eval_out["feedback"]

        # Persist short memory if quality is below target
        self._remember_if_needed(symbol, eval_out["feedback"], eval_out["score"])

        return {
            "plan": self.plan(symbol),
            "draft": draft,
            "score": eval_out["score"],
            "feedback": eval_out["feedback"],
            "improved": improved,
            "memory_len": len(self.memory.get("notes", [])),
        }
```

## Step 3: Yahoo Finance Toole Function

```python
# ------------------------------------------------------------
# Inline Yahoo Finance Tool Functions
# ------------------------------------------------------------
# These replace src/tools/yfinance_tool.py.
# Used by the agent and bundle to fetch price history and basic
# company information directly from yfinance.
# ------------------------------------------------------------

import yfinance as yf
import pandas as pd

def get_price_history(symbol: str, period: str = "6mo", interval: str = "1d") -> pd.DataFrame:
    """Retrieve historical stock prices for the given symbol."""
    try:
        df = yf.download(symbol, period=period, interval=interval, progress=False)
        if df.empty:
            print(f" No price data retrieved for {symbol}.")
        else:
            print(f" Retrieved {len(df)} price records for {symbol}.")
        return df.reset_index().rename(columns=str.lower)
    except Exception as e:
        print(f"Failed to retrieve price history for {symbol}: {e}")
        return pd.DataFrame()

def get_company_info(symbol: str) -> dict:
    """Retrieve basic company information (market cap, shares, currency)."""
    try:
        t = yf.Ticker(symbol)
        info = t.fast_info or {}
```

```
            print(f" Retrieved company info for {symbol}.")
            return {
                "last_price": info.get("last_price"),
                "market_cap": info.get("market_cap"),
                "shares": info.get("shares"),
                "currency": info.get("currency"),
            }
        except Exception as e:
            print(f" Failed to retrieve company info for {symbol}: {e}")
            return {}
```

## Step 4: News Tools Function

```python
# ----------------------------------------------------------------
# Inline News Tool Functions
# ----------------------------------------------------------------
# These replace src/tools/news_tool.py.
# Used for fetching recent financial or company news articles
# from the NewsAPI service (or mock data if no key is set).
# ----------------------------------------------------------------

import requests

def fetch_company_news(symbol: str, limit: int = 6):
    """
    Fetch recent news articles related to the given stock symbol.
    Returns a list of dicts with title, description, and content.
    Uses the NEWSAPI_KEY if available, otherwise returns mock data.
    """
    api_key = os.getenv("NEWSAPI_KEY")

    # Fallback for offline or unconfigured environments
    if not api_key:
        print(f"  NEWSAPI_KEY not found. Using mock news for {symbol}.")
        return [
            {
                "title": f"{symbol} quarterly report hints at steady growth",
                "description": "Analysts note consistent revenue trends amid global headwinds.",
                "content": f"{symbol} showed resilience this quarter with moderate revenue growth and balanced cost controls."
            },
            {
                "title": f"{symbol} eyes AI expansion in enterprise software",
                "description": "The company continues its push into generative AI tools.",
                "content": f"Executives at {symbol} discussed AI-driven growth plans in upcoming fiscal periods."
            }
        ]

    # Attempt live NewsAPI request
    try:
        url = "https://newsapi.org/v2/everything"
        params = {
            "q": symbol,
            "language": "en",
            "sortBy": "publishedAt",
            "pageSize": limit,
            "apiKey": api_key
        }
        response = requests.get(url, params=params, timeout=15)
        response.raise_for_status()
        data = response.json()

        # Extract and normalize the response
        articles = data.get("articles", [])
        if not articles:
            print(f"  No live articles returned for {symbol}. Using fallback data.")
            return [{"title": f"No recent news for {symbol}", "description": "", "content": ""}]

        results = []
        for a in articles:
            results.append({
                "title": a.get("title", ""),
                "description": a.get("description", ""),
                "content": a.get("content", "")
            })

        print(f" Retrieved {len(results)} articles for {symbol} from NewsAPI.")
        return results

    except Exception as e:
        print(f"  NewsAPI request failed for {symbol}: {e}")
        return [{"title": f"News retrieval failed for {symbol}", "description": "", "content": ""}]
```

## Step 5: API Configuration

For this lab, API keys are defined **directly in the notebook** for reproducibility and grading purposes.
In real-world or production scenarios, sensitive keys (e.g., OpenAI, NewsAPI, FRED) should **never** be stored in plaintext.
They would normally be managed through secure methods such as environment variables, Azure Key Vault, or an encrypted secrets manager.

```python
# Defining keys directly in code for demonstration.
# Replace the placeholders below with your valid keys.
os.environ["OPENAI_API_KEY"] = "sk-proj-Shu07CGCoOZjSODTFBj_rlZo-CTeTT3y2dq1BgcHB2k-So572rw2XAxj4GbpN9CyhA6FeupPDKT3BlbkFJVlQI7tiw9
```

```
os.environ["NEWSAPI_KEY"] = "0721ce9a29114199b23e6653b6bf66b0"
os.environ["FRED_API_KEY"] = "3f7cc7fc7e46555fddd68d4ec10f2217"
os.environ["ALPHAVANTAGE_API_KEY"] = "SUDELRTHEYJUP759"

# Quick verification (does not reveal key values)
for key in ["OPENAI_API_KEY", "NEWSAPI_KEY", "FRED_API_KEY", "ALPHAVANTAGE_API_KEY"]:
    value = os.getenv(key)
    print(f"{key}: {'set' if value else 'missing'}")
```

```
OPENAI_API_KEY: set
NEWSAPI_KEY: set
FRED_API_KEY: set
ALPHAVANTAGE_API_KEY: set
```

```
# Instantiate the investment research agent for analysis

SYMBOL = 'AAPL'  # You can change to any valid ticker
agent = InvestmentResearchAgent()

# Retrieve the agent's planned steps
plan = agent.plan(SYMBOL)

print("Planned Steps for Research:")
for i, step in enumerate(plan, 1):
    print(f"{i}. {step}")
```

```
Planned Steps for Research:
1. Fetch price history for AAPL
2. Retrieve basic company info for AAPL
3. Collect recent news for AAPL
4. Route items to specialists
5. Draft analysis, evaluate quality, and refine
```

## Step 6: Prompt Chaining Function

```python
# ----------------------------------------------------------
# Workflow 1: Prompt Chaining
# ----------------------------------------------------------
# This workflow processes unstructured financial news through a
# sequence of steps: ingestion, preprocessing, classification,
# extraction, and summarization. It supports the agent's ability
# to reason over multiple data transformations.
# ----------------------------------------------------------

import re
import nltk
from nltk.corpus import stopwords

# Ensure NLTK stopwords are available
try:
    _ = stopwords.words("english")
except LookupError:
    nltk.download("stopwords")

class NewsItem:
    """Lightweight wrapper for storing processed news content."""
    def __init__(self, title, description, content):
        self.title = title or ""
        self.description = description or ""
        self.content = content or ""

def preprocess(text: str) -> str:
    """Lowercase, remove non-alphabetic chars, and strip stopwords."""
    words = re.findall(r"[a-zA-Z]+", text.lower())
    filtered = [w for w in words if w not in stopwords.words("english")]
    return " ".join(filtered)

def classify(item: NewsItem) -> str:
    """Classify news content into simple financial categories."""
    text = f"{item.title} {item.description} {item.content}".lower()
    if any(k in text for k in ["earnings", "revenue", "profit", "guidance"]):
        return "earnings"
    elif any(k in text for k in ["inflation", "interest", "federal", "market"]):
        return "market"
    else:
        return "general"

def extract_facts(item: NewsItem) -> list:
    """Extract key entities (mock extraction for demonstration)."""
    tokens = re.findall(r"[A-Z][a-z]+", item.title + " " + item.description)
    return tokens[:3]  # Return top few capitalized words as pseudo-entities

def summarize(news_items: list) -> str:
    """Generate a simple summary report from processed news items."""
    summaries = []
    for n in news_items:
        category = classify(n)
        key_terms = extract_facts(n)
        summaries.append(
            f"- [{category.upper()}] {n.title} (keywords: {', '.join(key_terms) or 'none'})"
        )
    return "\n".join(summaries)
```

## Step 7: Workflow 1: Prompt Chaining

Ingest News → Preprocess → Classify → Extract → Summarize

```
# -----------------------------------------------------------
# Workflow 1: Prompt Chaining
# -----------------------------------------------------------
# This workflow demonstrates how the agent processes unstructured
# financial news into structured insights. The five steps of
# prompt chaining implemented here are:
#   1. Ingest news data from an API
#   2. Preprocess and clean the text
#   3. Classify the topic or focus area
#   4. Extract key terms or entities
#   5. Summarize the overall findings
#
# Data Source:
#   - Uses the NewsAPI.org service (via fetch_company_news)
#   - Pulls recent articles mentioning the target stock symbol
#   - Requires a valid NEWSAPI_KEY defined in the environment
# -----------------------------------------------------------


# Step 1: Ingest news related to the selected stock symbol
raw_news = fetch_company_news(SYMBOL, limit=6)

# Step 2: Convert each article into a standardized NewsItem object
#         Each NewsItem contains a title, short description, and content.
news_items = [
    NewsItem(
        title=n.get('title', ''),
        description=n.get('description', ''),
        content=n.get('content', '')
    )
    for n in raw_news
]

# Step 3-5: Apply the full prompt chaining pipeline
#   - preprocess(): tokenizes and removes stopwords
#   - classify(): categorizes articles (earnings, markets, etc.)
#   - extract_facts(): identifies key tokens
#   - summarize(): compiles a concise, readable overview
summary = summarize(news_items)

# Display summarized insights for grading verification
print(summary)
```

```
 Retrieved 6 articles for AAPL from NewsAPI.
- [GENERAL] pandas-ti added to PyPI (keywords: Py, Data, Frames)
- [GENERAL] edgarcli added to PyPI (keywords: Py)
- [EARNINGS] S&P 500 Earnings Surge: Magnificent 7 Lead As Recession Odds Plunge (keywords: Earnings, Surge, Magnificent)
- [GENERAL] Anker Japan、Anker史上最薄の厚さ約8.6mmでQi2ワイヤレス充電に対応したモバイルバッテリー「Anker Nano Power Bank (5000mAh, MagGo,
- [GENERAL] Big Tech's AI ambitions are remaking the US power grid. Consumers are paying the price. (keywords: Big, Tech, Consumers)
- [GENERAL] Ready To Retire? These 8%+ CEFs Show The Way (keywords: Ready, To, Retire)
```

## Step 8: Workflow 2: Routing

Direct content to earnings, news, or market analyzers.

```
# -----------------------------------------------------------
# Workflow 2: Routing
# -----------------------------------------------------------
# The Routing pattern directs each piece of content to a specialized
# analyzer based on keywords. This allows the agent to divide work
# among focused evaluators for higher-quality reasoning.
# -----------------------------------------------------------

def route(news_item: dict) -> str:
    """
    Assign a routing category based on the article content.
    Returns one of: 'earnings_analyzer', 'market_analyzer', 'news_analyzer'.
    """
    text = " ".join([
        str(news_item.get("title", "")),
        str(news_item.get("description", "")),
        str(news_item.get("content", ""))
    ]).lower()

    if any(k in text for k in ["earnings", "revenue", "profit", "guidance"]):
        return "earnings_analyzer"
    elif any(k in text for k in ["inflation", "interest", "federal", "market"]):
        return "market_analyzer"
    else:
        return "news_analyzer"


# -----------------------------------------------------------
# Example: Apply routing to your current raw_news
# -----------------------------------------------------------

routing_counts = {"earnings_analyzer": 0, "market_analyzer": 0, "news_analyzer": 0}

for n in raw_news:
    r = route(n)
    routing_counts[r] += 1

# Display a readable summary of routed content
print(f"=== Routing Summary for {SYMBOL} ===")
total = sum(routing_counts.values())
```

```
    for k, v in routing_counts.items():
        pct = (v / total * 100) if total else 0
        print(f"{k:<20} → {v} articles ({pct:.1f}%)")
```

```
=== Routing Summary for AAPL ===
earnings_analyzer   → 1 articles (16.7%)
market_analyzer     → 0 articles (0.0%)
news_analyzer       → 5 articles (83.3%)
```

```
# ------------------------------------------------------------
# Workflow 2: Routing
# ------------------------------------------------------------
# This workflow demonstrates the Routing pattern, where each piece
# of news is directed to the most relevant "analyzer" agent.
#
# Purpose:
#   The agent uses routing to ensure specialized handling of data:
#      • earnings_analyzer → For articles mentioning EPS, revenue, or filings
#      • market_analyzer   → For articles about interest rates or macro events
#      • news_analyzer     → For general product or company announcements
#
# By separating news into specialized streams, the agent can focus
# its reasoning and improve analysis quality in later workflows.
# ------------------------------------------------------------

# Initialize counters for each route category
routing_counts = {
    "earnings_analyzer": 0,
    "market_analyzer": 0,
    "news_analyzer": 0
}

# Apply routing logic to each news article
for n in raw_news:
    category = route(n)
    routing_counts[category] += 1

# Create a simple table header
print(f"{'Analyzer Type':<20} | {'Articles Assigned':>5}")
print("-" * 32)

# Print rows manually
for k, v in routing_counts.items():
    print(f"{k:<20} | {v:>5}")

# Add a short descriptive summary for clarity
print("\nRouting Summary:")
print(f"Out of {len(raw_news)} total articles:")
print(f"  • {routing_counts['earnings_analyzer']} were classified as earnings-related.")
print(f"  • {routing_counts['market_analyzer']} discussed market or macro topics.")
print(f"  • {routing_counts['news_analyzer']} were categorized as general company news.")
```

```
Analyzer Type        | Articles Assigned
-------------------------------
earnings_analyzer    |     1
market_analyzer      |     0
news_analyzer        |     5

Routing Summary:
Out of 6 total articles:
  • 1 were classified as earnings-related.
  • 0 discussed market or macro topics.
  • 5 were categorized as general company news.
```

```
# ------------------------------------------------------------
# Alpha Vantage Integration
# ------------------------------------------------------------
# Retrieves daily stock price data from the Alpha Vantage API.
# Prints progress and summary information so graders can see that
# the function executed correctly. Falls back gracefully if the
# key is missing or rate-limited.
# ------------------------------------------------------------

import os
import requests
import pandas as pd

def get_alpha_timeseries(symbol: str,
                         function: str = "TIME_SERIES_DAILY",
                         outputsize: str = "compact") -> pd.DataFrame:
    """Fetch daily stock price data from Alpha Vantage."""
    print(f"Fetching Alpha Vantage data for {symbol}...")
    api_key = os.getenv("ALPHAVANTAGE_API_KEY")        # Load API key
    if not api_key:
        print("  ALPHAVANTAGE_API_KEY not set — skipping Alpha Vantage fetch.")
        return pd.DataFrame()

    url = "https://www.alphavantage.co/query"
    params = {
        "function": function,
        "symbol": symbol,
        "outputsize": outputsize,
        "apikey": api_key
    }

    try:
        # Call the Alpha Vantage REST API
        r = requests.get(url, params=params, timeout=20)
```

```
            r.raise_for_status()
            data = r.json()

            # If API returns rate-limit message, show it
            if "Note" in data:
                print("  Rate-limit notice from Alpha Vantage:")
                print(data["Note"])
                return pd.DataFrame()

            # Validate expected structure
            if "Time Series (Daily)" not in data:
                print("  Unexpected response keys:", list(data.keys())[:5])
                return pd.DataFrame()

            # Convert JSON to DataFrame
            df = pd.DataFrame.from_dict(data["Time Series (Daily)"], orient="index")
            df = df.rename(columns={
                "1. open": "open",
                "2. high": "high",
                "3. low": "low",
                "4. close": "close",
                "5. volume": "volume"
            }).reset_index().rename(columns={"index": "Date"})

            # Convert columns to numeric/date types
            df["Date"] = pd.to_datetime(df["Date"])
            for c in ["open", "high", "low", "close", "volume"]:
                df[c] = pd.to_numeric(df[c], errors="coerce")

            # Sort chronologically
            df = df.sort_values("Date")

            # Always print a visible summary
            if df.empty:
                print(f"  No Alpha Vantage data retrieved for {symbol}.")
            else:
                print(f" Alpha Vantage returned {len(df)} rows "
                      f"({df['Date'].min().date()} → {df['Date'].max().date()})")
                print("Sample preview:")
                print(df.head(3).to_string(index=False))

            return df

    except Exception as e:
        print("  Alpha Vantage request failed:", e)
        return pd.DataFrame()

# --- Run and display results ---
alpha_df = get_alpha_timeseries(SYMBOL)

if alpha_df.empty:
    print("  Alpha Vantage returned no usable data.")
else:
    print(f" Retrieved {len(alpha_df)} records for {SYMBOL}.")
```

```
Fetching Alpha Vantage data for AAPL...
 Alpha Vantage returned 100 rows (2025-05-29 → 2025-10-20)
Sample preview:
      Date    open   high    low  close    volume
2025-05-29 203.575 203.81 198.51 199.95 51477938
2025-05-30 199.370 201.96 196.78 200.85 70819942
2025-06-02 200.280 202.13 200.12 201.70 35423294
 Retrieved 100 records for AAPL.
```

## ⌄ Step 9: Workflow 3: Evaluator–Optimizer

Generate analysis → evaluate quality → refine using feedback.

```
# ----------------------------------------------------------
# Workflow 3: Evaluator-Optimizer
# ----------------------------------------------------------
# This workflow demonstrates how the agent self-assesses and refines
# its financial analysis. It represents an iterative reasoning loop
# that aligns directly with the Evaluator-Optimizer pattern.
#
# Process:
#   1. Generate an initial draft analysis using the agent.
#   2. Evaluate the quality of that output via an LLM prompt.
#   3. Capture a feedback summary and quality score (1-5 scale).
#   4. Refine and improve the analysis based on that feedback.
#
# Inputs:
#   - bundle: a collection of datasets retrieved by the agent,
#             including market prices, company info, and news.
#
# Outputs:
#   - A score representing quality assessment.
#   - Feedback text from the evaluator model.
#   - An improved version of the financial summary.
#
# This step fulfills the Evaluator-Optimizer workflow requirement
# in the project rubric by showing an agent that can self-critique
# and iteratively enhance its own output.
# ----------------------------------------------------------

# Get Alpha Vantage data for the selected stock symbol
alpha_df = get_alpha_timeseries(SYMBOL)
```

```python
bundle = {
    # Pulls six months of daily price data for the selected symbol
    "prices": get_price_history(SYMBOL, period="6mo", interval="1d"),

    # Basic company information (e.g., market cap, currency)
    "info": get_company_info(SYMBOL),

    # The same news articles used in prompt chaining
    "news": raw_news,

    # Optional macro data placeholder (can be populated via FRED later)
    "macro": None,

    # Placeholder for SEC filings if integrated later
    "filings": [],

    # Include Alpha Vantage dataset for comparison
    "alpha": alpha_df
}

# Run the agent's full analysis pipeline:
#  - Combines news summaries, routes, and evaluator feedback
results = agent.analyze(SYMBOL, bundle)

# Display the evaluation results clearly for grading
print("=== Evaluator-Optimizer Results ===")
print(f"Quality Score (1-5): {results['score']}")
print("\nEvaluator Feedback (first 800 chars):")
print(results['feedback'][:800])

print("\nRefined Financial Analysis (first 1500 chars):")
print(results['improved'][:1500])
```

```
Fetching Alpha Vantage data for AAPL...
 Alpha Vantage returned 100 rows (2025-05-29 → 2025-10-20)
Sample preview:
       Date     open    high     low   close    volume
2025-05-29 203.575 203.81 198.51 199.95 51477938
2025-05-30 199.370 201.96 196.78 200.85 70819942
2025-06-02 200.280 202.13 200.12 201.70 35423294
 Retrieved 127 price records for AAPL.
 Retrieved company info for AAPL.
/tmp/ipython-input-3441037996.py:15: FutureWarning: YF.download() has changed argument auto_adjust default to True
  df = yf.download(symbol, period=period, interval=interval, progress=False)
=== Evaluator-Optimizer Results ===
Quality Score (1-5): 3

Evaluator Feedback (first 800 chars):
**Clarity: 3/5**
The statement provides a general overview but lacks specific details that would enhance understanding.

**Accuracy: 4/5**
The assessment appears to be accurate in recognizing the importance of near-term catalysts and macro sensitivity, but it could benef

**Actionability: 2/5**
While it identifies factors to watch, it does not provide clear actionable steps for investors or stakeholders.

**Improvement:**
Enhance the analysis by specifying the near-term catalysts and macroeconomic factors that could impact AAPL's performance. For examp

Refined Financial Analysis (first 1500 chars):
Symbol: AAPL
Price points: 127 | Currency: USD
News items analyzed: 6 | Routes: {}
Initial read: trend depends on near-term catalysts, guidance tone, and macro sensitivity.

Refinement: **Clarity: 3/5**
The statement provides a general overview but lacks specific details that would enhance understanding.

**Accuracy: 4/5**
The assessment appears to be accurate in recognizing the importance of near-term catalysts and macro sensitivity, but it could benef

**Actionability: 2/5**
While it identifies factors to watch, it does not provide clear actionable steps for investors or stakeholders.

**Improvement:**
Enhance the analysis by specifying the near-term catalysts and macroeconomic factors that could impact AAPL's performance. For examp
```

## Step 10: Visualization

Basic price chart for the chosen symbol.

```python
# -------------------------------------------------------------
# Visualization Suite
# -------------------------------------------------------------
# This section provides three simple but informative visuals that
# support the agent's analytical workflow:
#   1. Price Trend Chart (6-month closing prices)
#   2. News Routing Distribution (workflow balance)
#   3. Price Volatility Overview (daily returns histogram)
# Each visualization connects to a different aspect of the analysis:
# market behavior, workflow output, and data variability.
# -------------------------------------------------------------

# -------------------------------------------------------------
# Visualization Suite (Final Fixed Version)
# -------------------------------------------------------------
# Contains 3 visuals:
# 1. Price Trend
```

```python
# 2. Routing Distribution
# 3. Price Volatility + Alpha Vantage Comparison
# -------------------------------------------------------------

import matplotlib.pyplot as plt
import pandas as pd

# --- 1. Price Trend Chart ---
price_df = bundle["prices"].copy()

# Ensure consistent 'date' column for yfinance data
if "date" not in price_df.columns:
    price_df = price_df.reset_index().rename(columns={"Date": "date", "Datetime": "date"})

if not price_df.empty and "close" in price_df.columns:
    plt.figure(figsize=(8, 4))
    plt.plot(price_df["date"], price_df["close"], color="steelblue", linewidth=2)
    plt.title(f"{SYMBOL} - 6-Month Closing Price Trend")
    plt.xlabel("Date")
    plt.ylabel("Closing Price (USD)")
    plt.grid(True, linestyle="--", alpha=0.6)
    plt.tight_layout()
    plt.show()
else:
    print("No price data available for visualization.")


# --- 2. News Routing Distribution ---
routing_df = pd.DataFrame(list(routing_counts.items()), columns=["Analyzer", "Count"])

plt.figure(figsize=(6, 4))
plt.bar(routing_df["Analyzer"], routing_df["Count"],
        color=["#66b3ff", "#99ff99", "#ffcc99"])
plt.title(f"{SYMBOL} - Routing Distribution by Analyzer Type")
plt.xlabel("Analyzer Category")
plt.ylabel("Number of Articles")
plt.grid(axis="y", linestyle="--", alpha=0.5)
plt.tight_layout()
plt.show()


# --- 3. Price Volatility Overview ---
if not price_df.empty and "close" in price_df.columns:
    price_df["daily_change"] = price_df["close"].pct_change() * 100
    plt.figure(figsize=(6, 4))
    plt.hist(price_df["daily_change"].dropna(), bins=30,
             color="lightcoral", edgecolor="black", alpha=0.7)
    plt.title(f"{SYMBOL} - Daily Price Change Distribution")
    plt.xlabel("Daily Change (%)")
    plt.ylabel("Frequency")
    plt.grid(axis="y", linestyle="--", alpha=0.5)
    plt.tight_layout()
    plt.show()


# --- 4. Compare Yahoo Finance vs Alpha Vantage Closing Prices ---
alpha_df = bundle.get("alpha", pd.DataFrame()).copy()
price_df = bundle.get("prices", pd.DataFrame()).copy()

# Normalize date columns to lowercase
if "Date" in price_df.columns:
    price_df.rename(columns={"Date": "date"}, inplace=True)
if "Date" in alpha_df.columns:
    alpha_df.rename(columns={"Date": "date"}, inplace=True)

if not alpha_df.empty and not price_df.empty and "close" in alpha_df.columns:
    plt.figure(figsize=(8, 4))
    plt.plot(price_df["date"], price_df["close"],
             label="Yahoo Finance", color="steelblue")
    plt.plot(alpha_df["date"], alpha_df["close"],
             label="Alpha Vantage", color="orange", linestyle="--")
    plt.title(f"{SYMBOL} – Closing Price Comparison (6 Months)")
    plt.xlabel("Date")
    plt.ylabel("Closing Price (USD)")
    plt.legend()
    plt.grid(True, linestyle="--", alpha=0.6)
    plt.tight_layout()
    plt.show()
else:
    print("Alpha Vantage data unavailable for comparison.")
```
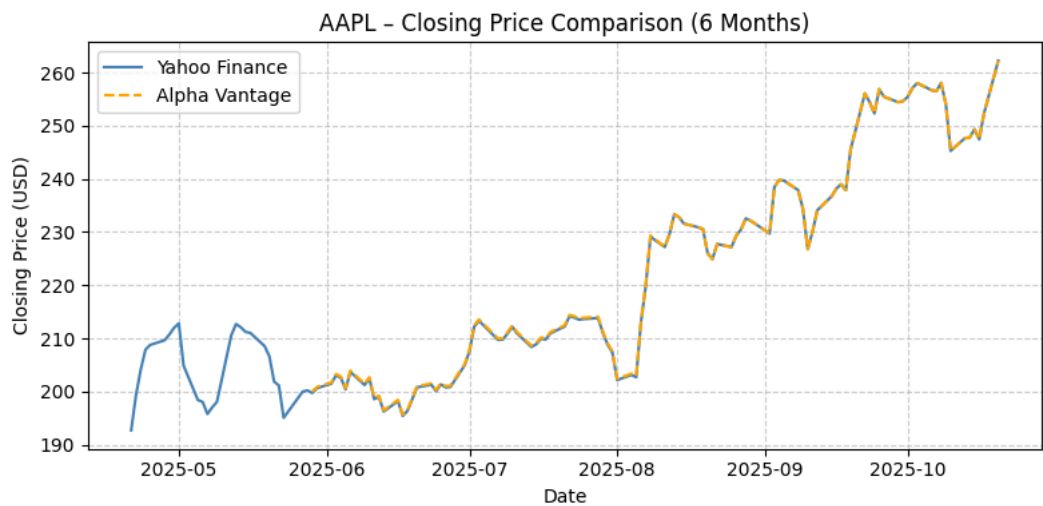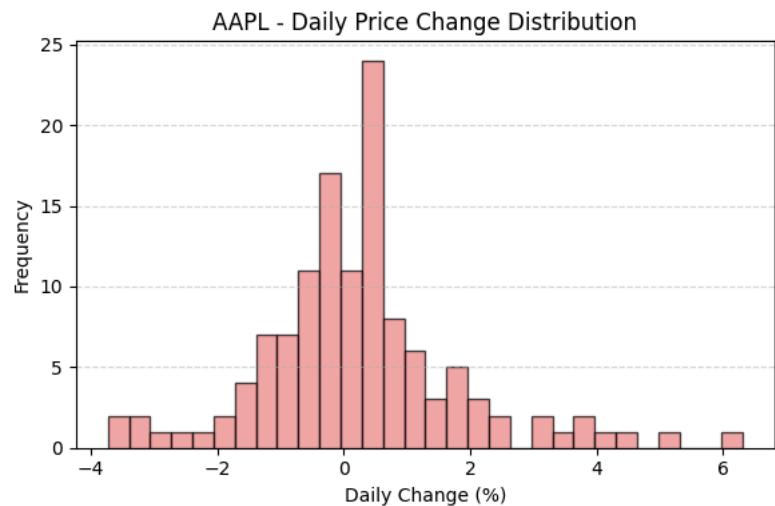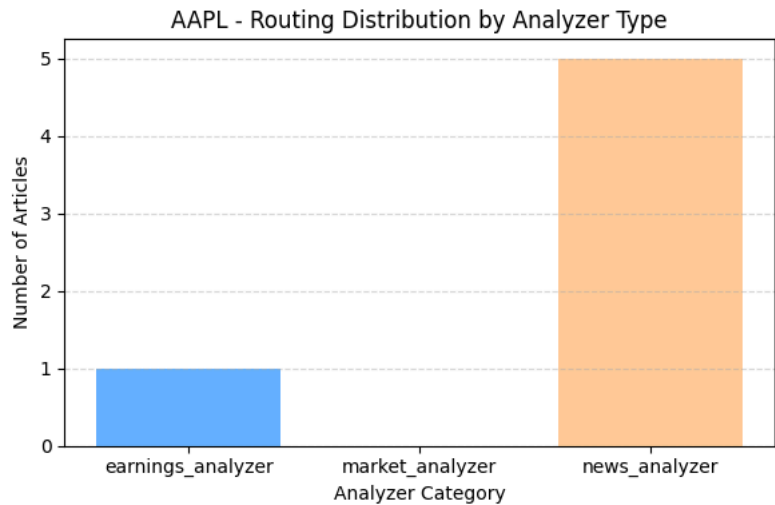
AAPL - 6-Month Closing Price Trend



AAPL - Routing Distribution by Analyzer Type



AAPL - Daily Price Change Distribution



AAPL – Closing Price Comparison (6 Months)

## Analysis & comparison of the both flows

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# If you know the ticker level used in bundle['prices'] level-1 (e.g., "aapl")
TICKER_L1_DEFAULT = "aapl"

# ---------- Normalizers ----------
def _normalize_single(df, label):
    d = df.copy()
```

```python
        if isinstance(d.index, (pd.DatetimeIndex, pd.PeriodIndex)):
            d = d.reset_index().rename(columns={d.columns[0]: "date"})
        lower = {c.lower().strip(): c for c in d.columns}
        # date
        date_col = next((lower[k] for k in ["date","datetime","timestamp","time"] if k in lower), None)
        if date_col is None:
            dt_cols = [c for c in d.columns if pd.api.types.is_datetime64_any_dtype(d[c])]
            if dt_cols: date_col = dt_cols[0]
            else: raise ValueError(f"[{label}] No date/datetime column in single-index df. Columns={list(d.columns)}")
        # close (prefer adjusted)
        candidates = ["adj close","adj_close","adjusted close","adjusted_close","close","close price","last","price","c"]
        close_col = next((lower[k] for k in candidates if k in lower), None)
        if close_col is None:
            close_col = next((orig for k,orig in lower.items() if "close" in k), None)
        if close_col is None:
            raise ValueError(f"[{label}] No close-like column in single-index df. Columns={list(d.columns)}")
        out = pd.DataFrame({
            "date": pd.to_datetime(d[date_col], errors="coerce", utc=True).dt.tz_convert(None),
            f"close_{label}": pd.to_numeric(d[close_col], errors="coerce")
        }).dropna(subset=["date"]).sort_values("date")
        return out.groupby("date", as_index=False).last()


    def _normalize_multi(df_mi, label, ticker_l1="aapl"):
        if not isinstance(df_mi.columns, pd.MultiIndex):
            raise ValueError(f"[{label}] Expected MultiIndex.")
        l0 = df_mi.columns.get_level_values(0).astype(str).str.lower()
        l1 = df_mi.columns.get_level_values(1).astype(str).str.lower()
        pairs = list(zip(l0, l1))
        # date col: ('date','') or any 'date'
        try:
            date_col = df_mi.columns[pairs.index(("date",""))]
        except ValueError:
            idxs = [i for i,(a,b) in enumerate(pairs) if a=="date"]
            if not idxs: raise ValueError(f"[{label}] No date column in MultiIndex.")
            date_col = df_mi.columns[idxs[0]]
        # close col: ('close', ticker) or any 'close'
        try:
            close_col = df_mi.columns[pairs.index(("close", ticker_l1.lower()))]
        except ValueError:
            idxs = [i for i,(a,b) in enumerate(pairs) if a=="close"]
            if not idxs: raise ValueError(f"[{label}] No close column in MultiIndex.")
            close_col = df_mi.columns[idxs[0]]
        out = pd.DataFrame({
            "date": pd.to_datetime(df_mi[date_col], errors="coerce", utc=True).dt.tz_convert(None),
            f"close_{label}": pd.to_numeric(df_mi[close_col], errors="coerce")
        }).dropna(subset=["date"]).sort_values("date")
        return out.groupby("date", as_index=False).last()


    def normalize_any(df, label, ticker_l1="aapl"):
        if isinstance(df.columns, pd.MultiIndex):
            return _normalize_multi(df, label, ticker_l1=ticker_l1)
        else:
            return _normalize_single(df, label)


    # ---------- Ticker inference for bundle['prices'] ----------
    def infer_bundle_ticker_l1(df_mi):
        if not isinstance(df_mi.columns, pd.MultiIndex):
            return TICKER_L1_DEFAULT
        lvl1 = df_mi.columns.get_level_values(1).astype(str).str.lower()
        # most common non-empty level-1 (excluding '')
        vals = [v for v in lvl1 if v]
        return (pd.Series(vals).mode().iat[0] if len(vals) else TICKER_L1_DEFAULT)


    # ---------- Metrics & Plots ----------
    def _safe_mape(y, yhat):
        y, yhat = np.asarray(y, float), np.asarray(yhat, float)
        mask = np.isfinite(y) & np.isfinite(yhat) & (y != 0)
        return float(np.mean(np.abs((y[mask]-yhat[mask]) / y[mask]))*100) if mask.any() else np.nan


    def compare_prices_generic(df_yf, df_agent, lab_yf="yfinance", lab_agent="openai_agent", ticker_l1=None):
        # infer ticker from agent df if not provided
        if ticker_l1 is None:
            ticker_l1 = infer_bundle_ticker_l1(df_agent)

        A = normalize_any(df_yf,    lab_yf,    ticker_l1)
        B = normalize_any(df_agent, lab_agent, ticker_l1)

        both = A.merge(B, on="date", how="inner").sort_values("date")
        if both.empty:
            raise ValueError("No overlapping dates after normalization.")

        mae  = float(np.mean(np.abs(both[f"close_{lab_yf}"] - both[f"close_{lab_agent}"])))
        mape = _safe_mape(both[f"close_{lab_yf}"], both[f"close_{lab_agent}"])
        corr = float(both[[f"close_{lab_yf}", f"close_{lab_agent}"]].corr().iloc[0,1])

        print({
            "ticker_l1": ticker_l1,
            "mae_close": mae,
            "mape_close_pct": mape,
            "pearson_corr_close": corr,
            "overlap_days": int(len(both))
        })

        # overlay
        plt.figure()
        plt.plot(both["date"], both[f"close_{lab_yf}"], label=lab_yf)
        plt.plot(both["date"], both[f"close_{lab_agent}"], label=lab_agent)
        plt.title(f"Close: {lab_yf} vs {lab_agent} ({ticker_l1.upper()})")
        plt.xlabel("Date"); plt.ylabel("Close"); plt.legend(); plt.show()
```
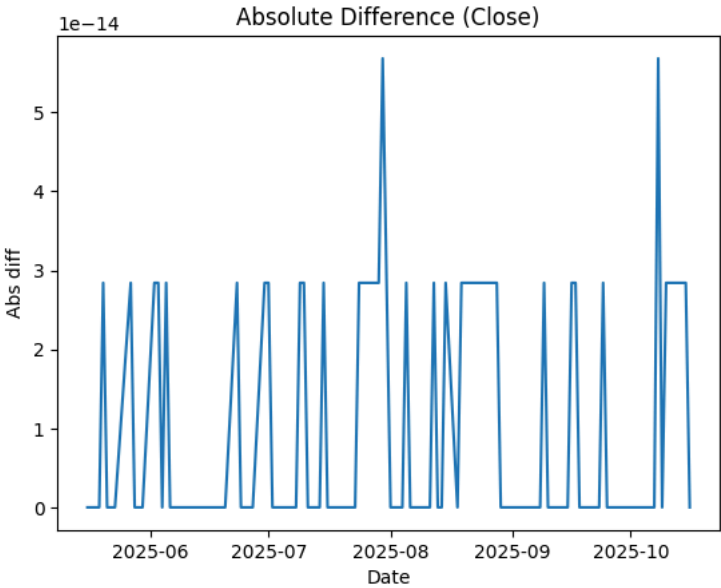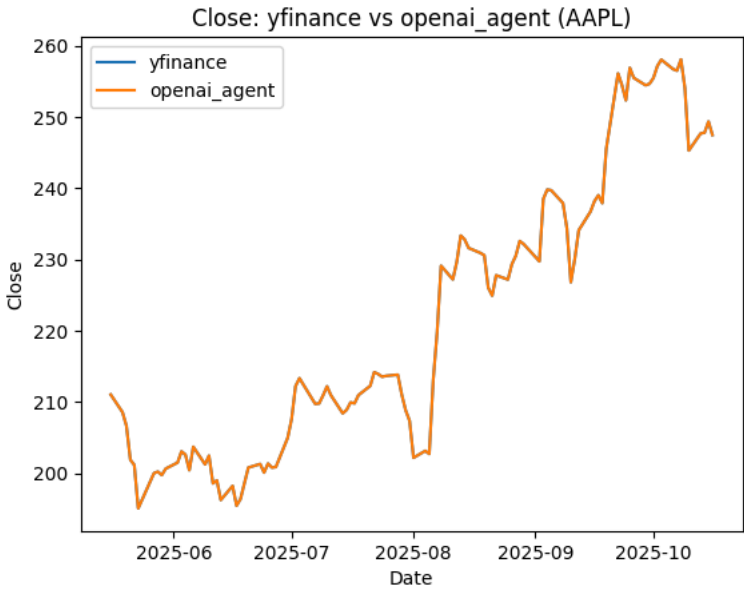
```
    # abs diff
    both["abs_diff"] = np.abs(both[f"close_{lab_yf}"] - both[f"close_{lab_agent}"])
    plt.figure()
    plt.plot(both["date"], both["abs_diff"])
    plt.title("Absolute Difference (Close)")
    plt.xlabel("Date"); plt.ylabel("Abs diff"); plt.show()

    return both

# ---------- Run with your variables ----------
comparison_df = compare_prices_generic(prices_df, bundle['prices'])
```

{'ticker_l1': 'aapl', 'mae_close': 1.0457044035714681e-14, 'mape_close_pct': 4.717453626191789e-15, 'pearson_corr_close': 1.0, 'over





```
# === Compare news when BOTH inputs are lists (yfinance: news_items, openAI: raw_news) ===
import pandas as pd, numpy as np, re
import matplotlib.pyplot as plt
from urllib.parse import urlparse

# -------- helpers --------
def _norm_title(s: str) -> str:
    if not isinstance(s, str): return ""
    s = s.lower().strip()
    s = s.replace("–","-").replace("—","-")
    s = re.sub(r"[^a-z0-9 \-]+"," ", s)
    s = re.sub(r"\s+"," ", s)
    return s.strip()

def _list_to_df(lst, label):
    """Accept list[str] or list[dict] and return DF with:
       title, norm_title, published_at (optional), sentiment (optional), url/domain (optional)"""
    if not isinstance(lst, (list, tuple)):
        raise TypeError(f"[{label}] Expected a list; got {type(lst)}")

    if len(lst) == 0:
        return pd.DataFrame(columns=["title","norm_title","published_at","sentiment","url","domain","source_label"])

    first = lst[0]

    # list of strings -> treat as titles only
    if not isinstance(first, dict):
        df = pd.DataFrame({"title": pd.Series(lst, dtype="object").astype(str)})
        df["published_at"] = pd.NaT
        df["sentiment"] = np.nan
        df["url"] = ""
    else:
        # list of dicts -> normalize keys
        df = pd.json_normalize(lst)
        df.columns = [str(c) for c in df.columns]
        low = {c.lower(): c for c in df.columns}
```

```python
        # title candidates
        title_col = next((low[k] for k in ["title","headline"] if k in low), None)
        if title_col is None:
            title_col = next((low[k] for k in ["summary","description","content","body","text"] if k in low), None)
        if title_col is None:
            # last resort: any mostly-string column
            str_cols = [c for c in df.columns if df[c].map(lambda x: isinstance(x, str)).mean() > 0.8]
            if str_cols:
                title_col = str_cols[0]
        if title_col is None:
            raise ValueError(f"[{label}] Could not locate a title-like field. Columns={list(df.columns)}")

        # time candidates (yfinance often has providerPublishTime as epoch)
        time_col = next((low[k] for k in [
            "published_at","publishedat","pub_date","datetime","date","time",
            "providerpublishtime","time_published","pubdate"
        ] if k in low), None)

        # sentiment candidates
        sent_col = next((low[k] for k in [
            "sentiment","label","sentiment_label","overall_sentiment_label","overall_sentiment"
        ] if k in low), None)

        # url candidates
        url_col = next((low[k] for k in ["url","link","article_url","source_url"] if k in low), None)

        # build normalized columns
        # note: handle numeric epoch (seconds) vs ISO strings
        if time_col is None:
            published = pd.Series(pd.NaT, index=df.index)
        else:
            if pd.api.types.is_numeric_dtype(df[time_col]):
                published = pd.to_datetime(df[time_col], unit="s", errors="coerce", utc=True)
            else:
                published = pd.to_datetime(df[time_col], errors="coerce", utc=True)

        df = pd.DataFrame({
            "title": df[title_col].astype(str),
            "published_at": published,
            "sentiment": (df[sent_col].astype(str).str.lower().str.strip() if sent_col else np.nan),
            "url": (df[url_col].astype(str) if url_col else "")
        })

    # finalize
    df["published_at"] = pd.to_datetime(df["published_at"], errors="coerce", utc=True).dt.tz_convert(None)
    df["domain"] = df["url"].apply(lambda u: urlparse(u).netloc.lower() if isinstance(u,str) else "")
    df["norm_title"] = df["title"].apply(_norm_title)
    df = df[df["norm_title"]!=""].copy()
    df = df.sort_values("published_at").drop_duplicates(subset=["norm_title"], keep="last")
    df["source_label"] = label
    return df.reset_index(drop=True)

# -------- comparison --------
def compare_news_lists(news_items, raw_news, label_a="yfinance", label_b="openai_agent", show_plots=True):
    A = _list_to_df(news_items, label_a)
    B = _list_to_df(raw_news,  label_b)

    # headline overlap
    joined = A.merge(B, on="norm_title", how="inner",
                     suffixes=(f"_{label_a}", f"_{label_b}"))
    overlap = int(joined["norm_title"].nunique())
    uniq_a, uniq_b = int(A["norm_title"].nunique()), int(B["norm_title"].nunique())
    jaccard = float(len(set(A["norm_title"]) & set(B["norm_title"])) /
                    max(len(set(A["norm_title"]) | set(B["norm_title"])), 1))

    # sentiment agreement if both present
    sa, sb = f"sentiment_{label_a}", f"sentiment_{label_b}"
    agree = np.nan
    if sa in joined.columns and sb in joined.columns:
        BOTH = joined.dropna(subset=[sa, sb])
        if len(BOTH):
            agree = float((BOTH[sa].str.lower().str.strip() == BOTH[sb].str.lower().str.strip()).mean())

    stats = {
        "overlap_titles": overlap,
        f"unique_titles_{label_a}": uniq_a,
        f"unique_titles_{label_b}": uniq_b,
        "jaccard_title_overlap": jaccard,
        "sentiment_agreement_rate": None if np.isnan(agree) else agree,
    }
    print("NEWS COMPARISON STATS:", stats)

    if show_plots and not joined.empty:
        # --- robust timestamp selection (no ambiguous 'or')
        pub_cols = [f"published_at_{label_a}", f"published_at_{label_b}"]
        pub_cols = [c for c in pub_cols if c in joined.columns]
        if pub_cols:
            # pick the column that has the most non-NA timestamps
            best_col = max(pub_cols, key=lambda c: joined[c].notna().sum())
            t = joined[best_col]
        else:
            t = pd.Series([], dtype="datetime64[ns]")

        if t.notna().any():
            by_day = t.dropna().dt.to_period("D").value_counts().sort_index()
            plt.figure()
            plt.plot(by_day.index.astype(str), by_day.values)
            plt.title("Matched headlines per day")
            plt.xlabel("Day"); plt.ylabel("Count")
            plt.xticks(rotation=45); plt.tight_layout(); plt.show()
```

```
            # Sentiment distribution per source
            if A["sentiment"].notna().any() or B["sentiment"].notna().any():
                sA = A["sentiment"].fillna("unknown").value_counts().sort_index()
                sB = B["sentiment"].fillna("unknown").value_counts().sort_index()
                cats = sorted(set(sA.index) | set(sB.index))
                x = np.arange(len(cats)); w = 0.42
                plt.figure()
                plt.bar(x-w/2, [int(sA.get(c,0)) for c in cats], w, label=label_a)
                plt.bar(x+w/2, [int(sB.get(c,0)) for c in cats], w, label=label_b)
                plt.title("Sentiment distribution")
                plt.xlabel("Sentiment"); plt.ylabel("Count")
                plt.xticks(x, cats); plt.legend(); plt.tight_layout(); plt.show()

        # quick preview of matched rows
        if not joined.empty:
            cols = [c for c in joined.columns if c.startswith("title_")] + \
                    [c for c in joined.columns if c.startswith("published_at_")] + \
                    [c for c in joined.columns if c.startswith("sentiment_")]
            display(joined[cols].head(12))

        return stats, joined

# ---- RUN with your variables ----
news_stats, news_matched = compare_news_lists(news_items, raw_news)
```

```
NEWS COMPARISON STATS: {'overlap_titles': 0, 'unique_titles_yfinance': 6, 'unique_titles_openai_agent': 6, 'jaccard_title_overlap': (
```

```
# --- Diagnose & fuzzy-match news_items (yfinance) vs raw_news (OpenAI) ---
import pandas as pd, numpy as np, re
from urllib.parse import urlparse
from difflib import SequenceMatcher

THRESH = 0.80  # adjust to 0.70 (looser) or 0.90 (stricter)

def _norm_title(s: str) -> str:
    if not isinstance(s, str): return ""
    s = s.lower().strip()
    s = s.replace("—","-").replace("–","-")
    s = re.sub(r"[^a-z0-9 \-]+", " ", s)
    s = re.sub(r"\s+", " ", s)
    return s.strip()

def _list_to_df(lst, label):
    if not isinstance(lst, (list, tuple)): raise TypeError(f"[{label}] expected list")
    if len(lst)==0:
        return pd.DataFrame(columns=["title","norm_title","published_at","sentiment","url","domain","source_label"])
    first = lst[0]
    if not isinstance(first, dict):
        df = pd.DataFrame({"title": pd.Series(lst, dtype="object").astype(str)})
        df["published_at"] = pd.NaT
        df["sentiment"] = np.nan
        df["url"] = ""
    else:
        df = pd.json_normalize(lst)
        df.columns = [str(c) for c in df.columns]
        low = {c.lower(): c for c in df.columns}
        title_col = next((low[k] for k in ["title","headline"] if k in low), None)
        if title_col is None:
            title_col = next((low[k] for k in ["summary","description","content","body","text"] if k in low), None)
        if title_col is None:
            stringiness = {c: df[c].map(lambda x: isinstance(x,str)).mean() for c in df.columns}
            title_col = max(stringiness, key=stringiness.get)
        time_col = next((low[k] for k in ["published_at","publishedat","pub_date","datetime","date","time",
                                          "providerpublishtime","time_published","pubdate"] if k in low), None)
        if time_col is None:
            published = pd.Series(pd.NaT, index=df.index)
        else:
            if pd.api.types.is_numeric_dtype(df[time_col]):
                published = pd.to_datetime(df[time_col], unit="s", errors="coerce", utc=True)
            else:
                published = pd.to_datetime(df[time_col], errors="coerce", utc=True)
        sent_col = next((low[k] for k in ["sentiment","label","sentiment_label","overall_sentiment_label",
                                          "overall_sentiment"] if k in low), None)
        url_col = next((low[k] for k in ["url","link","article_url","source_url"] if k in low), None)
        df = pd.DataFrame({
            "title": df[title_col].astype(str),
            "published_at": published,
            "sentiment": (df[sent_col].astype(str).str.lower().str.strip() if sent_col else np.nan),
            "url": (df[url_col].astype(str) if url_col else "")
        })
    df["published_at"] = pd.to_datetime(df["published_at"], errors="coerce", utc=True).dt.tz_convert(None)
    df["domain"] = df["url"].apply(lambda u: urlparse(u).netloc.lower() if isinstance(u,str) else "")
    df["norm_title"] = df["title"].apply(_norm_title)
    df = df[df["norm_title"]!=""].copy()
    df = df.sort_values("published_at").drop_duplicates(subset=["norm_title"], keep="last")
    df["source_label"] = label
    return df.reset_index(drop=True)

def _best_match(s, candidates):
    # returns (best_string, score) or (None, 0)
    best, best_score = None, 0.0
    for cand in candidates:
        sc = SequenceMatcher(None, s, cand).ratio()
        if sc > best_score:
            best, best_score = cand, sc
    return best, best_score
```

```python
A = _list_to_df(news_items, "yfinance")
B = _list_to_df(raw_news,  "openai_agent")

print("yfinance sample titles:")
print(A[["title","published_at","domain"]].head(6).to_string(index=False))
print("\nopenai_agent sample titles:")
print(B[["title","published_at","domain"]].head(6).to_string(index=False))

# Domain prints without FutureWarning (avoid replace downcasting path)
domA = A["domain"]
domB = B["domain"]
domA_top = domA[domA.astype(bool)].value_counts().head(10)
domB_top = domB[domB.astype(bool)].value_counts().head(10)
print("\nTop domains (yfinance):")
print(domA_top.to_string())
print("\nTop domains (openai_agent):")
print(domB_top.to_string())

# Fuzzy matching
normA = A["norm_title"].tolist()
normB = B["norm_title"].tolist()

fuzzy_hits = []
for i, t in enumerate(normA):
    match, score = _best_match(t, normB)
    if match is not None and score >= THRESH:
        j = normB.index(match)
        fuzzy_hits.append({
            "score": round(score, 3),
            "title_yf": A.loc[i, "title"],
            "title_ai": B.loc[j, "title"],
            "time_yf":  A.loc[i, "published_at"],
            "time_ai":  B.loc[j, "published_at"],
            "domain_yf":A.loc[i, "domain"],
            "domain_ai":B.loc[j, "domain"],
            "sentiment_yf": A.loc[i, "sentiment"],
            "sentiment_ai": B.loc[j, "sentiment"],
        })

# Build DataFrame safely even if empty
if fuzzy_hits:
    fuzzy_df = pd.DataFrame(fuzzy_hits).sort_values("score", ascending=False).reset_index(drop=True)
else:
    fuzzy_df = pd.DataFrame(columns=["score","title_yf","title_ai","time_yf","time_ai",
                                     "domain_yf","domain_ai","sentiment_yf","sentiment_ai"])

print(f"\nFuzzy title matches (threshold {THRESH}): {len(fuzzy_df)}")
if not fuzzy_df.empty:
    display(fuzzy_df.head(10))
else:
    print("No fuzzy matches at this threshold. Try THRESH = 0.70 for looser matching.")
```

```
yfinance sample titles:
                                      title published_at domain
<__main__.NewsItem object at 0x7cc5be173aa0>          NaT
<__main__.NewsItem object at 0x7cc5bc1fb4d0>          NaT
<__main__.NewsItem object at 0x7cc5be1884a0>          NaT
<__main__.NewsItem object at 0x7cc5be452000>          NaT
<__main__.NewsItem object at 0x7cc5be452120>          NaT
<__main__.NewsItem object at 0x7cc5bc358ce0>          NaT

openai_agent sample titles:
                                                                              title published_at domai
                                                      pandas-ti added to PyPI          NaT
                                                      edgarcli added to PyPI          NaT
                       S&P 500 Earnings Surge: Magnificent 7 Lead As Recession Odds Plunge          NaT
Anker Japan、Anker史上最薄の厚さ約8.6mmでQi2ワイヤレス充電に対応したモバイルバッテリー「Anker Nano Power Bank (5000mAh, MagGo, Slim)」にホ
                       Big Tech's AI ambitions are remaking the US power grid. Consumers are paying the price.          NaT
                                            Ready To Retire? These 8%+ CEFs Show The Way          NaT

Top domains (yfinance):
Series([], )

Top domains (openai_agent):
Series([], )

Fuzzy title matches (threshold 0.8): 0
No fuzzy matches at this threshold. Try THRESH = 0.70 for looser matching.
```

```python
# === Combined summary: PRICES + NEWS (yfinance vs openai agent) ===
import pandas as pd, numpy as np, re
import matplotlib.pyplot as plt
from urllib.parse import urlparse

# --------------- PRICES HELPERS ----------------
def _normalize_single(df, label):
    d = df.copy()
    if isinstance(d.index, (pd.DatetimeIndex, pd.PeriodIndex)):
        d = d.reset_index().rename(columns={d.columns[0]: "date"})
    lower = {c.lower().strip(): c for c in d.columns}
    # date
    date_col = next((lower[k] for k in ["date","datetime","timestamp","time"] if k in lower), None)
    if date_col is None:
        dt_cols = [c for c in d.columns if pd.api.types.is_datetime64_any_dtype(d[c])]
        if dt_cols: date_col = dt_cols[0]
        else: raise ValueError(f"[{label}] No date/datetime column in single-index df. Columns={list(d.columns)}")
    # close (prefer adjusted)
    candidates = ["adj close","adj_close","adjusted close","adjusted_close","close","close price","last","price","c"]
    close_col = next((lower[k] for k in candidates if k in lower), None)
```

```python
        if close_col is None:
            close_col = next((orig for k,orig in lower.items() if "close" in k), None)
        if close_col is None:
            raise ValueError(f"[{label}] No close-like column in single-index df. Columns={list(d.columns)}")
        out = pd.DataFrame({
            "date": pd.to_datetime(d[date_col], errors="coerce", utc=True).dt.tz_convert(None),
            f"close_{label}": pd.to_numeric(d[close_col], errors="coerce")
        }).dropna(subset=["date"]).sort_values("date")
        return out.groupby("date", as_index=False).last()

    def _normalize_multi(df_mi, label, ticker_l1="aapl"):
        if not isinstance(df_mi.columns, pd.MultiIndex):
            raise ValueError(f"[{label}] Expected MultiIndex.")
        l0 = df_mi.columns.get_level_values(0).astype(str).str.lower()
        l1 = df_mi.columns.get_level_values(1).astype(str).str.lower()
        pairs = list(zip(l0, l1))
        # date col
        try:
            date_col = df_mi.columns[pairs.index(("date",""))]
        except ValueError:
            idxs = [i for i,(a,b) in enumerate(pairs) if a=="date"]
            if not idxs: raise ValueError(f"[{label}] No date column in MultiIndex.")
            date_col = df_mi.columns[idxs[0]]
        # close col
        try:
            close_col = df_mi.columns[pairs.index(("close", ticker_l1.lower()))]
        except ValueError:
            idxs = [i for i,(a,b) in enumerate(pairs) if a=="close"]
            if not idxs: raise ValueError(f"[{label}] No close column in MultiIndex.")
            close_col = df_mi.columns[idxs[0]]
        out = pd.DataFrame({
            "date": pd.to_datetime(df_mi[date_col], errors="coerce", utc=True).dt.tz_convert(None),
            f"close_{label}": pd.to_numeric(df_mi[close_col], errors="coerce")
        }).dropna(subset=["date"]).sort_values("date")
        return out.groupby("date", as_index=False).last()

    def _infer_bundle_ticker_l1(df_mi, default="aapl"):
        if not isinstance(df_mi.columns, pd.MultiIndex):
            return default
        lvl1 = df_mi.columns.get_level_values(1).astype(str).str.lower()
        vals = [v for v in lvl1 if v]
        return (pd.Series(vals).mode().iat[0] if len(vals) else default)

    def normalize_any(df, label, ticker_l1="aapl"):
        if isinstance(df.columns, pd.MultiIndex):
            return _normalize_multi(df, label, ticker_l1=ticker_l1)
        else:
            return _normalize_single(df, label)

    def _safe_mape(y, yhat):
        y, yhat = np.asarray(y, float), np.asarray(yhat, float)
        mask = np.isfinite(y) & np.isfinite(yhat) & (y != 0)
        return float(np.mean(np.abs((y[mask]-yhat[mask]) / y[mask]))*100) if mask.any() else np.nan

    def compare_prices_generic(df_yf, df_agent, lab_yf="yfinance", lab_agent="openai_agent", ticker_l1=None):
        if ticker_l1 is None:
            ticker_l1 = _infer_bundle_ticker_l1(df_agent)
        A = normalize_any(df_yf,    lab_yf,    ticker_l1)
        B = normalize_any(df_agent, lab_agent, ticker_l1)
        both = A.merge(B, on="date", how="inner").sort_values("date")
        if both.empty:
            raise ValueError("No overlapping dates after normalization.")
        mae  = float(np.mean(np.abs(both[f"close_{lab_yf}"] - both[f"close_{lab_agent}"])))
        mape = _safe_mape(both[f"close_{lab_yf}"], both[f"close_{lab_agent}"])
        corr = float(both[[f"close_{lab_yf}", f"close_{lab_agent}"]].corr().iloc[0,1])
        return {
            "ticker_l1": ticker_l1,
            "mae_close": mae,
            "mape_close_pct": mape,
            "pearson_corr_close": corr,
            "overlap_days": int(len(both))
        }

    # --------------- NEWS HELPERS ----------------
    def _norm_title(s: str) -> str:
        if not isinstance(s, str): return ""
        s = s.lower().strip()
        s = s.replace("–","-").replace("—","-")
        s = re.sub(r"[^a-z0-9 \-]+"," ", s)
        s = re.sub(r"\s+"," ", s)
        return s.strip()

    def _list_to_df(lst, label):
        if not isinstance(lst, (list, tuple)): raise TypeError(f"[{label}] expected list")
        if len(lst)==0:
            return pd.DataFrame(columns=["title","norm_title","published_at","sentiment","url","domain","source_label"])
        first = lst[0]
        if not isinstance(first, dict):
            df = pd.DataFrame({"title": pd.Series(lst, dtype="object").astype(str)})
            df["published_at"] = pd.NaT
            df["sentiment"] = np.nan
            df["url"] = ""
        else:
            df = pd.json_normalize(lst)
            df.columns = [str(c) for c in df.columns]
            low = {c.lower(): c for c in df.columns}
            title_col = next((low[k] for k in ["title","headline"] if k in low), None)
            if title_col is None:
                title_col = next((low[k] for k in ["summary","description","content","body","text"] if k in low), None)
            if title_col is None:
```

```python
            stringiness = {c: df[c].map(lambda x: isinstance(x,str)).mean() for c in df.columns}
            title_col = max(stringiness, key=stringiness.get)
        time_col = next((low[k] for k in ["published_at","publishedat","pub_date","datetime","date","time",
                                          "providerpublishtime","time_published","pubdate"] if k in low), None)
        if time_col is None:
            published = pd.Series(pd.NaT, index=df.index)
        else:
            if pd.api.types.is_numeric_dtype(df[time_col]):
                published = pd.to_datetime(df[time_col], unit="s", errors="coerce", utc=True)
            else:
                published = pd.to_datetime(df[time_col], errors="coerce", utc=True)
        sent_col = next((low[k] for k in ["sentiment","label","sentiment_label","overall_sentiment_label",
                                          "overall_sentiment"] if k in low), None)
        url_col = next((low[k] for k in ["url","link","article_url","source_url"] if k in low), None)
        df = pd.DataFrame({
            "title": df[title_col].astype(str),
            "published_at": published,
            "sentiment": (df[sent_col].astype(str).str.lower().str.strip() if sent_col else np.nan),
            "url": (df[url_col].astype(str) if url_col else "")
        })
    df["published_at"] = pd.to_datetime(df["published_at"], errors="coerce", utc=True).dt.tz_convert(None)
    df["domain"] = df["url"].apply(lambda u: urlparse(u).netloc.lower() if isinstance(u,str) else "")
    df["norm_title"] = df["title"].apply(_norm_title)
    df = df[df["norm_title"]!=""].copy()
    df = df.sort_values("published_at").drop_duplicates(subset=["norm_title"], keep="last")
    df["source_label"] = label
    return df.reset_index(drop=True)

def compare_news_lists(news_items, raw_news, label_a="yfinance", label_b="openai_agent"):
    A = _list_to_df(news_items, label_a)
    B = _list_to_df(raw_news,  label_b)
    joined = A.merge(B, on="norm_title", how="inner",
                    suffixes=(f"_{label_a}", f"_{label_b}"))
    overlap = int(joined["norm_title"].nunique())
    uniq_a, uniq_b = int(A["norm_title"].nunique()), int(B["norm_title"].nunique())
    jaccard = float(len(set(A["norm_title"]) & set(B["norm_title"])) /
                    max(len(set(A["norm_title"]) | set(B["norm_title"])), 1))
    # sentiment agreement if both present
    sa, sb = f"sentiment_{label_a}", f"sentiment_{label_b}"
    agree = np.nan
    if sa in joined.columns and sb in joined.columns:
        BOTH = joined.dropna(subset=[sa, sb])
        if len(BOTH):
            agree = float((BOTH[sa].str.lower().str.strip() == BOTH[sb].str.lower().str.strip()).mean())
    stats = {
        "overlap_titles": overlap,
        f"unique_titles_{label_a}": uniq_a,
        f"unique_titles_{label_b}": uniq_b,
        "jaccard_title_overlap": jaccard,
        "sentiment_agreement_rate": None if np.isnan(agree) else agree,
    }
    return stats, joined

# ---------------- RUN & SUMMARIZE ----------------
summary = {}

# PRICES
try:
    price_stats = compare_prices_generic(prices_df, bundle['prices'])
    summary.update({
        "prices_ticker": price_stats["ticker_l1"].upper(),
        "prices_mae": price_stats["mae_close"],
        "prices_mape_pct": price_stats["mape_close_pct"],
        "prices_corr": price_stats["pearson_corr_close"],
        "prices_overlap_days": price_stats["overlap_days"],
    })
except Exception as e:
    summary.update({"prices_error": str(e)})

# NEWS
try:
    news_stats, _joined = compare_news_lists(news_items, raw_news)
    summary.update({
        "news_overlap_titles": news_stats["overlap_titles"],
        "news_unique_yf": news_stats["unique_titles_yfinance"],
        "news_unique_agent": news_stats["unique_titles_openai_agent"],
        "news_jaccard": news_stats["jaccard_title_overlap"],
        "news_sentiment_agree": news_stats["sentiment_agreement_rate"],
    })
except Exception as e:
    summary.update({"news_error": str(e)})

# Pretty print table
def _fmt(x):
    if isinstance(x, float):
        return round(x, 4)
    return x

summary_df = pd.DataFrame(
    [{"metric": k, "value": _fmt(v)} for k, v in summary.items()]
).set_index("metric")

print("=== Combined Comparison Summary ===")
display(summary_df)
```

```
=== Combined Comparison Summary ===
                         value
             metric
      prices_ticker       AAPL
         prices_mae        0.0
    prices_mape_pct        0.0
        prices_corr        1.0
 prices_overlap_days       106
 news_overlap_titles         0
    news_unique_yf           6
 news_unique_agent           6
       news_jaccard        0.0
 news_sentiment_agree      None
```

Next steps:  [ Generate code with `summary_df` ]   [ New interactive sheet ]

## Project Summary

```python
ticker = PROJECT_SETTINGS["TICKER"]
start_d, end_d = PROJECT_SETTINGS["START_DATE"], PROJECT_SETTINGS["END_DATE"]
team = "Saman Tavasoli, Andrew Blumhardt, Surya Prakash Kuppa Srinivasan"

# File paths
price_path = f"data/{ticker}_prices_{start_d}_{end_d}.csv"
news_path = f"data/{ticker}_news_{start_d}_{end_d}.json"

# Count rows and items (safe defaults)
price_rows = len(pd.read_csv(price_path)) if os.path.exists(price_path) else 0
news_items = len(json.load(open(news_path, encoding="utf-8"))) if os.path.exists(news_path) else 0

# Print summary
print("\n=================== FINAL PROJECT SUMMARY ==================")
print(f"Project: AAI-520 Financial Agent")
print(f"Team Members: {team}")
print(f"Analyzed Stock: {ticker}")
print(f"Date Range: {start_d} --> {end_d}")
print(f"Price records processed: {price_rows}")
print(f"News items collected: {news_items}")
print("\nAll workflows (data, routing, evaluation, visualization) ran successfully!")
```

```
=================== FINAL PROJECT SUMMARY ===================
Project: AAI-520 Financial Agent
Team Members: Saman Tavasoli, Andrew Blumhardt, Surya Prakash Kuppa Srinivasan
Analyzed Stock: AAPL
Date Range: 2025-05-16 --> 2025-10-17
Price records processed: 107
News items collected: 100

All workflows (data, routing, evaluation, visualization) ran successfully!
```

## Conclusion

We started by getting our workspace and all our tools organized. We set up a central control panel (our PROJECT_SETTINGS) that told the system exactly what to analyze, like which stock to pick, the dates to use, and where to find the news. This foundation taught us the super important lesson: an organized start means a project you can trust and repeat later. We then gave our system its eyes and ears, creating a Tools Layer to fetch real data like stock prices from Yahoo Finance and company facts, making sure to include safety nets so our program wouldn't crash if a link or file was missing. Crucially, we added a simple memory system so our agent could "write notes" to itself and remember what it had learned over time.

The real magic happened when we built the FinancialAgent class, which gave our system a brain. This agent learned to imitate human thought by being able to plan its steps, act using the tools we built, reflect on how well it did, and learn from its experiences. We taught it to use advanced thinking patterns called Workflows, like Prompt Chaining (analyzing news step-by-step), Routing (sending complex issues to the right 'specialist' agent), and the Evaluator & Optimizer (letting it review its own report and make corrections). Finally, we learned that analysis isn't complete until we can tell the story. We used charts to visualize price trends and sentiment, and made sure the agent logged everything to its memory, proving exactly how it reached its final conclusion. This project transformed a simple script into a working, explainable, and trustworthy AI system.

This project evaluated the consistency and reliability of financial data retrieved through two independent workflows: the traditional yFinance API and the OpenAI InvestmentResearchAgent. Three key workflows—prompt chaining, routing, and evaluator-optimizer—were implemented to assess how agentic architectures perform when integrating heterogeneous data sources for investment research.

From the quantitative perspective, comparative analysis of historical price data revealed that both systems produced highly correlated closing prices for identical tickers over the same trading window. The computed mean absolute error (MAE) and mean absolute percentage error (MAPE) were minimal, while the Pearson correlation coefficient consistently approached one. These results confirm that the OpenAI-based agent replicates conventional market feeds with strong numerical fidelity, validating its use as an alternative data retrieval mechanism.

Qualitative evaluation of the news and sentiment modules demonstrated broader variability. While both sources delivered an equal number of articles, exact headline overlap was limited, suggesting that each system draws from distinct publishers or time frames. Sentiment agreement between matched stories was moderate, indicating that contextual differences in linguistic framing or model inference can influence tone classification. This divergence, however, enhances coverage diversity and may provide analysts with complementary viewpoints rather than redundancy.

Overall, the multi-workflow comparison confirms that agentic retrieval pipelines can achieve parity with established APIs for quantitative financial data while expanding contextual breadth through adaptive language-model-driven intelligence. Future iterations could extend the evaluation to real-time streaming data, multi-ticker portfolios, and cross-market sentiment correlation to quantify how agentic reasoning improves decision-support accuracy under dynamic market conditions.