## Summary of Work Completed

1. **Data Loading and EDA**:
   - The dataset is downloaded using `kagglehub`, and an EDA report is generated using `ydata_profiling` to understand the data distribution and characteristics.
2. **Feature Engineering**:
   - The data is preprocessed by normalizing pixel values, reshaping images, and scaling features using `MinMaxScaler`. Additional features are engineered using edge detection with the Sobel filter.
3. **Model Training and Evaluation**:
   - A baseline `RandomForestClassifier` is trained and evaluated. SHAP (SHapley Additive exPlanations) is used to explain the model's predictions and identify important features.
4. **AutoML with H2O**:
   - The H2O AutoML framework is used to automatically select the best model from a range of algorithms, optimizing for accuracy.
5. **Hyperparameter Tuning with Optuna**:
   - Optuna is used to fine-tune the hyperparameters of the `RandomForestClassifier`, maximizing accuracy.
6. **Model Tracking and Monitoring with MLflow**:
   - MLflow is used to log parameters, metrics, and artifacts for each run. Drift detection is implemented to monitor model performance over time.

## Justification for Choice of Libraries/Tools

1. **Kagglehub**:
   - Used for downloading the Fashion MNIST dataset, ensuring easy access to the data.
2. **YData Profiling**:
   - Generates an EDA report, providing insights into the dataset's structure and characteristics, which is crucial for understanding the data before modeling.
3. **TensorFlow and Scikit-learn**:
   - TensorFlow is used for loading the Fashion MNIST dataset, while Scikit-learn provides tools for preprocessing, model training, and evaluation.
4. **SHAP**:
   - Provides explainability by visualizing feature importance, helping to understand the model's decision-making process.
5. **H2O AutoML**:
   - Automates the process of model selection and hyperparameter tuning, saving time and effort while ensuring optimal model performance.
6. **Optuna**:
   - A powerful hyperparameter optimization library that efficiently searches the hyperparameter space to improve model accuracy.
7. **MLflow**:

o Tracks experiments, logs parameters and metrics, and stores models, ensuring reproducibility and facilitating model monitoring.

**Explanation of How Each Functionality Contributes to MLOps Best Practices**

1. **Data Loading and EDA**:
   o Ensures data quality and understanding, which is critical for building reliable models. The EDA report helps identify potential issues like missing values or imbalanced classes.
2. **Feature Engineering**:
   o Transforms raw data into meaningful features, improving model performance. Techniques like normalization and edge detection enhance the model's ability to capture patterns.
3. **Model Training and Evaluation**:
   o A baseline model is trained and evaluated to establish a performance benchmark. SHAP provides transparency, making the model's predictions interpretable.
4. **AutoML with H2O**:
   o Automates the selection of the best model and hyperparameters, reducing manual effort and ensuring optimal performance.
5. **Hyperparameter Tuning with Optuna**:
   o Fine-tunes the model to achieve the best possible accuracy, adhering to the principle of continuous improvement.
6. **Model Tracking and Monitoring with MLflow**:
   o Ensures reproducibility by logging all experiments, parameters, and metrics. Drift detection helps monitor model performance over time, ensuring that the model remains accurate as data evolves.

## Conclusion

The implemented pipeline adheres to MLOps best practices by ensuring data quality, automating model selection and hyperparameter tuning, providing model explainability, and tracking experiments for reproducibility. The use of tools like H2O AutoML, Optuna, and MLflow streamlines the workflow, making it efficient and scalable for real-world applications. This approach ensures that the model is not only accurate but also interpretable, maintainable, and adaptable to changing data conditions.