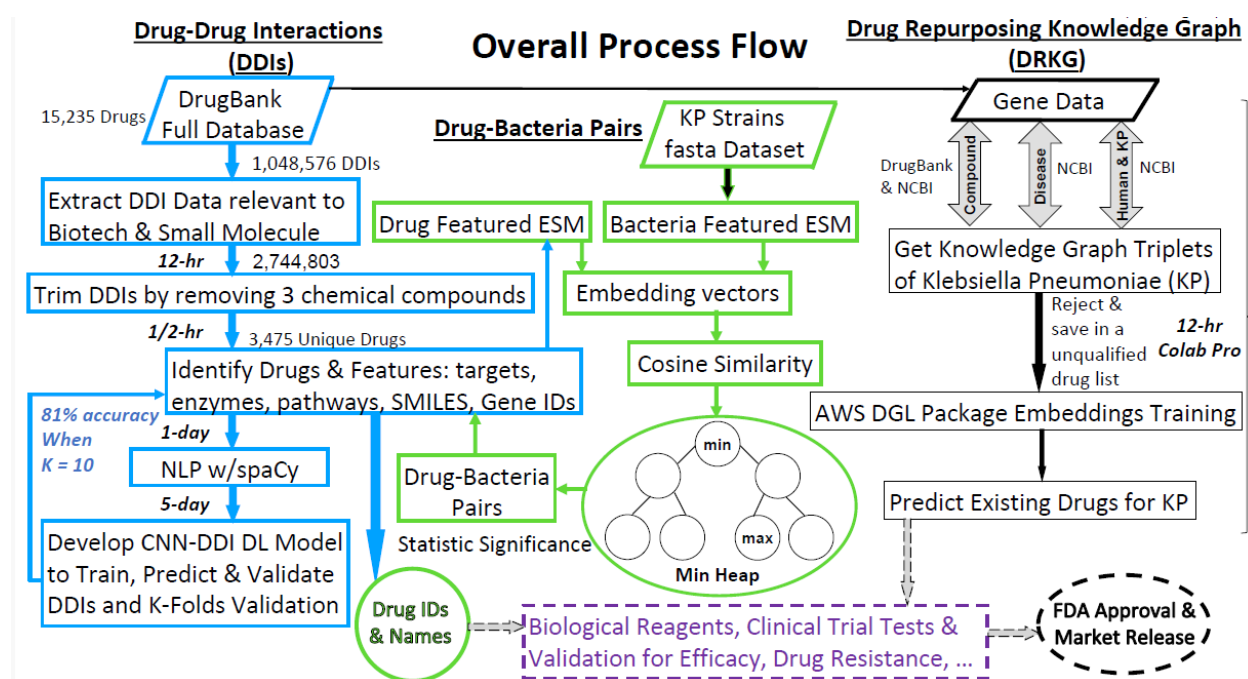# Overview:

This project employs deep learning techniques to accelerate drug discovery for combating *Klebsiella pneumoniae* (KP) infections. KP, a highly antibiotic-resistant pathogen, poses a significant global health challenge. Our approach integrates convolutional neural networks (CNN) and evolutionary scale modeling (ESM) to identify promising drug candidates from a dataset of biotech and small-molecule drugs extracted from DrugBank.

The figure below illustrates the overall process flow, including DDIs and unique drug extraction, NLP, CNN modeling, ESM with drug-bacteria similarity analysis, and the Drug Repurposing Knowledge Graph (DRKG), used in the deep learning framework to identify potential drug candidates against Klebsiella pneumoniae (KP).



# Dataset

1. DrugBank 5.1.11 version
   https://go.drugbank.com/releases/latest
2. Klebsiella Pneumonia (KP) strain fasta dataset
   https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9769640/pdf/spectrum.02306-22.pdf
3. National Center for Biotechnology Information (NCBI) databases

# Key Features

- Data-Driven Drug Discovery: Utilizes 3,475 biotech and small-molecule drugs from DrugBank, focusing on potential candidates for KP treatment.
- Deep Learning Models:
  - Convolutional Neural Networks (CNN): Applied for classification and interaction prediction.
  - Evolutionary Scale Modeling (ESM): Used for similarity analysis of drug molecules to KP strains.
- Drug-Drug Interaction Analysis: Incorporates drug-drug interaction (DDI) data for improved predictions.
- Similarity Metrics: Employs cosine similarity to identify candidates with >85% molecular similarity to KP strains.

## Results

1. Accuracy: Achieved approximately 72% validation accuracy for the CNN model.
2. Candidate Identification: Five promising drugs were identified with molecular similarities exceeding 85% to KP strains.
3. Efficiency: Reduced the drug discovery timeline from years to months.
4. Robustness: PairRE outperformed other models across evaluation metrics for predictive ranking.

## Significance

This project provides a scalable framework for drug discovery, offering a faster and more efficient alternative to traditional laboratory-based approaches. It holds potential for addressing bacterial pathogens like KP, aiding in effective treatment and minimizing side effects.

## Repository Contents

- **/data/**: Contains DrugBank-derived dataset used for training and evaluation.
- **/models/**: Includes implementations of CNN and ESM models.
- **/scripts/**: Python scripts for data preprocessing, training, and evaluation.
- **/results/**: Output files, including accuracy metrics and identified drug candidates.
- **README.md**: This file, describing the project.