



Background & Motivation

Clinical Context: Cytokine Release Syndrome (CRS) is a major safety concern for therapies like epcoritamab. Large pharmacovigilance databases (FAERS, EudraVigilance, JADER) capture real-world safety events at scale, but the data are heterogeneous, sparse, and noisy.

Problem: Traditional drug–AE analysis pipelines are slow, difficult to scale across oncology agents, and often generate outputs that are not easily interpretable by clinicians. As a result, rare but clinically important CRS signals may be overlooked.

Goal: Develop a generalized and explainable pharmacovigilance pipeline capable of analyzing any drug–AE pair, using epcoritamab and CRS as the motivating case study. Our system integrates anomaly detection, disproportionality statistics, and interpretable ML to produce clinician outputs.

Primary Databases: FAERS, EudraVigilance, JADER

Task 1: Causal Inference and NLP Methods for Multi-Database CRS Pharmacovigilance in Epcoritamab

Overview: In this task, we develop unified framework combining causal inference and NLP-driven feature extraction to characterize CRS risk across multiple pharmacovigilance databases. The goal is to quantify the effect of key treatments and clinical factors while leveraging unstructured narratives to improve severe CRS detection.

Methods:

1 Causal Inference

- DAG separating exposure (dose), confounders (age, disease stage, prior therapies), effect modifiers (steroids, tocilizumab), and colliders.
- Association tests: odds ratios, p-values for dose, steroids, tocilizumab, co-medications.
- Propensity score modeling to estimate causal effect of steroids on severe CRS.
- Sensitivity analysis using E-values for robustness to unmeasured confounding.

2 NLP & Narrative Features

- Rule-based extraction: ICU/hypotension, intubation, vasopressors.
- Identify steroid/tocilizumab administration, CRS grade, time-to-onset.
- Classifier predicts severe CRS using narrative-derived features.

3 Key Data Summary

- High-dose epcoritamab (≥ 24 mg) shows a strong association with severe CRS (OR = 7.18; robust E -value = 13.84).
- Steroid use is protective (OR = 2.84, $p = 0.017$); propensity analysis suggests a +10% reduction in severe CRS.
- Polypharmacy correlates with CRS severity (OR/SD = 1.78), likely reflecting disease complexity rather than causality.
- Significant cross-database heterogeneity ($p = 0.0002$) underscores the need for multi-source adjustment.

Task 2: Scalable Survival Analysis for Epcoritamab-Associated CRS

Objective: Develop a generalized, parameterized survival analysis pipeline that can evaluate *any* drug–adverse event combination. Using epcoritamab–CRS as the case study, we identify real-world CRS risk factors and characterize temporal onset patterns.

Cox Proportional Hazards Model Results

Variable	HR	95% CI	p-value
Weight (per kg)	0.992	0.985–1.000	0.037
Age (per year)	0.995	0.984–1.006	0.347
Polypharmacy (≥ 3)	0.616	0.153–2.482	0.495
Prior hospitalization	1.123	0.892–1.413	0.321

Table 1. Cox model showing weight as the only significant predictor of CRS (C-index = 0.58, * $p < 0.05$).

Weight-Based CRS Risk Stratification

Weight Category	Patients	CRS Rate	Clinical Action
<60 kg (low weight)	148	42.6%	High risk; inpatient dosing, enhanced monitoring
60–80 kg	404	30.7%	Moderate risk; standard 24h observation
>80 kg (high weight)	252	28.6%	Lower risk; may discharge after 24h if stable

Table 2. Low-weight patients (<60kg) have 1.49× higher CRS risk than >80kg.

Rare & Unexpected Signal Detection

We detect rare and unexpected drug–AE relationships using Isolation Forest anomaly detection combined with a 4-step filtering pipeline: (1) statistical anomaly identification, (2) FDA label filtering, (3) indication term removal, and (4) frequency-based filtering. Results are validated using disproportionality metrics (PRR>2, IC025>0, Chi-square>4). Our pipeline identified 1,386 rare signals across 37 oncology drugs. For each signal, BERT-based clinical feature analysis identifies demographic and medical history risk factors.

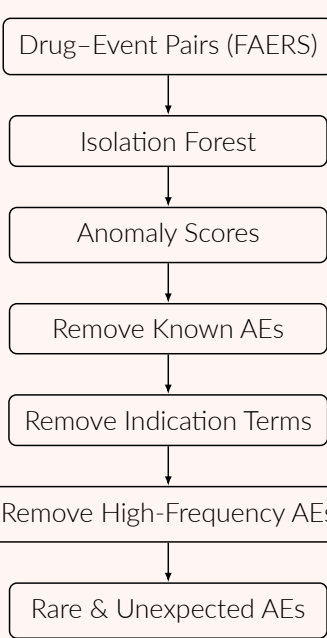


Figure 4. 4-Step Filtering Pipeline: From raw data to rare signals

Drug–Event: Epcoritamab + Haemorrhagic gastroenteritis

Status: RARE & UNEXPECTED

Report Count: 1

Statistical Metrics:
PRR (Proportional Reporting Ratio): 111.69 (threshold: >2)
IC025 (Information Component): 3.602 (threshold: >0)
Chi-square: 41.14 (threshold: >4)

Clinical Impact:
Death Rate: 100.0%
Hospitalization Rate: 100.0%
Serious Rate: 100.0%

Assessment:
In No-Cap Results: Yes
Passes PRR Test (>2): Yes
Passes IC025 Test (>0): Yes
Passes Chi² Test (>4): Yes
Passes All 3 Tests: Yes

CONCLUSION:
RARE & UNEXPECTED (passed IF + all 3 statistical tests)

Drug–Event: Epcoritamab + Renal impairment

Status: NOT RARE/UNEXPECTED

Report Count: 8

Statistical Metrics:
PRR (Proportional Reporting Ratio): 2.84 (threshold: >2)
IC025 (Information Component): 0.562 (threshold: >0)
Chi-square: 7.68 (threshold: >4)

Clinical Impact:
Death Rate: 137.5%
Hospitalization Rate: 175.0%
Serious Rate: 100.0%

Assessment:
In No-Cap Results: No
Passes PRR Test (>2): Yes
Passes IC025 Test (>0): Yes
Passes Chi² Test (>4): Yes
Passes All 3 Tests: Yes

CONCLUSION:
NOT RARE/UNEXPECTED (did not pass Isolation Forest)

(a) Rare & Unexpected (b) Not Rare/Unexpected

(c) Risk Factor Analysis

Figure 5. Three Example Outputs from Task 3 Analysis

CRS Mortality Modeling for Epcoritamab Using FAERS

A parameterized FAERS based pipeline was developed to identify CRS cases for Epcoritamab and predict CRS related mortality using machine learning. The workflow integrates engineered clinical features and produces interpretable risk factors through SHAP analysis and granular risk stratification. The pipeline is parameterized rather than tied to a single drug or event, so it can be quickly retargeted to other therapies and safety signals in FAERS, serving as a reusable template for broader drug safety monitoring.

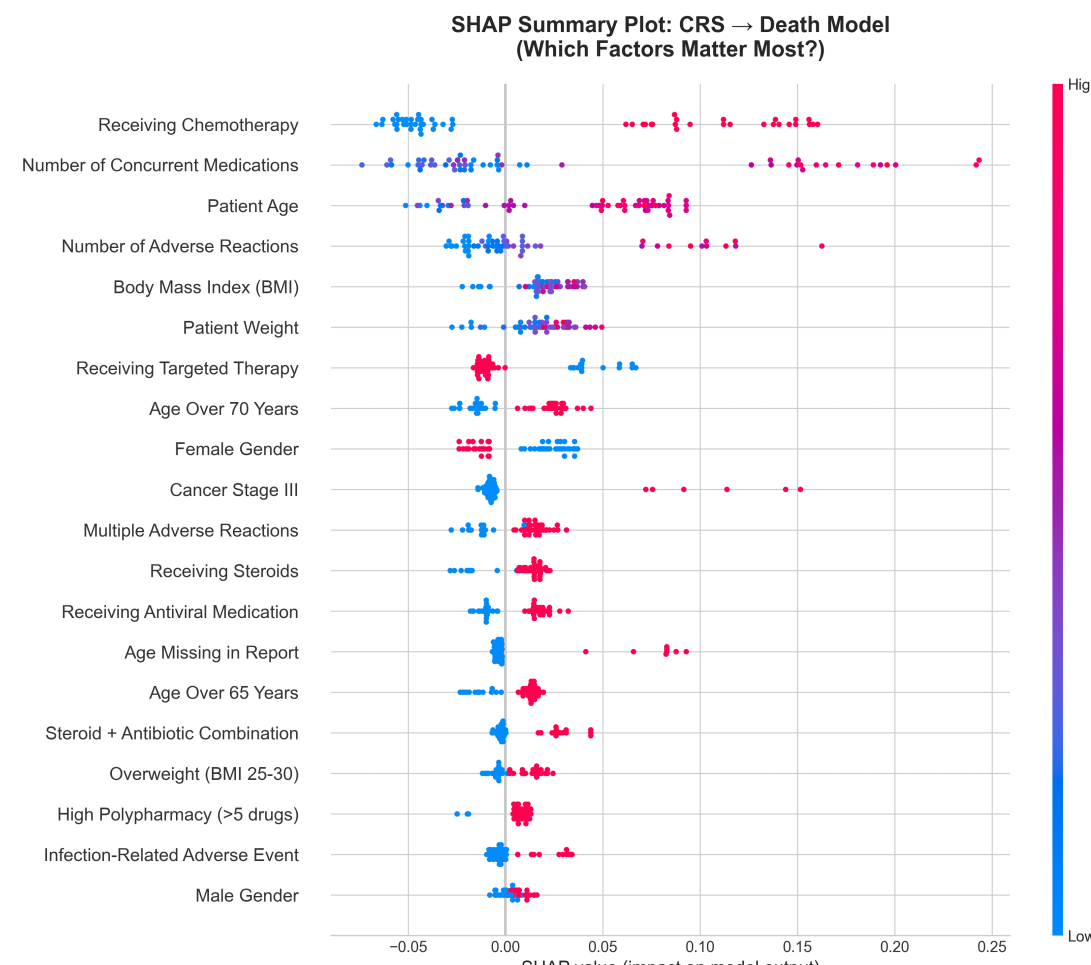
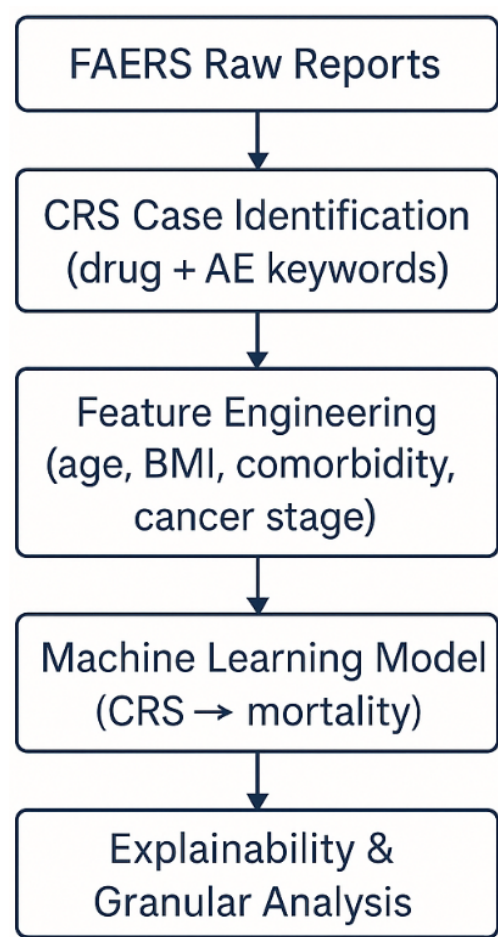


Figure X.

Explains which clinical features most strongly drive the CRS to mortality predictions.

Figure X. FAERS-based CRS mortality modeling pipeline.

A parameterized FAERS-based workflow identifies CRS cases for Epcoritamab, extracts cancer stage and engineered clinical features, and models CRS-related mortality using machine learning. The pipeline provides explainability through SHAP and granular clinical risk stratification.

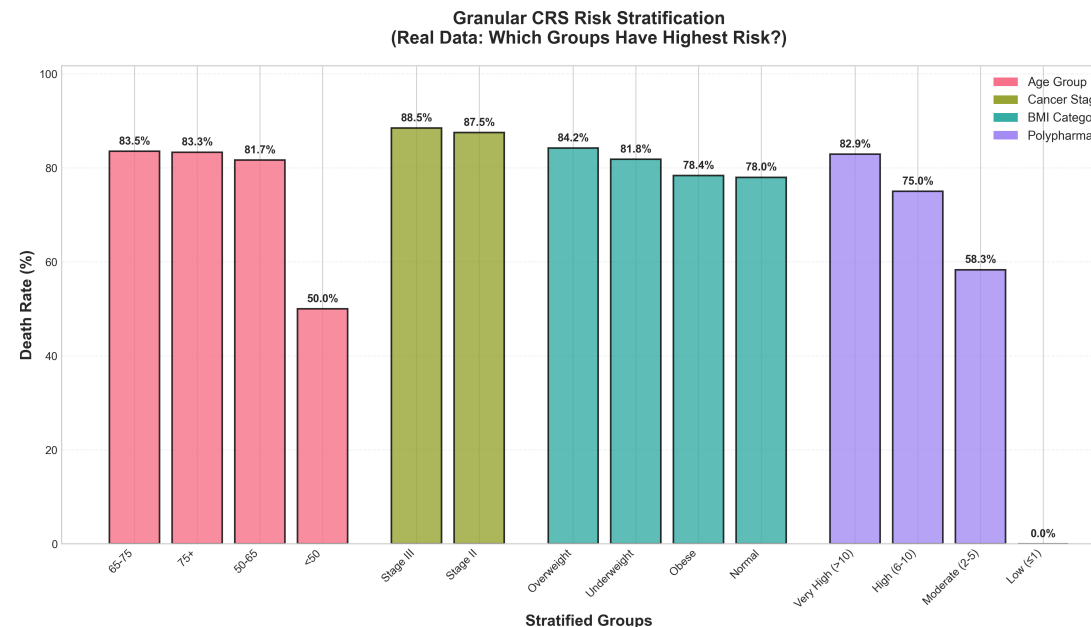


Figure X.

Shows real-world CRS mortality differences across age groups, cancer stage, BMI categories, and levels of polypharmacy.

Conclusion & Future Work

[CONCLUSION TEXT PLACEHOLDER - 100 words]

- [Future Work Item 1]
- [Future Work Item 2]
- [Future Work Item 3]



Figure 1. [Summary Caption]

Key Results

- [Key Result 1]
- [Key Result 2]
- [Key Result 3]
- [Key Result 4]