

## 1. Sumber Dataset Publik (Legal dan Bebas Lisensi)

Karena data medis bersifat sensitif, pengembangan sistem **AI Symptom Checker + Smart Triage** harus menggunakan dataset yang legal, terbuka, dan tidak mengandung identitas pasien.

Berikut beberapa sumber dataset publik yang bisa digunakan untuk melatih model klasifikasi gejala-penyakit:

| Dataset  | Deskripsi   | Format     | Link  |
|--|---|------------|---|
| <b>SymCAT<br/>(Symptom–<br/>Disease<br/>Associations )</b>       | Basis data berisi lebih dari 600 penyakit dan 1000 gejala beserta probabilitas keterkaitannya . Sangat berguna untuk membuat peta <i>symptom - to-disease mapping</i> . | JSON / CSV | <a href="https://www.symcat.com/diseases">https://www.symcat.com/diseases</a>   |
| <b>Human<br/>Symptoms<br/>Dataset<br/>(Kaggle)</b>               | Dataset 13.000 entri keluhan medis umum, cocok untuk melatih model klasifikasi teks dasar.  | CSV        | <a href="https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset">https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset</a> |
| <b>Medical<br/>Dialogue<br/>Dataset<br/>(MedDialog -<br/>EN)</b> | Dataset percakapan pasien–dokter dalam bahasa Inggris , bisa digunakan untuk ekstraksi keluhan awal pasien.   | JSON       | <a href="https://github.com/UCSD-AI4H/Medical-Dialogue-System">https://github.com/UCSD-AI4H/Medical-Dialogue-System</a>   |
| <b>ClinicalBERT<br/>Datasets<br/>(MIMIC-III<br/>Derived)</b>     | Dataset klinis hasil de-identifikasi untuk penelitian NLP medis. Dapat digunakan untuk fine-tuning model bahasa medis.  | SQL / CSV  | <a href="https://physionet.org/content/mimiciii-demo/">https://physionet.org/content/mimiciii-demo/</a>   |

 **Catatan Etika & Legalitas**

*Gunakan hanya data yang **tidak mengandung informasi pribadi (PII)**.*

*Jangan melakukan web scraping ke situs medis tanpa izin eksplisit.*

*Dataset ini cocok untuk tahap riset, prototipe, dan edukasi, **bukan** untuk keputusan medis klinis langsung.*