# DATASET

## Metadata

### Object

| doi | title | authors | publication_date | type | paper_id | concept |
|-----|-------|---------|------------------|------|----------|---------|

### Integer

| publication_year | cited_by_count |
|------------------|----------------|

## Paper Database — on .txt

4354 Paper

## Train & Test

|       | paper | referenced_paper |
|-------|-------|------------------|
| Train | 773   | 3834             |
| Test  | 773   | 3834             |

## Label On Train



### is_referenced

value_counts()

| 0 | 406399 |
|---|--------|
| 1 | 4292   |

imbalanced?

# PROCESSING TIMELINE

**Data Loading & Preprocessing**

**Document-Level Embedding & FE**

**Chunking & Chunk-Level Feature Extraction**

**Feature Assembly**

**Model Training & Tuning**

# Choosing Embedder

## Fine Tuning timeline of Specter



**Triplet-loss Fine-Tuning**

**In-Batch Negative Sampling**

**Sinyal Relasi Kutipan Eksplisit**

**Projection Head untuk Retrieval**

## Structure of BERT (12 layered)

**BERT** —— Pretrained on ——➤ Wikipedia, Books



**SciBERT** —— Trained from scratch from ——➤ Korpus ilmiah, tokenizer dengan **SciVocab**

**SPECTER** —— Fine tuned with ——➤ Citation aware objective

**Document-Level Embedding & FE**

# Embedding process

## For each paper on paper_db

Melewati tahapan →

### Inisialisasi

Specter

Hugging Face

### Batch Inference

batch_size = 16
max_length = 512

**Tokenize**

**Padding**

**Truncation**

### Ekstraksi Embedding

**CLS**

classification token

**768 dimensi**

Local MCC CV 0.372

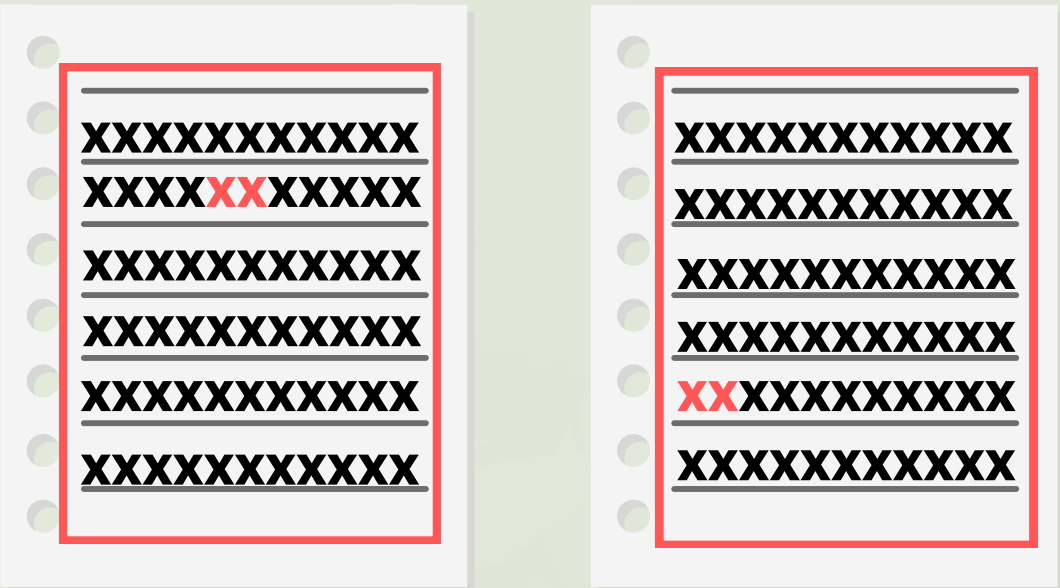**Document-Level Embedding & FE**

# Document level FE Creation

cosine similarity

year_diff same_year

can_cite

**Paper Database**


paper


paper ref

citations count

authors overlap

title Jaccard

# Pendekatan chunking kami

**x** = similaritas tinggi
**x** = similaritas rendah

## Document-Level Embedding

A          B



Cosine similarity ≈ 0.35

## Chunk Level Embedding

A          B



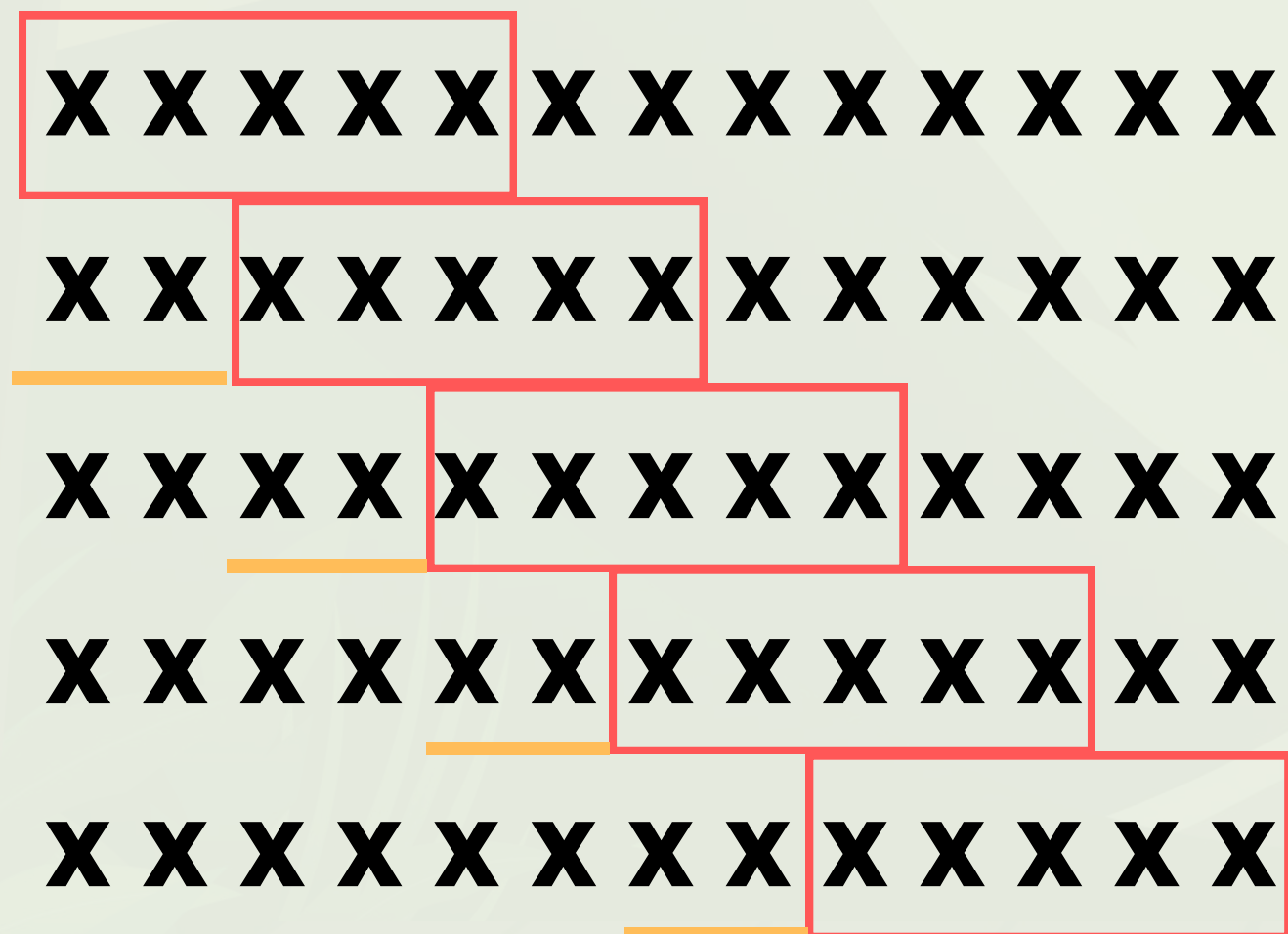Chunking with iterative approach on MiniLM-L6-v2

Chunk-level menemukan kesamaan lokal di bagian metodologi paper A dengan isi paper B.

Diekstrak menjadi

- Max similarity
- Mean Similarity
- Std Similarity
- Fraction Above 0.8

Chunk di line 2 (paper A dan line 5 (paper B) similarity ≈ 0.82

**Chunking & Chunk-Level Feature Extraction**
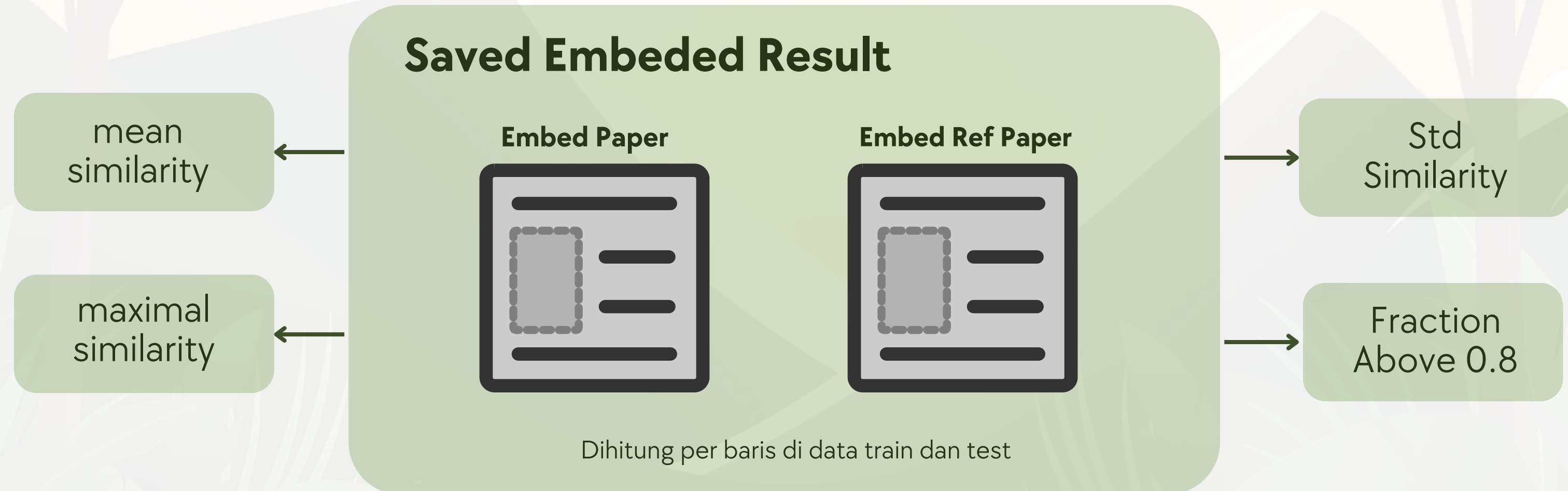
# Pendekatan chunking kami

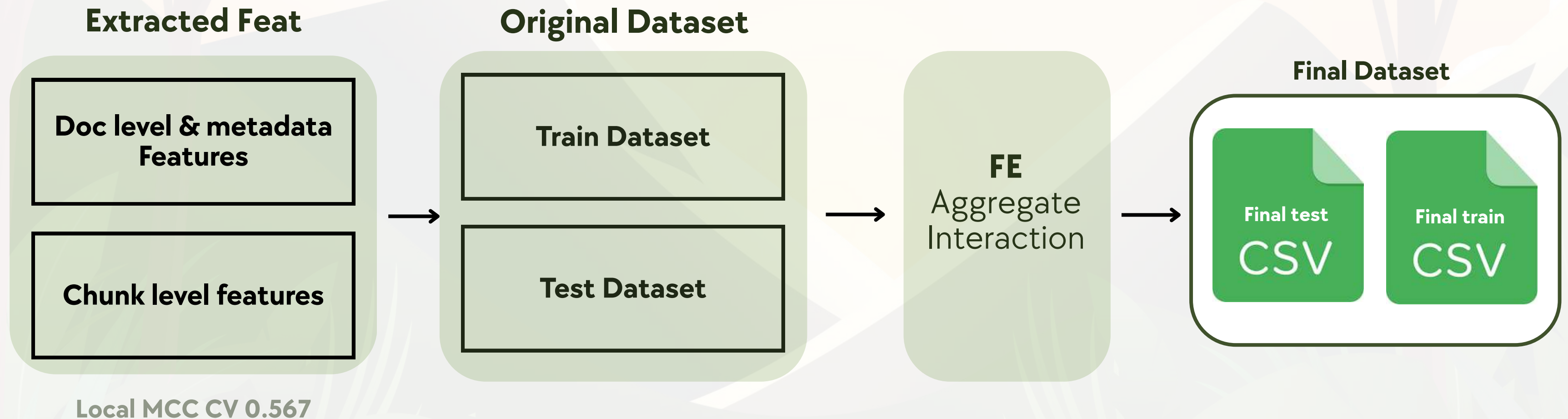$\boxed{\phantom{XXXX}}$ = chunk_size

———— = stride

Disimpan menjadi data embed untuk **feature extraction lanjutan di train dan test.**
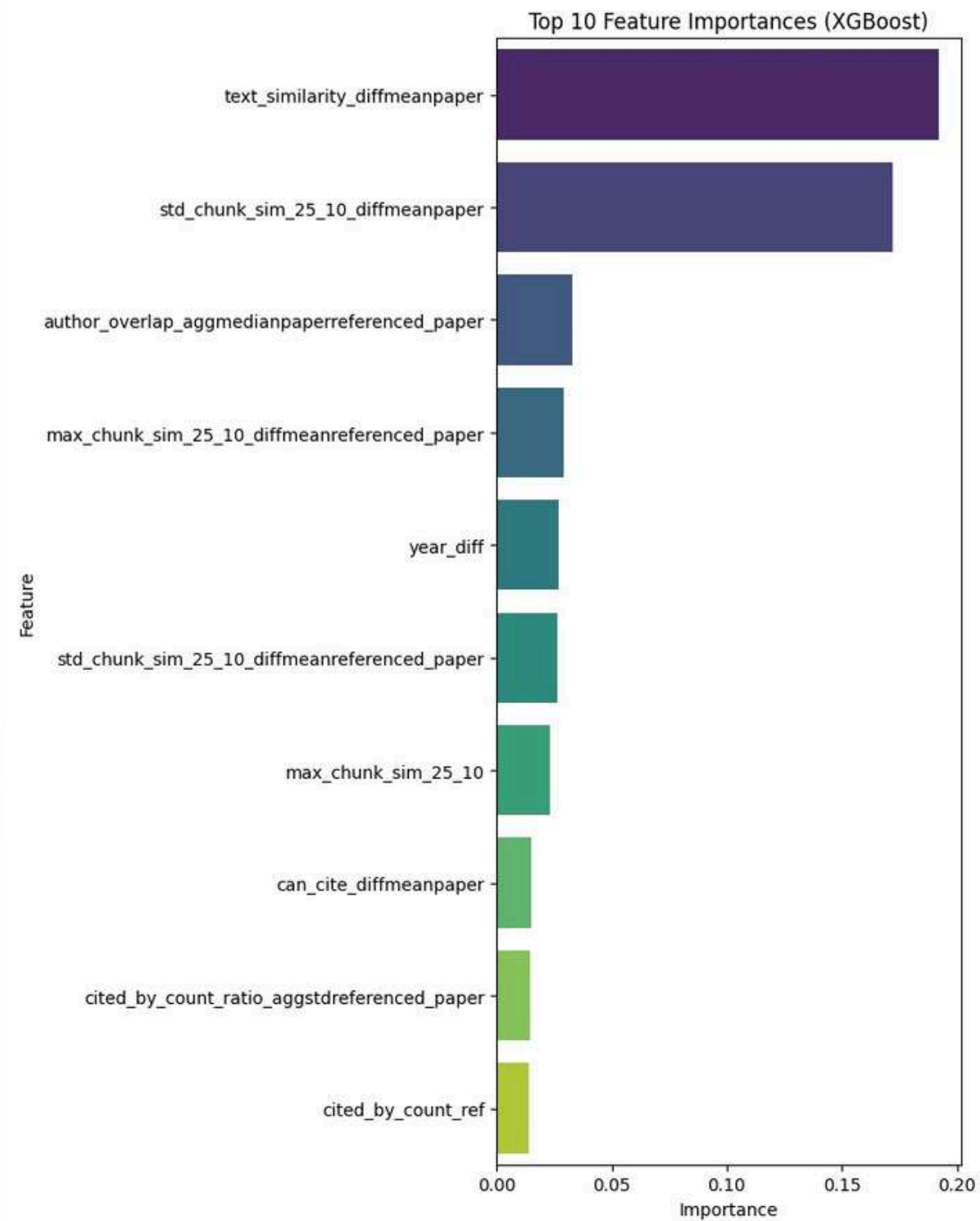
Local MCC CV 0.510

Chunking &
Chunk-Level
Feature Extraction

**Top 10 most important feat (XGboost)**

**Trained on final Dataset**

**Default model (untuned)**

Feature Extraction based on the **top 200 features** by **Feature Importance**
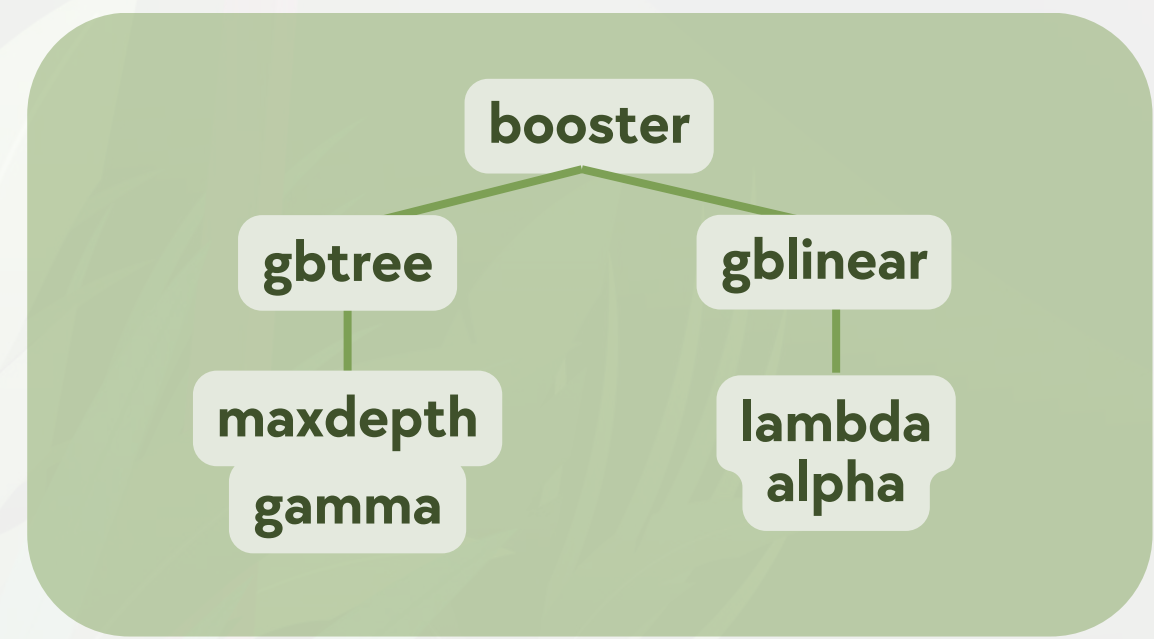
**Model Training & Tuning**

# Hyperparameter Tuning

## Tree-Structured Parzen Estimator (TPE)

$$\mathrm{EI}_{y^\star}[\boldsymbol{x}|\mathcal{D}] := \int_{-\infty}^{y^\star} (y^\star - y)p(y|\boldsymbol{x}, \mathcal{D})dy.$$

Function of expected improvement

booster
gbtree    gblinear
maxdepth    lambda
gamma       alpha

simple hyperparameter space example

Bayes-opt library: membangun dua model probabilistik
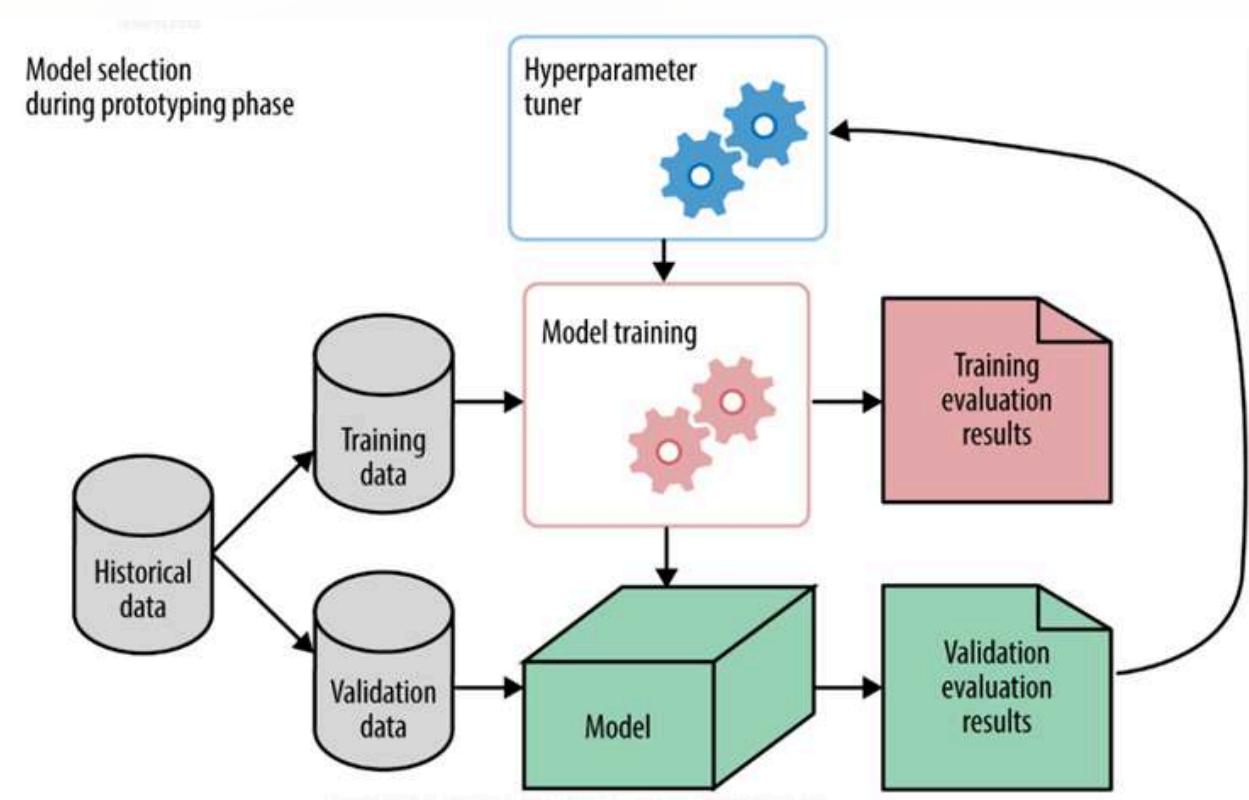
**l(x)**
"Better" high reward parameter

**g(x)**
"lesser" parameter

$$\mathbf{argmax}(\frac{l(x)}{g(x)})$$

**Model Training & Tuning**

# Hyperparameter Tuning

hyperparameter tuning process with optuna

## Tuned Hyperparameter

'lambda'                 'n_estimators'
'alpha'                  'max_depth'
'colsample_bytree'       'min_child_weight'
'subsample'              'gamma'
'learning_rate'

## Optimalisasi MCC

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$
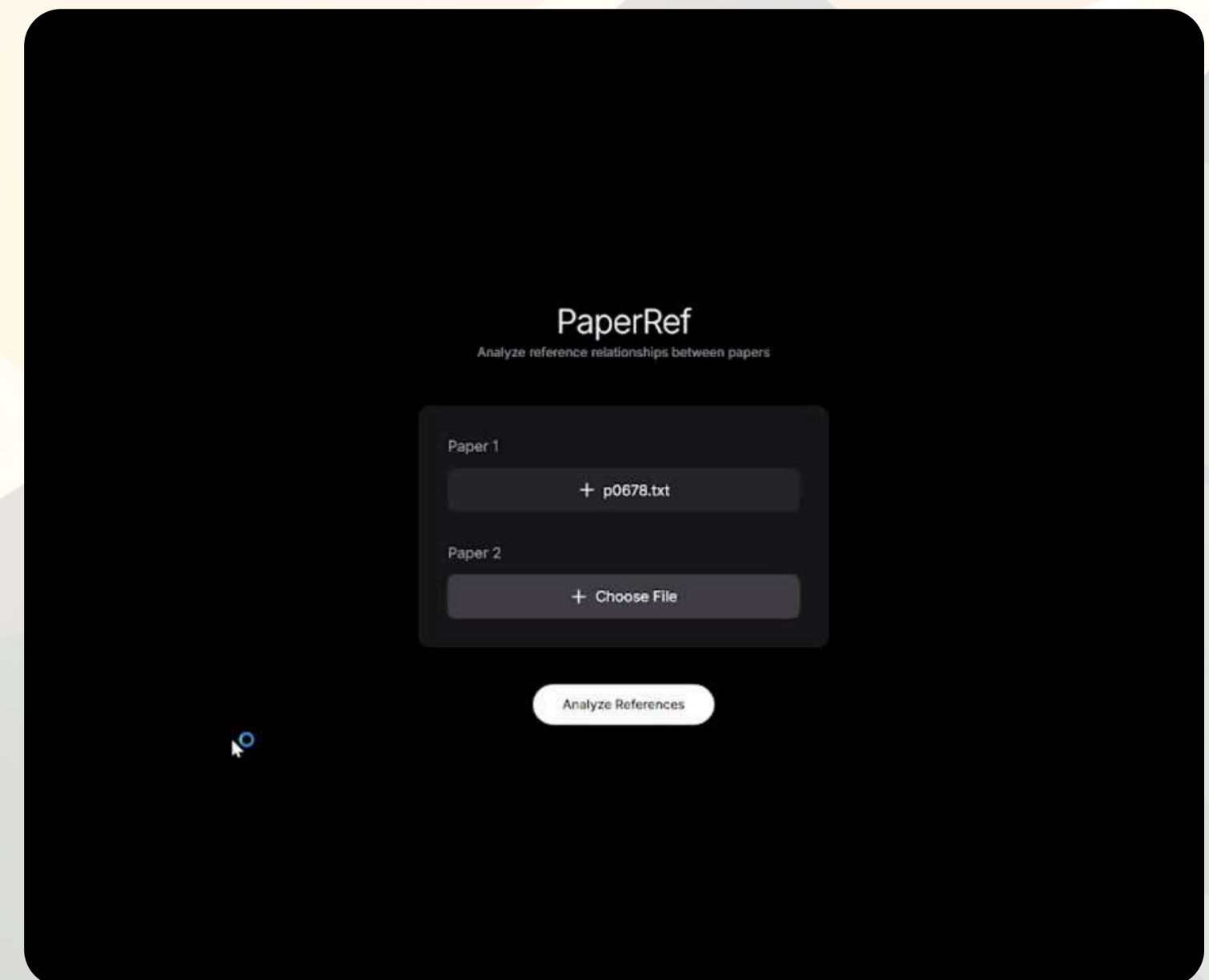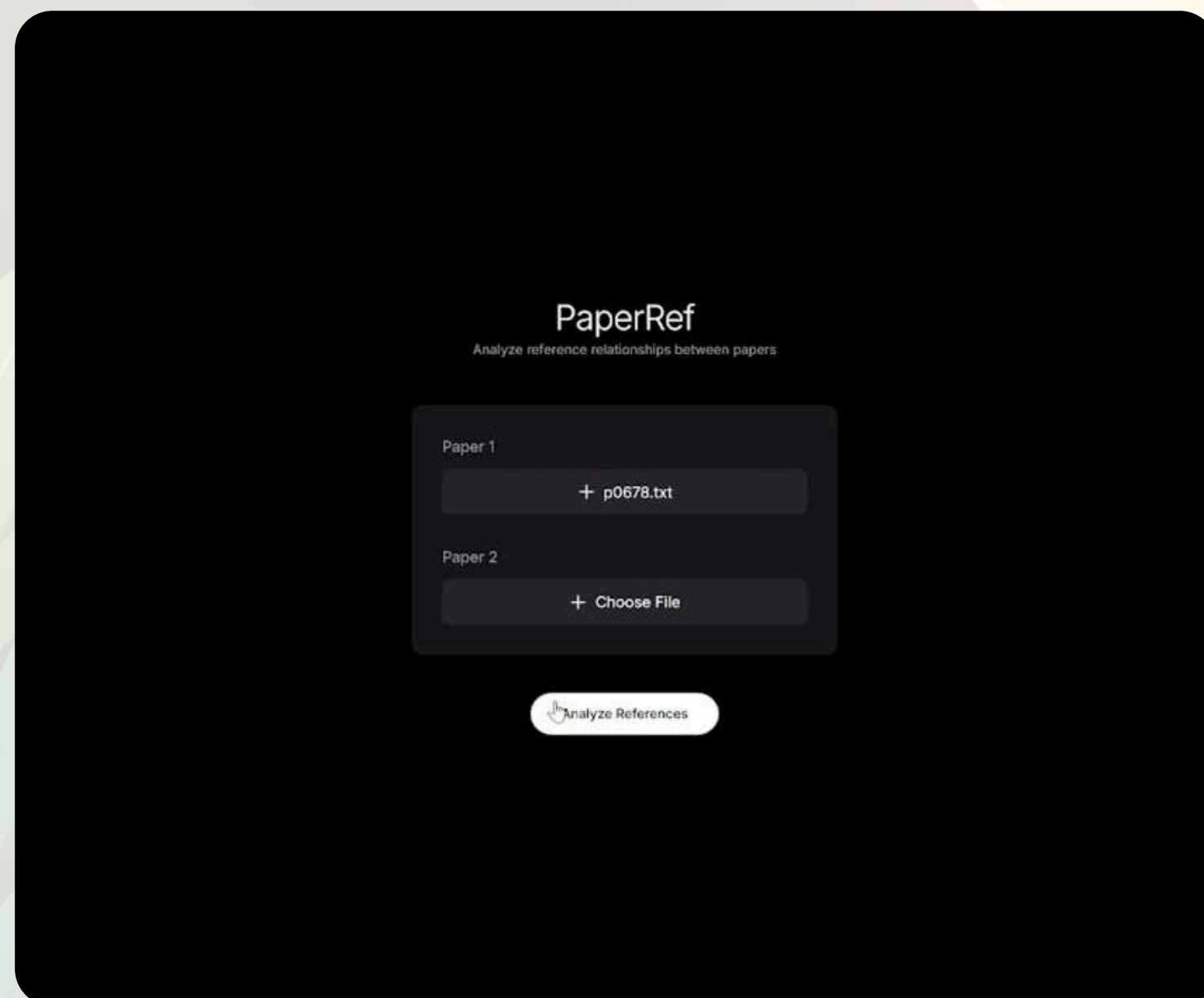
**Model Training & Tuning**

5

Run on local     Responsif ke pasangan baru     Modular Approach     Tidak perlu retrain boosted model

p0678 referensi ke p0508

p0678 tidak referensi ke p4101





https://github.com/Nadekoooo/refchecker.git

Sinergi antara pemahaman global, lokal, dan konteks bibliografis secara signifikan meningkatkan akurasi prediksi kutipan.
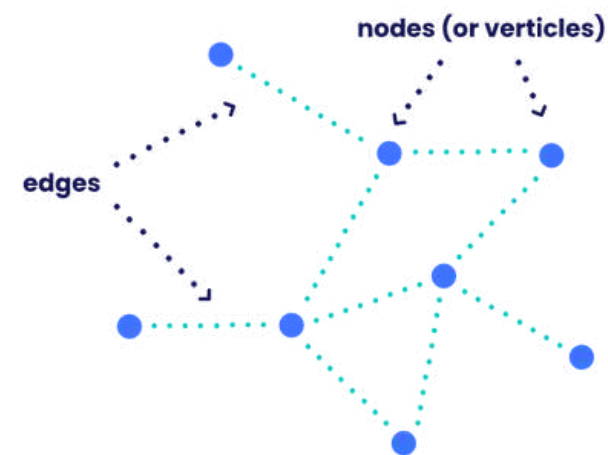
# REKOMENDASI

**Eksplorasi teknik ensembel lanjutan**



hillclimbing ensemble graph searching global max

**Integrasi fitur graf sitasi**



Ilustrasi GNN simpel

**Penggunaan model embedding yang lebih kuat**



scincl at Huggingface

## Kesimpulan dan Rekomendasi