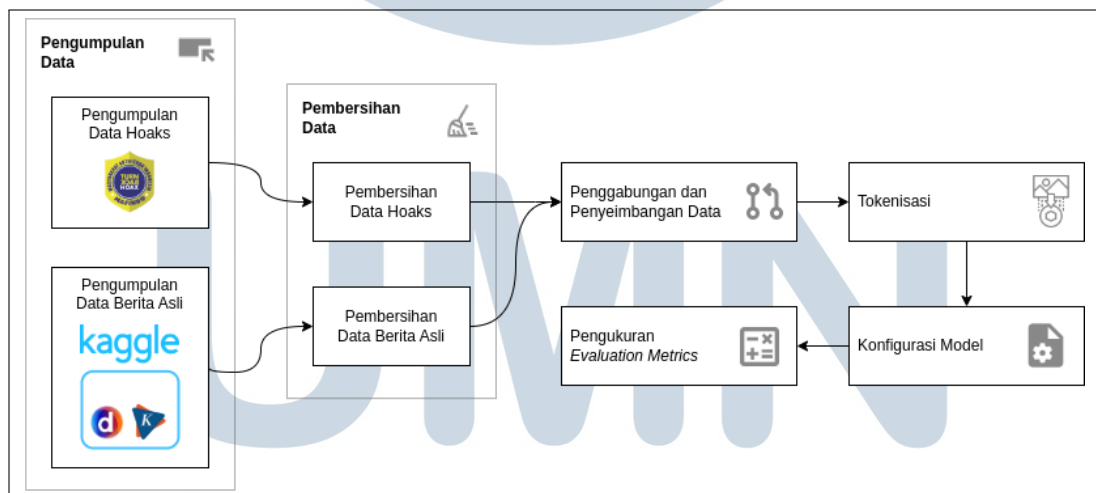


BAB 3 METODOLOGI

Dalam penelitian ini, pendekatan sistematis diterapkan untuk mendeteksi berita hoaks dalam bahasa Indonesia menggunakan model berbasis transformer. Metodologi ini mencakup alur kerja yang komprehensif mulai dari pengumpulan data hingga evaluasi model, yang menangani kompleksitas tugas klasifikasi teks dalam dunia nyata.

Langkah-langkah yang diuraikan dalam bagian ini menjelaskan bagaimana data dikumpulkan, dipra-proses, dan digunakan untuk melakukan fine-tuning model IndoBERT untuk tugas klasifikasi biner, seperti yang diilustrasikan dalam Gambar 3.1. Metodologi ini juga mencakup bagaimana dataset dibagi menjadi kategori berita hoaks dan berita asli, diproses melalui tahapan pembersihan individual, dan digabungkan untuk menciptakan dataset yang seimbang. Proses tokenisasi dan evaluasi model juga diuraikan secara detail untuk memberikan pemahaman lengkap tentang bagaimana model disiapkan untuk deteksi hoaks.



Gambar 3.1. *Flowchart* proses penelitian

3.1 Pengumpulan Data

Penelitian ini menggunakan strategi pengumpulan data multi-sumber untuk memastikan dataset yang kuat dan representatif untuk deteksi hoaks. Oleh karena itu, digunakan dua sumber dataset utama berikut.

3.1.1 Dataset Hoaks

Sumber utama data hoaks adalah API Mafindo, yang disediakan oleh organisasi anti-hoaks Indonesia Mafindo. Pendekatan pengambilan data secara programatik diimplementasikan menggunakan modul requests Python untuk mengakses endpoint API pada <https://yudistira.turnbackhoax.id/api/antihoax/> dimana implementasinya dapat dilihat pada Kode 3.1.

```
1 import requests
2 import csv
3
4 url = "https://yudistira.turnbackhoax.id/api/antihoax/"
5 batch_size = 1000
6 total_records = 22712
7
8 offset = 0
9
10 with open('api_response_data.csv', mode='w', newline='', encoding='
    'utf-8') as file:
11     writer = csv.writer(file)
12     headers_written = False
13     while offset < total_records:
14         data = {
15             'key': '...',
16             'limit': batch_size,
17             'offset': offset
18         }
19         response = requests.post(url, data=data)
20         if response.status_code == 200:
21             json_data = response.json()
22             if json_data and isinstance(json_data, list):
23                 if not headers_written:
24                     writer.writerow(json_data[0].keys())
25                     headers_written = True
26
27                 for entry in json_data:
28                     writer.writerow(entry.values())
29                 offset += batch_size
30
31                 print(f"Fetch batch starting at offset {offset}")
32
33     else:
```

```

33         print("No more data returned, stopping.")
34         break
35     else:
36         print(f"Failed to fetch data. Status code: {response.
status_code}, Response: {response.text}")
37         break

```

Kode 3.1: Hoax Data Fetching

Proses pengambilan data menghasilkan 14.552 entri, di mana 14.057 entri dilabeli sebagai hoaks dan 495 entri dilabeli sebagai berita asli. Ketidakseimbangan yang signifikan dalam dataset (dimana berita asli hanya sekitar 3% dari total data) mengharuskan pengumpulan sumber berita asli tambahan untuk menciptakan dataset yang lebih seimbang.

3.1.2 Dataset Berita Asli

Untuk melengkapi data berita asli yang terbatas, penelitian ini juga menggunakan "*Indonesia News Dataset (2024)*" dari Kaggle [22]. Dataset tambahan ini menyediakan koleksi artikel berita yang komprehensif dari media massa Indonesia yang faktual; yaitu Kompas.com, Tempo.co dan Detik.com. Dataset ini memiliki sekitar 45.234 artikel yang berkisar dari Januari 2024 hingga September 2024.

Kombinasi data hoaks API Mafindo dan dataset berita asli Kaggle membentuk korpus dasar untuk penelitian deteksi hoaks ini. Sumber yang beragam memastikan representasi yang komprehensif dari konten berita hoaks dan berita sah dalam bahasa Indonesia.

Tabel 3.1. Sumber dan Karakteristik Pengumpulan Data

Sumber Data	Total Entri	Entri Hoaks	Entri Berita Asli
API Mafindo	14.552	14.057	495
Dataset Kaggle	45.234	0	45.234
Total	59.786	14.057	45.729

Proses pengumpulan data dirancang untuk mengatasi potensi bias dan keterbatasan dalam dataset deteksi hoaks yang ada dengan memasukkan berbagai sumber dan memastikan representasi yang komprehensif dari konten berita hoaks dan berita sah.

3.2 Pembersihan Data

Proses pembersihan data merupakan langkah *preprocessing* yang kritis yang dirancang untuk meningkatkan kualitas dan keandalan dataset untuk deteksi hoaks. Proses ini melibatkan teknik normalisasi teks yang canggih yang diterapkan secara terpisah pada dataset hoaks dan berita asli untuk memastikan konsistensi dan menghilangkan noise yang berpotensi menyesatkan model pembelajaran mesin.

3.2.1 Dataset Hoaks

Untuk dataset hoaks yang dikumpulkan dari API Mafindo, proses pembersihan dimulai dengan penyaringan label yang cermat. Banyak entri mengandung metadata atau label status dari API Mafindo yang tidak secara langsung relevan dengan tugas klasifikasi hoaks. Secara khusus, label seperti "Klarifikasi", "Berita", "Dalam Proses", "Edukasi", dan kategorisasi serupa dihapus secara sistematis. Penyaringan selektif ini memastikan bahwa hanya entri hoaks atau berita asli yang jelas yang dipertahankan.

Prosedur pembersihan teks menerapkan beberapa transformasi kunci yang diimplementasikan dalam bentuk Python pada Kode 3.2, dan dijabarkan sebagai berikut.

- Penghapusan karakter khusus dan tanda baca yang dapat memperkenalkan noise ke dalam data tekstual
- Penghilangan URL dan tautan web yang tidak berkontribusi pada konten semantik
- Penghapusan spasi berlebih dan karakter baris baru

```
1 import pandas as pd
2 data = pd.read_csv('api_response_data.csv')
3
4 df_filtered = data[["status", "classification", "title", "content"
5 ]]
6 df_filtered = df_filtered[df_filtered["classification"] != "
7 CekFakta"]
8
9 # List of values to be removed
10 remove_values = ['Klarifikasi', 'Berita', 'Dalam Proses',
```

```

9         'Edukasi', 'Clarification', 'News', 'Education',
10        'Sebagian Benar']
11 # Filter DataFrame to exclude rows where 'status' is in the
12 # remove_values list
13 df_filtered = df_filtered[~df_filtered['status'].isin(
14     remove_values)]
15
16 df_filtered["status"].value_counts()
17
18 # Dictionary for mapping status values to 1 or 0
19 status_mapping = {
20     'Salah': 1,
21     'False': 1,
22     'Belum ada bukti': 1,
23     'Benar': 0,
24     'True': 0
25 }
26
27 # Apply the mapping to the 'status' column
28 df_filtered['status_mapped'] = df_filtered['status'].map(
29     status_mapping)
30
31 # Use str.replace() with a regular expression to remove the label
32 # at the beginning
33 df_filtered['title'] = df_filtered['title'].str.replace(r'
34     ^\[.*?\]\s*?:?^[A-Z]+\s*?:?^[A-Z]+\s*:\s*', '', regex=True)
35 df_filtered['title'] = df_filtered['title'].str.replace(
36     r'^(CEK FAKTA|Cek Fakta|CEK FAKTA Debat|\(CEK FAKTA Debat\))\s
37     *[:,\[\]\(\)]*', '', regex=True
38 )
39 # These two should replace unicode quotes into regular ASCII
40 # quotes, but the unicode quotes cannot be shown inside LaTeX
41 # df_filtered['title'] = df_filtered['title'].str.replace(r
42     '['\"']', '', regex=True)
43 # df_filtered['content'] = df_filtered['content'].str.replace(r
44     '['\"'\"'\"']', '', regex=True)
45
46 df_filtered['content'] = df_filtered['content'].str.replace(r'[\r\
47     n]+', ' ', regex=True)
48
49 df_filtered = df_filtered.drop(columns=['status', 'classification'
50 ])

```

```

40 df_filtered.rename(columns={'status_mapped': 'is_hoax'}, inplace=
    True)
41
42
43 # remove all missing values
44 df_filtered = df_filtered.dropna()
45 df_filtered
46 df_filtered["is_hoax"].value_counts()
47
48 # export df_filtered to csv
49 df_filtered.to_csv('filtered_data.csv', index=False)

```

Kode 3.2: Hoax Data Cleaning

3.2.2 Dataset Berita Asli

Dataset berita asli menjalani proses pembersihan yang serupa, namun sedikit dimodifikasi. Perhatian khusus diberikan untuk menghapus informasi pengantar situs web dan tag lokasi yang dapat memperkenalkan bias. Misalnya, awalan seperti "KOMPAS.com" atau indikator lokasi secara sistematis dihilangkan dari awal artikel. Pendekatan ini memastikan bahwa model akan fokus pada konten aktual daripada metadata periferal. Secara garis besar, dilakukan pembersihan teks sebagai berikut.

- Menghapus nama situs web dan informasi publikasi di awal
- Menghilangkan karakter khusus dan tanda baca yang tidak perlu
- Menormalisasi kasus teks dan menghapus spasi berlebih

Langkah-langkah tersebut kemudian diimplementasikan dalam Python dengan menggunakan regex. Adapun kode implementasinya adalah sebagai berikut.

```

1 def clean_website(text):
2     # Check if the input is a string, if not return it as-is (e.g
    ., handle NaN)
3     if isinstance(text, str):
4         # Check if the text starts with a news website name
    followed by additional information
5         if re.match(r'^(KOMPAS.com|TEMPO.CO|detik.com|CNN
    Indonesia|ANTARA News|Reuters|[A-Z\s,]+)', text):
6             # Remove everything up to and including the first dash
    character

```

```

7         text = re.sub(r'^[A-Z\s,]*(KOMPAS.com|TEMPO.CO|detik.
      com)[^ -]*[- ]\s*', '', text)
8         # Generalized version to remove text before the first
      dash if it's a location or publisher
9         text = re.sub(r'^[A-Z\s,.-]+[- ]\s*', '', text)
10        # Remove any leading/trailing whitespace
11        text = text.strip()
12    return text

```

Kode 3.3: Real News Data Cleaning

3.3 Penggabungan dan Penyeimbangan Data

Penggabungan dan penyeimbangan dataset dilakukan dengan pertimbangan cermat terhadap ketidakseimbangan kelas awal. Dataset Mafindo asli mengandung perbedaan yang signifikan, dengan entri hoaks (14.057) jauh melebihi entri berita asli (495). Untuk mengatasi ketidakseimbangan ini, pendekatan pengambilan sampel acak diterapkan. Dari Indonesia News Dataset yang lebih besar, yang berisi 45.234 artikel berita asli, 15.000 artikel di-*sampling* secara acak dan digabungkan dengan dataset hoaks Mafindo seperti yang tertera pada tabel 3.2. Hanya dua kolom yang akhirnya digunakan untuk pelatihan model, yaitu *content* (Teks lengkap artikel berita) dan *is_hoax* dengan label biner; 0 untuk berita asli, 1 untuk hoaks. Dataset ini kemudian dibagi menjadi subset *train*, *test*, dan *validation*. Pemisahan ini dilakukan untuk memastikan bahwa setiap subset mempertahankan representasi proporsional artikel hoaks dan berita asli. Dalam Python, pemisahan ini dilakukan dengan menggunakan fungsi *train_test_split* yang terdapat pada *library* *sklearn.model_selection*, seperti yang dijabarkan oleh Kode 3.4.

Tabel 3.2. Jumlah serta Rasio Pembagian Dataset yang digunakan

(a) Jumlah Data Berita		(b) Rasio Pembagian Data		
Berita Asli	Berita Hoaks	Training	Test	Validation
15.495	14.057	80%	10%	10%
Total: 29.552 artikel				

```

1 from sklearn.model_selection import train_test_split
2 # Split the data into train and test sets (80% train, 10%
      validation, 10% test)
3 train_texts, temp_texts, train_labels, temp_labels =
      train_test_split(

```

```

4     news_data['cleaned_content'], labels, test_size=0.2,
      random_state=42, stratify=labels
5 )
6
7 # Further split temp into validation and test sets (50% validation
  , 50% test from remaining 20%)
8 val_texts, test_texts, val_labels, test_labels = train_test_split(
9     temp_texts, temp_labels, test_size=0.5, random_state=42,
      stratify=temp_labels
10 )

```

Kode 3.4: Train, test, and validation dataset splitting process

3.4 Tokenisasi

Dalam penelitian ini, proses tokenisasi diimplementasikan menggunakan *library HuggingFace Transformers* (trainer), khususnya memanfaatkan tokenizer dari model `indobenchmark/indobert-base-p1`, yang dioptimalkan untuk bahasa Indonesia. Proses tokenisasi dari model ini dijabarkan secara algoritmik dengan menggunakan Python pada Kode 3.5.

```

1 from transformers import AutoTokenizer
2 tokenizer = AutoTokenizer.from_pretrained("indobenchmark/indobert-
  base-p1")
3
4 def tokenize_text(texts, tokenizer, max_len=512):
5     return tokenizer(
6         texts.tolist(),
7         truncation=True,
8         padding=True,
9         max_length=max_len,
10        return_tensors='np'
11    )
12
13 # Tokenize the train, validation, and test data
14 train_encodings = tokenize_text(train_texts, tokenizer)
15 val_encodings = tokenize_text(val_texts, tokenizer)
16 test_encodings = tokenize_text(test_texts, tokenizer)

```

Kode 3.5: Tokenizer Function

Di mana dengan menggunakan kode tersebut, proses tokenisasi pada dataset penelitian ini dapat dicontohkan sebagai berikut.


```
1 Teks Input: "Berita hoax merugikan masyarakat "  
2 Tokenisasi: [CLS] berita hoax merugikan masyarakat [SEP]  
3 Token: [101, 2034, 4231, 3211, 2098, 102]
```

3.5 Konfigurasi Model

Model deteksi hoaks memanfaatkan arsitektur transformer BERT. Model ini kemudian diadaptasi (*fine-tuning*) dalam penelitian ini untuk klasifikasi berita biner.

3.5.1 Arsitektur

Penelitian ini menggunakan model IndoBERT; secara spesifik varian `indobenchmark/indobert-base-p1`, yang menyediakan representasi linguistik bahasa Indonesia yang kuat. Konfigurasi model ini dirancang mengikuti nilai yang direkomendasikan oleh peneliti dari BERT untuk mengoptimalkan kinerja untuk klasifikasi biner, seperti yang diilustrasikan dalam Tabel 3.3.

Tabel 3.3. Parameter Konfigurasi Model

Parameter	Nilai
Base Model	IndoBERT
Number of Labels	2
Dropout Probability	0.3
Label Mapping	0: Legitimate News 1: Hoax News

3.5.2 Konfigurasi

Untuk membantu dalam proses konfigurasi *hyperparameter* pada penelitian ini, digunakan library `optuna` untuk melakukan variasi-variasi dalam *range* yang terdapat pada tabel 3.4. Nilai parameter yang digunakan merupakan nilai-nilai yang direkomendasikan oleh Devlin et al. dalam penelitiannya mengenai BERT [13].

Tabel 3.4. Hyperparameter Value

Hyperparameter	Nilai
Optimizer	AdamW
Learning Rate	5e-5, 3e-5 dan 2e-5
Weight Decay	1e-2
Ukuran Batch	32
Epoch	2 dan 3

Adapun dalam pengimplementasiannya, dapat diilustrasikan pada fungsi Python yang ada pada Kode 3.6.

```

1 def objective(trial):
2     learning_rate = trial.suggest_categorical('learning_rate', [5e
3     -5, 3e-5, 2e-5])
4
5     epoch = trial.suggest_categorical('epoch', [2, 3, 4])
6
7     (train_encodings, train_labels,
8     val_encodings, val_labels,
9     test_encodings, test_labels,
10    tokenizer) = prepare_data()
11
12    train_dataset = NewsDataset(train_encodings, train_labels)
13    val_dataset = NewsDataset(val_encodings, val_labels)
14
15    config = BertConfig.from_pretrained("indobenchmark/indobert-
16    base-p1")
17    config.num_labels = 2
18
19    model = BertForSequenceClassification.from_pretrained(
20        "indobenchmark/indobert-base-p1",
21        config=config
22    )
23
24    def compute_metrics(pred):
25        labels = pred.label_ids
26        preds = pred.predictions.argmax(-1)
27        precision, recall, f1, _ = precision_recall_fscore_support
28        (labels, preds, average='binary')
29        acc = accuracy_score(labels, preds)
30        return {
31            'accuracy': acc,
32            'f1': f1,

```

```

29         'precision': precision,
30         'recall': recall
31     }
32
33     training_args = TrainingArguments(
34         output_dir='./results',
35         num_train_epochs=epoch,
36         per_device_train_batch_size=32,
37         per_device_eval_batch_size=32,
38         learning_rate=learning_rate,
39         weight_decay=0.01,
40         eval_strategy="epoch",
41         save_strategy="no",
42         logging_dir='./logs',
43     )
44
45     trainer = Trainer(
46         model=model,
47         args=training_args,
48         train_dataset=train_dataset,
49         eval_dataset=val_dataset,
50         compute_metrics=compute_metrics,
51     )
52
53     trainer.train()
54
55     eval_results = trainer.evaluate()
56
57     return eval_results['eval_accuracy']

```

Kode 3.6: Hyperparameter Function

3.6 Pengukuran *Evaluation Metrics*

Pada penelitian ini, proses pengukuran *Evaluation Metrics* dibantu dengan menggunakan library `sklearn.metrics` dimana pengimplementasian library ini dapat dilihat pada Kode 3.7.

```

1 def compute_metrics(pred):
2     labels = pred.label_ids
3     preds = pred.predictions.argmax(-1)
4     precision, recall, f1, _ = precision_recall_fscore_support(
        labels, preds, average='binary')

```

```
5 acc = accuracy_score(labels, preds)
6 return {
7     'accuracy': acc,
8     'f1': f1,
9     'precision': precision,
10    'recall': recall
11 }
```

Kode 3.7: Model Evaluation Process

