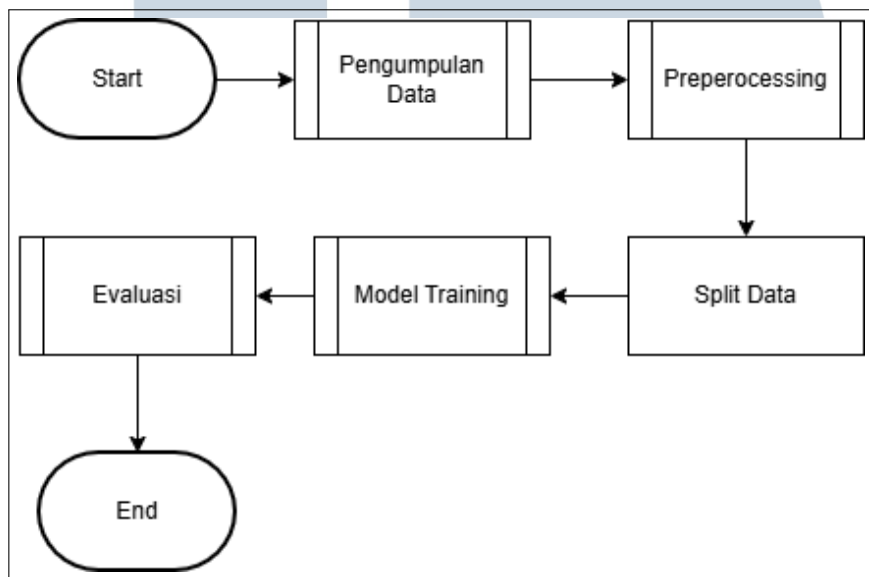


BAB 3 METODOLOGI PENELITIAN

3.1 Gambaran Umum Penelitian

Penelitian ini mengikuti prosedur yang terstruktur untuk memastikan bahwa proses berlangsung secara logis, seperti yang ditunjukkan pada Gambar 3.1.



Gambar 3.1. Diagram Alur Metodologi Penelitian

Tahap awal dalam penelitian ini dimulai dengan pengumpulan data terkait opini publik terhadap Program Makan Bergizi Gratis (MBG). Data yang diperoleh berasal dari media sosial X yang kemudian dilabeli menjadi 3 kelas yaitu sentimen negatif, sentimen positif, dan netral terhadap program MBG. Setelah data terkumpul, dilakukan proses pembentukan *dataset* terstruktur melalui tahapan *preprocessing*, yang mencakup pembersihan teks, normalisasi, dan penghapusan elemen yang tidak relevan agar siap digunakan dalam pemodelan.

Selanjutnya, data yang telah diproses digunakan untuk melatih model IndoBERTweet, yaitu model bahasa berbasis *Transformer* yang dioptimalkan untuk teks Twitter berbahasa Indonesia. *Dataset* dibagi menjadi tiga bagian: 70% untuk pelatihan, 15% untuk validasi, dan 15% untuk pengujian. Proses pelatihan bertujuan untuk menyesuaikan parameter model dengan karakteristik opini masyarakat terhadap MBG, sehingga model mampu mengenali pola-pola sentimen secara kontekstual dan akurat. Tahap akhir adalah pengujian model

terhadap data yang belum pernah dilihat sebelumnya, guna mengevaluasi performa klasifikasi sentimen. Evaluasi dilakukan menggunakan metrik seperti *F1-score* dan akurasi yang merefleksikan kemampuan model dalam mendeteksi opini publik secara tepat. Alur keseluruhan proses penelitian ini digambarkan secara visual pada Gambar 3.1.

3.2 Studi Literatur

Tahap ini mencakup kajian mendalam terhadap literatur ilmiah yang membahas model bahasa berbasis arsitektur *Transformer*, khususnya BERT dan pengembangannya dalam konteks bahasa Indonesia melalui IndoBERT dan IndoBERTweet. Tujuan dari studi ini adalah untuk memperoleh pemahaman komprehensif mengenai struktur model, mekanisme *pre-training* dan *fine-tuning*, serta efektivitasnya dalam tugas klasifikasi teks. Selain itu, ditinjau pula berbagai penelitian yang berfokus pada deteksi ujaran kebencian, terutama yang menggunakan data dari *platform* Twitter berbahasa Indonesia. Hasil kajian ini menjadi dasar konseptual dalam proses desain sistem, pemilihan model yang sesuai, dan implementasi algoritma klasifikasi berbasis IndoBERT dan IndoBERTweet dalam penelitian ini.

3.3 Arsitektur

Penelitian ini menggunakan model IndoBERTweet; secara spesifik varian *indolem/indobertweet-base-uncased*, yaitu model pralatih yang dirancang untuk teks Twitter berbahasa Indonesia dan mampu menghasilkan representasi linguistik yang kuat pada domain media sosial. Konfigurasi model ditetapkan dengan mengacu pada pengaturan yang direkomendasikan dalam pengembangan BERT untuk mengoptimalkan kinerja pada tugas klasifikasi, sebagaimana ditunjukkan pada Tabel 3.1 [5].

Selain parameter dasar tersebut, proses pelatihan dilakukan dengan *grid search* pada nilai *learning rate* 2e-5, 3e-5, dan 4e-5, yang mengacu pada pengembangan BERT, di mana nilai *learning rate* kecil pada orde 10^{-5} terbukti memberikan stabilitas dan akurasi tinggi pada berbagai tugas NLP. Dengan demikian, konfigurasi ini diharapkan mampu menghasilkan model yang optimal untuk klasifikasi sentimen publik terhadap Program MBG [20].

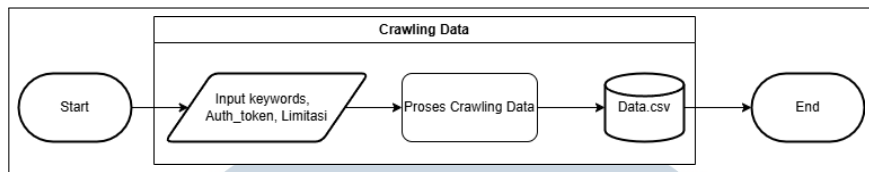
Tabel 3.1. Parameter Konfigurasi Model

Parameter	Nilai
<i>Base Model</i>	IndoBERTweet
<i>Number of Labels</i>	3
<i>Dropout Probability</i>	0.3
<i>Learning Rate</i>	2e-5, 3e-5, 4e-5
<i>Label Mapping</i>	0: Negatif 1: Netral 2: Positif

3.4 Pengumpulan Data

Data yang digunakan dalam penelitian ini diperoleh dari *platform X* (Twitter) menggunakan *Tweet Harvest*. Pengumpulan data dilakukan dengan ketentuan berikut:

- Sumber Data: *Tweet post* yang membahas Program Makan Gratis.
- Bahasa: Hanya mengambil *tweet* yang berbahasa Indonesia.
- Periode Data: Mulai dari 20 Oktober 2024 sampai dengan 20 Oktober 2025 dengan menggunakan *filter*.



Gambar 3.2. *Flowchart* Pengumpulan Data

Gambar 3.2 menampilkan *flowchart* untuk proses pengambilan data menggunakan metode *crawling*. Proses ini mencakup memasukkan *input* berupa *keywords*, *auth_token*, dan *limit*, kemudian melakukan *crawling* dengan Google Colab, dan menyimpan hasilnya dalam *file* *data.csv*.

3.4.1 Perancangan Model

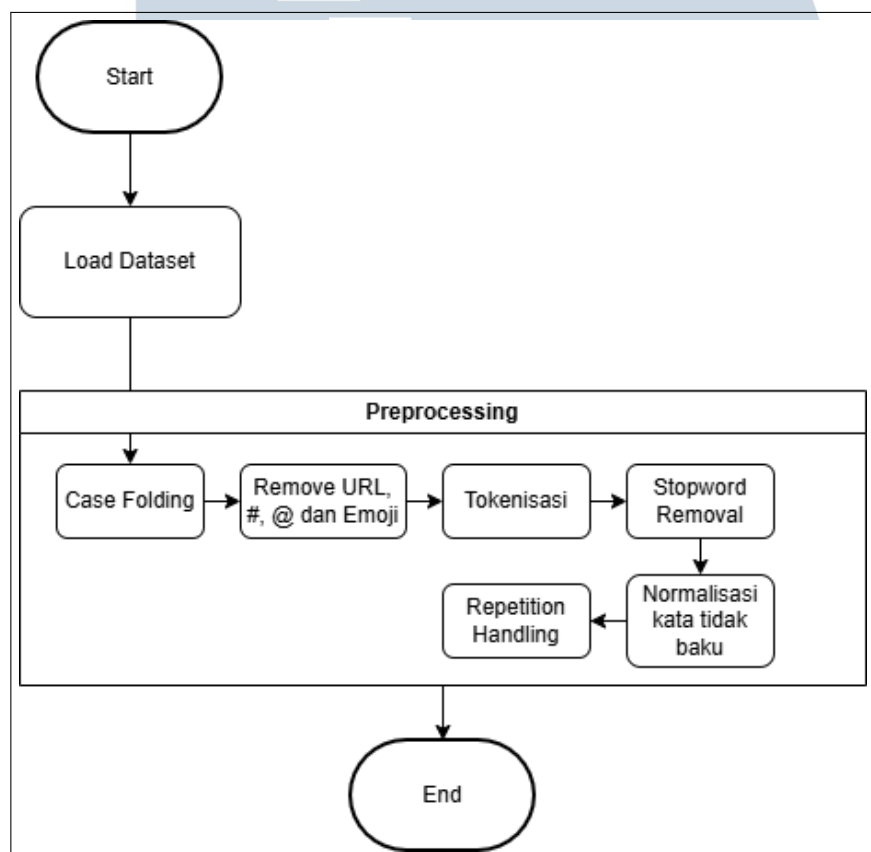
Data penelitian diperoleh dari unggahan publik di *platform* X (sebelumnya Twitter) dengan fokus pada opini masyarakat terhadap Program Makan Bergizi Gratis (MBG). Setelah proses pengumpulan berbasis kata kunci dan *filter* bahasa Indonesia, diperoleh 17.480 *tweet* yang memenuhi kriteria.

Untuk keperluan klasifikasi, setiap *tweet* dilabeli menjadi tiga kelas: Negatif, Netral, dan Positif. *Labeling* dilakukan melalui kombinasi anotasi manual (*seed*) dan strategi *pseudo-labeling* untuk memperbesar jumlah data berlabel. Distribusi kelas dievaluasi dan, bila perlu, dilakukan penyeimbangan (misalnya *undersampling/oversampling* atau penyesuaian bobot *loss*) agar model tidak bias terhadap kelas mayoritas. *Dataset* akhir dibagi menjadi *set* pelatihan, validasi, dan pengujian dengan rasio 70% / 15% / 15% untuk proses *fine-tuning* IndoBERTweet dan evaluasi performa.

UNIVERSITAS
MULTIMEDIA
NUSANTARA

3.4.2 Preprocessing

Sebelum data digunakan dalam pelatihan model, dilakukan proses *text preprocessing* untuk membersihkan dan menormalisasi teks komentar dari media sosial. Teks dari *platform* seperti X (atau Twitter) sering kali mengandung bahasa informal, simbol tidak relevan, dan struktur kalimat yang tidak baku. Oleh karena itu, *preprocessing* menjadi tahap krusial untuk meningkatkan kualitas *input* model.



Gambar 3.3. Flowchart Text Preprocessing

Gambar 3.3 memperlihatkan alur *preprocessing* yang terdiri dari tahapan berikut:

1. Case Folding

Gambar 3.4 menunjukkan proses pengubahan semua huruf dalam teks menjadi huruf kecil (*lowercase*) untuk menyamakan format kata dan menghindari redundansi makna. Contohnya, kata "Perempuan" dan "perempuan" akan dianggap sama oleh model.

```
1. BEFORE: @prabowo Program MBG keren bangettt! 😊 https://mbg.go.id #MBG2024
   AFTER:  @prabowo program mbg keren bangettt! 😊 https://mbg.go.id #mbg2024

2. BEFORE: Kasihan anak2 Papua keracunan makanan 🤢 @Kemendikbud tolong dong!
   AFTER:  kasihan anak2 papua keracunan makanan 🤢 @kemendikbud tolong dong!

3. BEFORE: Pemerintah alokasikan Rp 71T untuk program ini www.detik.com/news
   AFTER:  pemerintah alokasikan rp 71t untuk program ini www.detik.com/news

4. BEFORE: Gak bagussss program ini, ngabisin duit rakyat ajaaa 🙄
   AFTER:  gak bagussss program ini, ngabisin duit rakyat ajaaa 🙄

5. BEFORE: Mantapppp! Programnya okeh bgt, semoga merata yaa 👍
   AFTER:  mantapppp! programnya okeh bgt, semoga merata yaa 👍
```

Gambar 3.4. Contoh kata dari *case folding*

2. Remove URL, #, Mention, dan Emoji

Gambar 3.5 menunjukkan penghapusan elemen-elemen yang tidak memiliki nilai semantik seperti tautan (URL), *mention* (@username), *hashtag* (#topik), dan emoji. Langkah ini bertujuan untuk membersihkan teks dari simbol yang tidak relevan.

```
1. BEFORE: @prabowo program mbg keren bangettt! 😊 #mbg2024
   AFTER:  program mbg keren bangettt  mbg

2. BEFORE: kasihan anak2 papua keracunan makanan 🤢 @kemendikbud tolong dong!
   AFTER:  kasihan anak papua keracunan makanan  tolong dong

3. BEFORE: pemerintah alokasikan rp 71t untuk program ini
   AFTER:  pemerintah alokasikan rp  t untuk program ini

4. BEFORE: gak bagussss program ini, ngabisin duit rakyat ajaaa 🙄
   AFTER:  gak bagussss program ini  ngabisin duit rakyat ajaaa

5. BEFORE: mantapppp! programnya okeh bgt, semoga merata yaa 👍
   AFTER:  mantapppp  programnya okeh bgt  semoga merata yaa
```

Gambar 3.5. Contoh untuk menghapus URL, #, *mention*, dan emoji

3. Tokenisasi

Gambar 3.6 menjelaskan pemecahan teks menjadi unit kata (*token*) menggunakan *tokenizer* yang sesuai dengan struktur Bahasa Indonesia. Tokenisasi membantu model memahami struktur kalimat secara granular.

```
1. BEFORE: program mbg keren banget mbg
   AFTER: ['program', 'mbg', 'keren', 'banett', 'mbg']
   Total tokens: 5

2. BEFORE: kasihan anak papua keracunan makanan tolong dong
   AFTER: ['kasihan', 'anak', 'papua', 'keracunan', 'makanan', 'tolong', 'dong']
   Total tokens: 7

3. BEFORE: pemerintah alokasikan rp t untuk program ini
   AFTER: ['pemerintah', 'alokasikan', 'rp', 't', 'untuk', 'program', 'ini']
   Total tokens: 7

4. BEFORE: gak baguss program ini ngabisin duit rakyat ajaa
   AFTER: ['gak', 'baguss', 'program', 'ini', 'ngabisin', 'duit', 'rakyat', 'ajaa']
   Total tokens: 8

5. BEFORE: mantapp programnya okeh bgt semoga merata yaa
   AFTER: ['mantapp', 'programnya', 'okeh', 'bgt', 'semoga', 'merata', 'yaa']
   Total tokens: 7
```

Gambar 3.6. Contoh tokenisasi

4. *Stopword Removal*

Gambar 3.7 menunjukkan penghapusan kata-kata umum yang tidak berkontribusi pada analisis sentimen, seperti "yang", "dan", "itu", menggunakan daftar *stopword* Bahasa Indonesia.

```
1. BEFORE: ['program', 'mbg', 'bagus', 'banett', 'mbg']
   AFTER: ['program', 'mbg', 'bagus', 'banett', 'mbg']
   Tokens: 5 → 5

2. BEFORE: ['kasihan', 'anak', 'papua', 'keracunan', 'makanan', 'tolong']
   AFTER: ['kasihan', 'anak', 'papua', 'keracunan', 'makanan', 'tolong']
   Tokens: 6 → 6

3. BEFORE: ['pemerintah', 'alokasikan', 'rp', 't', 'untuk', 'program', 'ini']
   AFTER: ['pemerintah', 'alokasikan', 'program']
   Removed: ['rp', 't', 'untuk', 'ini']
   Tokens: 7 → 3

4. BEFORE: ['tidak', 'baguss', 'program', 'ini', 'ngabisin', 'duit', 'rakyat', 'ajaa']
   AFTER: ['baguss', 'program', 'ngabisin', 'duit', 'rakyat', 'ajaa']
   Removed: ['tidak', 'ini']
   Tokens: 8 → 6

5. BEFORE: ['mantapp', 'programnya', 'baik', 'sangat', 'semoga', 'merata', 'yaa']
   AFTER: ['mantapp', 'programnya', 'semoga', 'merata', 'yaa']
   Removed: ['baik', 'sangat']
```

Gambar 3.7. Contoh *stopword*

5. Normalisasi Kata Tidak Baku

Gambar 3.8 menunjukkan perubahan kata-kata *slang* atau tidak baku seperti "gk", "bgt", "tdk" menjadi bentuk formal seperti "tidak" dan "banget" menggunakan kamus normalisasi.

```
1. BEFORE: ['program', 'mbg', 'keren', 'bangett', 'mbg']
   AFTER:  ['program', 'mbg', 'bagus', 'bangett', 'mbg']
   Changed: keren→bagus

2. BEFORE: ['kasihan', 'anak', 'papua', 'keracunan', 'makanan', 'tolong', 'dong']
   AFTER:  ['kasihan', 'anak', 'papua', 'keracunan', 'makanan', 'tolong']

3. BEFORE: ['pemerintah', 'alokasikan', 'rp', 't', 'untuk', 'program', 'ini']
   AFTER:  ['pemerintah', 'alokasikan', 'rp', 't', 'untuk', 'program', 'ini']

4. BEFORE: ['gak', 'baguss', 'program', 'ini', 'ngabisin', 'duit', 'rakyat', 'ajaa']
   AFTER:  ['tidak', 'baguss', 'program', 'ini', 'ngabisin', 'duit', 'rakyat', 'ajaa']
   Changed: gak→tidak

5. BEFORE: ['mantapp', 'programnya', 'okeh', 'bgt', 'semoga', 'merata', 'yaa']
   AFTER:  ['mantapp', 'programnya', 'baik', 'sangat', 'semoga', 'merata', 'yaa']
   Changed: okeh→baik, bgt→sangat
```

Gambar 3.8. Contoh menormalisasikan kata baku

6. Menghapus Kata yang Berulang

Gambar 3.9 menunjukkan penghapusan kata-kata yang berulang seperti "bagusss", "bangett", "burukkk", "hebataat" menjadi bentuk normal seperti "bagus", "banget", "buruk", "hebat".

```
1. SEBELUM : program ini bagusss sekali
   SESUDAH : program ini bagus sekali

2. SEBELUM : mantap bangett
   SESUDAH : mantap banget

3. SEBELUM : burukkk banget
   SESUDAH : buruk banget

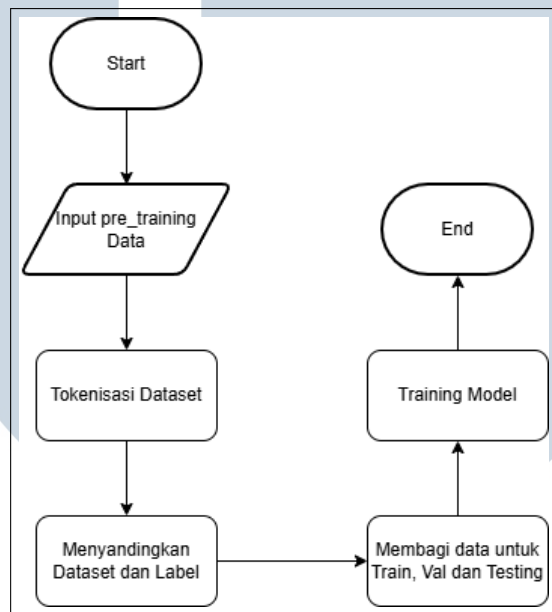
4. SEBELUM : hebataat luar biasaaa
   SESUDAH : hebat luar biasa

5. SEBELUM : tidak setujuuu
   SESUDAH : tidak setuju
```

Gambar 3.9. Contoh kata berulang

3.4.3 Model Training

Gambar 3.10 menunjukkan alur pelatihan model IndoBERTweet untuk klasifikasi ujaran kebencian. *Dataset* dibagi menjadi tiga bagian: 70% untuk pelatihan, 15% validasi, dan 15% pengujian. Setelah itu, ditentukan *hyperparameter* seperti jumlah *epoch*, ukuran *batch*, dan *learning rate*.



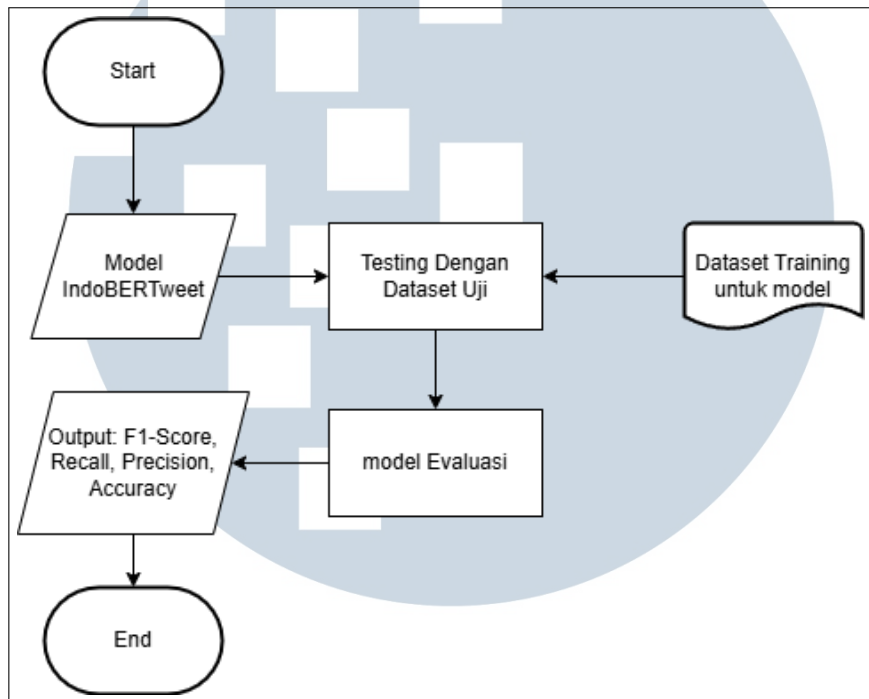
Gambar 3.10. Flowchart Pelatihan Model

Pemilihan model terbaik dilakukan berdasarkan nilai *F1-score* tertinggi pada data validasi. Hal ini dikarenakan *F1-score* merupakan metrik yang menggabungkan *precision* dan *recall* dalam satu ukuran, sehingga lebih representatif dibandingkan akurasi terutama pada data yang tidak seimbang. Dengan begitu, model yang dipilih tidak hanya memiliki kemampuan prediksi yang tepat, tetapi juga bisa mendeteksi seluruh kelas secara seimbang.

UNIVERSITAS
MULTIMEDIA
NUSANTARA

3.5 Evaluasi

Gambar 3.11 menunjukkan alur proses pengujian model IndoBERTweet. Model terbaik yang telah diperoleh dari tahap pelatihan kemudian diuji menggunakan *dataset* yang sudah disiapkan.



Gambar 3.11. *Flowchart Testing*

Tujuan dari pengujian ini adalah untuk mengevaluasi kemampuan generalisasi model terhadap data baru. Setelah proses *testing* selesai, performa model diukur menggunakan metrik evaluasi seperti *F1-score*, *recall*, *precision*, dan akurasi guna mengetahui seberapa efektif IndoBERTweet dalam mengidentifikasi ujaran kebencian pada data aktual.

UNIVERSITAS
MULTIMEDIA
NUSANTARA