

## DAFTAR LAMPIRAN

2021230044_Rezza Maulana			
ORIGINALITY REPORT			
19%	18%	7%	%
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS
PRIMARY SOURCES			
1	journal.uinjkt.ac.id Internet Source		1%
2	jurnal.penerbitdaarulhuda.my.id Internet Source		1%
3	ejurnal.esaunggul.ac.id Internet Source		1%
4	eprints.upj.ac.id Internet Source		1%
5	amikmahaputra.ac.id Internet Source		1%
6	repository.uin-suska.ac.id Internet Source		1%
7	text-id.123dok.com Internet Source		1%
8	jurnal.unived.ac.id Internet Source		1%
9	repo.darmajaya.ac.id Internet Source		1%
10	tunasbangsa.ac.id Internet Source		<1%
11	Muhammad Yusuf, Indah Purnama Sari. "Sistem Pakar Mencegah Stunting dengan		<1%



**UNIVERSITAS DARMA PERSADA**  
**UPT PERPUSTAKAAN**  
Gedung Rektorat Lantai 3,  
Jl. Taman Malaka Selatan, Pondok Kelapa – Jakarta Timur 13450

**SURAT KETERANGAN**  
**HASIL PENGECEKAN TURNITIN**

UPT Perpustakaan Universitas Darma Persada menerangkan telah selesai melakukan pemeriksaan duplikasi/*similarity* menggunakan perangkat lunak Turnitin terhadap hasil karya sebagai berikut:

Judul : Penerapan Natural Language Processing pada Pengelolaan Berkas Digital menggunakan Chatbot di Pusdatin Dinas PPKUKM

Penulis : Rezza Maulana

NIM : 2021230044

Tgl pemeriksaan : 25 Juli 2025

Dengan hasil Tingkat Kesamaan (*similarity index*) **19%**

Demikian Surat Keterangan kami buat, untuk dipergunakan sebagaimana mestinya.

Jakarta, 25 Juli 2025

Ka.UPT Perpustakaan Unsada

Yus Rusmiyati, SS., MM

Batas maksimal similarity 30% untuk Fakultas Sastra dan Ekonomi

Batas maksimal similarity 25% untuk Fakultas Teknik, Kelautan dan Pasca Sarjana

## Source code

### Nlp\_api.py

```
1. from fastapi import FastAPI, Query
2. from pydantic import BaseModel
3. import re
4. import os
5. import fitz # PyMuPDF
6. import torch
7. from transformers import AutoTokenizer,
   AutoModelForSequenceClassification, pipeline
8.
9. app = FastAPI()
10.
11. # Load intent classification model (BERT)
12. intent_model_path = "bert-intent-model" # Folder hasil
   training
13. intent_model =
   AutoModelForSequenceClassification.from_pretrained(
14.     intent_model_path,
15.     use_safetensors=True
16. )
17. intent_tokenizer =
   AutoTokenizer.from_pretrained(intent_model_path)
18. label_map = intent_model.config.id2label
19.
20. # Load summarization pipeline
21. summarizer = pipeline("summarization",
   model="facebook/bart-large-cnn")
22.
23. # Model input
24. class QueryModel(BaseModel):
25.     question: str
26.
27. # Fungsi klasifikasi intent
28. def classify_intent(text: str) -> str:
29.     try:
30.         inputs = intent_tokenizer(text,
   return_tensors="pt", truncation=True, padding=True)
31.         with torch.no_grad():
32.             outputs = intent_model(**inputs)
33.             predicted_class_id = torch.argmax(outputs.logits,
   dim=1).item()
34.             return label_map[predicted_class_id]
35.     except Exception as e:
```

```

36.         print(f"✘ classify_intent error: {e}")
37.         return "unknown"
38.
39.     # Daftar sinonim topik
40.     TOPIC_SYNONYMS = {
41.         "prestasi": ["prestasi", "promosi", "penghargaan",
42. "kenaikan jabatan"],
43.         "cuti": ["cuti", "izin tidak masuk", "izin kerja",
44. "libur"],
45.         "mutasi": ["mutasi", "perpindahan", "rotasi"],
46.         "absensi": ["absensi", "kehadiran", "presensi"],
47.         "Permohonan": ["zoom meeting", "jaringan", "desain
48. grafis"],
49.     # Tambahkan sesuai kebutuhan
50.     }
51.     # Endpoint NLP utama
52.     @app.post("/nlp")
53.     def get_entities(query: QueryModel):
54.         text = query.question
55.         print(f"👉 Pertanyaan user: {text}")
56.         intent = classify_intent(text)
57.         print(f"👈 Intent dikenali: {intent}")
58.         nama = bulan = tahun = file_name = topic = None
59.         matching_files = []
60.         # Deteksi nama file PDF
61.         if intent in ["read_pdf_content", "summarize_pdf"]:
62.             match_file =
63.             re.search(r"(?:file|berkas|dokumen) (?:"
64. bernama|)?\s+(.*?\pdf)", text.lower())
65.             if match_file:
66.                 file_name =
67.                 match_file.group(1).strip().replace(" ", "%20")
68.         # Deteksi berdasarkan kriteria
69.         if intent == "search_by_criteria":
70.             match_nama =
71.             re.search(r"(?:atas
72. nama|milik|punya|untuk|dari|berkas (?:"
73. dari)?) (?:\s+(?:pak|ibu))\s+([a-zA-Z]+(?:\s[a-zA-Z]+)*)",
74. text.lower())
75.             if match_nama:
76.                 nama = match_nama.group(1).strip().lower()
77.         month_map = {

```

```

73.         "januari": "01", "februari": "02", "maret":
"03", "april": "04",
74.         "mei": "05", "juni": "06", "juli": "07",
"agustus": "08",
75.         "september": "09", "oktober": "10",
"november": "11", "desember": "12"
76.     }
77.     for month_name, month_code in month_map.items():
78.         if month_name in text.lower():
79.             bulan = month_code
80.             break
81.
82.     for y in range(2020, 2031):
83.         if str(y) in text:
84.             tahun = str(y)
85.             break
86.
87.     # Deteksi berdasarkan topik
88.     if intent == "search_by_topic":
89.         topic = None
90.         if isinstance(text, str):
91.             text = text.lower()
92.
93.         match_topic = re.search(
94.             r"(?:tentang|topik|dengan
topik|berkas(?:\s+terkait)?|dokumen|file|mengenai|seputar)(?:\
s+yang)?(?:\s+membahas)?\s+([a-zA-Z\s]+)",
95.             text
96.         )
97.
98.         if match_topic:
99.             topic = match_topic.group(1).strip().lower()
100.        else:
101.            match_simple = re.search(r"tentang\s+([a-zA-
Z\s]+)", text)
102.            if match_simple:
103.                topic =
match_simple.group(1).strip().lower()
104.
105.            if topic:
106.                # Ekspansi keyword berdasarkan sinonim
107.                expanded_keywords = []
108.                for key, synonyms in TOPIC_SYNONYMS.items():
109.                    if topic in synonyms:
110.                        expanded_keywords = synonyms
111.                    break

```

```

112.         if not expanded_keywords:
113.             expanded_keywords = [topic]
114.
115.         folder_path = os.path.join("storage", "app",
    "public")
116.         for file in os.listdir(folder_path):
117.             if file.endswith(".pdf"):
118.                 normalized_filename =
    file.lower().replace("_", " ").replace("-", " ")
119.                 for keyword in expanded_keywords:
120.                     if keyword in
    normalized_filename:
121.                         matching_files.append(file)
122.                         break
123.
124.         return {
125.             "intent": intent,
126.             "nama": nama,
127.             "bulan": bulan,
128.             "tahun": tahun,
129.             "file_name": file_name,
130.             "topic": topic,
131.             "matching_files": matching_files
132.         }
133.
134.     # Endpoint baca isi PDF
135.     @app.get("/read-pdf")
136.     def read_pdf(file_name: str = Query(...,
    alias="file_name")):
137.         try:
138.             decoded_name = file_name.replace('%20', ' ')
139.             file_path = os.path.join("storage", "app",
    "public", decoded_name)
140.
141.             if not os.path.exists(file_path):
142.                 return {"error": "File tidak ditemukan."}
143.
144.             with fitz.open(file_path) as doc:
145.                 text = "\n".join([page.get_text() for page in
    doc])
146.                 return {"content": text if text.strip() else
    "Isi dokumen kosong atau tidak dapat dibaca."}
147.             except Exception as e:
148.                 return {"error": f"Gagal membaca isi dokumen:
    {str(e)}"}
149.

```

```

150.     # Endpoint ringkasan PDF
151.     @app.get("/summarize-pdf")
152.     def summarize_pdf(file_name: str = Query(...,
153.         alias="file_name")):
154.         try:
155.             decoded_name = file_name.replace('%20', ' ')
156.             file_path = os.path.join("storage", "app",
157.                 "public", decoded_name)
158.             if not os.path.exists(file_path):
159.                 return {"error": "File tidak ditemukan."}
160.             with fitz.open(file_path) as doc:
161.                 text = "\n".join([page.get_text() for page in
162.                     doc])
163.                 if not text.strip():
164.                     return {"summary": "Isi dokumen kosong atau
165.                         tidak dapat dibaca."}
166.                 max_len = 1024
167.                 text = text.replace("\n", " ")
168.                 chunks = [text[i:i + max_len] for i in range(0,
169.                     len(text), max_len)][:3]
170.                 summaries = []
171.                 for chunk in chunks:
172.                     try:
173.                         result = summarizer(chunk,
174.                             max_length=130, min_length=30, do_sample=False)
175.                         summaries.append(result[0]['summary_text'])
176.                     except Exception as e:
177.                         summaries.append("[Gagal meringkas
178.                             sebagian isi dokumen]")
179.                 return {"summary": " ".join(summaries)}
180.             except Exception as e:
181.                 return {"error": f"Gagal meringkas isi dokumen:
182.                     {str(e)}"}
183.     # Endpoint untuk pencarian berdasarkan topik langsung
184.     @app.get("/search-by-topic")
185.     def search_by_topic(topic: str = Query(...,
186.         alias="topic")):

```

```

186.         try:
187.             topic = topic.lower()
188.             matching_files = []
189.
190.             # Ekspansi keyword berdasar sinonim
191.             expanded_keywords = []
192.             for key, synonyms in TOPIC_SYNONYMS.items():
193.                 if topic in synonyms:
194.                     expanded_keywords = synonyms
195.                     break
196.             if not expanded_keywords:
197.                 expanded_keywords = [topic]
198.
199.             folder_path = os.path.join("storage", "app",
    "public")
200.             for file in os.listdir(folder_path):
201.                 if file.endswith(".pdf"):
202.                     normalized_filename =
    file.lower().replace("_", " ").replace("-", " ")
203.                     for keyword in expanded_keywords:
204.                         if keyword in normalized_filename:
205.                             matching_files.append(file)
206.                             break
207.
208.             if not matching_files:
209.                 return {"message": f"Tidak ditemukan berkas
    terkait topik '{topic}'."}
210.
211.             return {
212.                 "topic": topic,
213.                 "matching_files": matching_files
214.             }
215.
216.         except Exception as e:
217.             return {"error": f"Gagal mencari berkas
    berdasarkan topik: {str(e)}"}

```

## Source code

### Train\_intent\_classifier.py

```
1. from datasets import load_dataset
2. from transformers import AutoTokenizer, AutoModelForSequenceClassification, Trainer, TrainingArguments
3. import torch

4. # Model IndoBERT
5. model_name = "indobenchmark/indobert-base-p1"

6. # Load dataset
7. dataset = load_dataset("json",
    data_files="intent_dataset.jsonl", split="train")

8. # Label encoding
9. unique_labels = sorted(set(example["label"] for example in
    dataset))
10. label2id = {label: idx for idx, label in
    enumerate(unique_labels)}
11. id2label = {idx: label for label, idx in
    label2id.items()}

12. # Map label string ke ID numerik
13. dataset = dataset.map(lambda x: {"label":
    label2id[x["label"]]])

14. # Tokenizer
15. tokenizer = AutoTokenizer.from_pretrained(model_name)

16. # Tokenisasi + padding dan truncation dengan panjang
    maksimal 128
17. tokenized_dataset = dataset.map(
18. lambda x: tokenizer(x["text"], truncation=True,
    padding="max_length", max_length=128),
19. batched=True
20. )

21. # Load model untuk klasifikasi
22. model = AutoModelForSequenceClassification.from_pretrained(
23. model_name,
24. num_labels=len(label2id),
25. id2label=id2label,
26. label2id=label2id
27. )

28. # Argumen pelatihan
29. training_args = TrainingArguments(
```

```
30.     output_dir="./bert-intent-model",
31.     per_device_train_batch_size=8,
32.     num_train_epochs=5,
33.     logging_dir="./logs",
34.     save_strategy="epoch", # Simpan model setiap epoch
35.     save_total_limit=2,
36.     load_best_model_at_end=False, # Tidak perlu karena
tidak ada eval dataset
37.     logging_steps=50
38. )

39. # Trainer
40. trainer = Trainer(
41.     model=model,
42.     args=training_args,
43.     train_dataset=tokenized_dataset
44. )

45. # Jalankan training
46. trainer.train()

47. # Simpan model dan tokenizer
48. model.save_pretrained("bert-intent-model")
49. tokenizer.save_pretrained("bert-intent-model")
```