# AudiTREE

Albuddy's Solution for:
**BPK HACKATHON** 2020

Jan 13 2021

# Character Introduction

Audi is a professional auditor

He works more than 8 years in a well-known Audit Board in a country

During audit period, he spends a considerable amount of time and energy to process files and documents from auditee

He thinks whether his works should always last like this...

# Then he meets his pal, Albuddy ..

He states his problems are:

- His physical files and documents are **SCATTERED** everywhere.
- He mostly uses his intuitive on **TEXT REPORTS** for Compliance Analysis.
- He search files and documents based on **TOKEN MATCHING**.
- He find it so hard to predict fraud/**FAKE NUMERICAL** value.
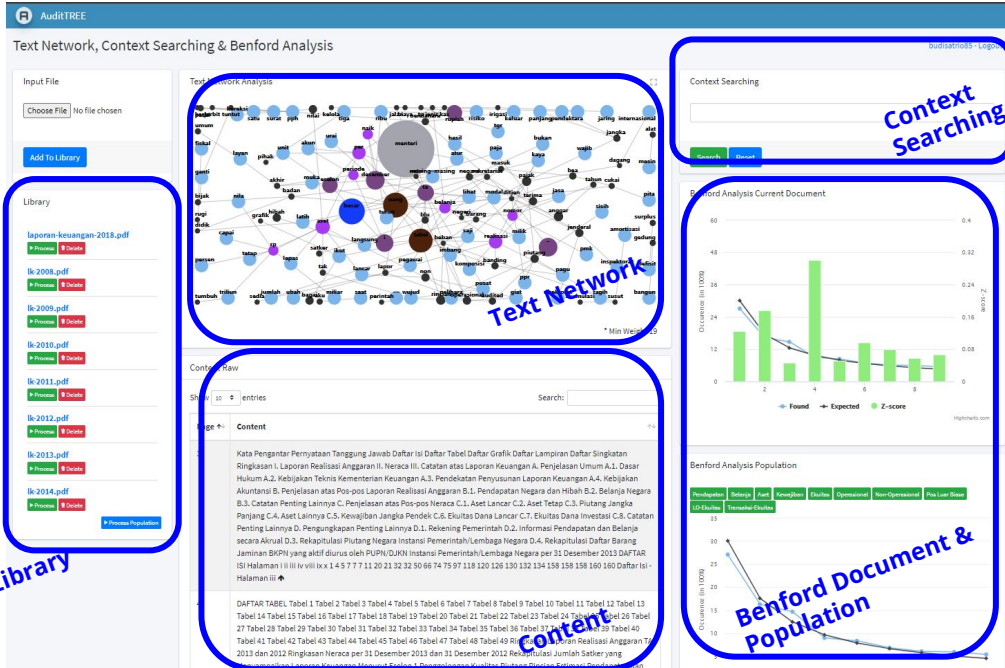
# Albuddy introduces AudiTREE

AudiTREE has the solutions of his problems:

- An **ELECTRONIC LIBRARY** of files and documents
- A text processor to enrich analysis with **NETWORK GRAPH VISUALIZATION**.
- A searching tool based on **CONTEXT MATCHING** of his files and documents.
- An intuitive chart and value of **BENFORD ALGORITHM** to predict fake numerical value.

# Audi doing trial

- He gathers PDF document of Annual **FINANCIAL REPORT**S
- Documents are sampled from **21 MINISTRIES** where extractable
- He inserts financial report of a ministry from 2008 to 2014, and 2018 as his trial

# Jos!, said Audi to the results



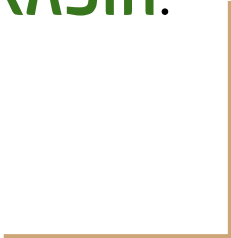Audi finds AudiTREE really helps him accelerate his job

He tells his boss and colleagues about it and also proposes it as a central tools to be used when auditing

# Audi now designs the Future Works

- **OPTIMIZE** the **CONTEXT** search by fine tuning on financial domain
- **ENRICH**  data **GATHERING** to include scanned images using OCR
- Add **NER** (Named Entity Recognition) into Text Network Analysis

Ke Sabah membeli rotan,
pulang petang bawa selasih.
Ada salah mohon benarkan,
ini sekian **TERIMA KASIH**.

# Technology Stack

- Python
- Flask
- AdminLTE by AppSeed
- Networkx
- Highcharts
- Faizz
- Benford-Py

# Data Pipeline

1) **Search Tool**

   PDF → Convert to Text → Cleansing → Sentence/Paragraph/Summary list → BERT Embedding list → FAISS Clustering → Model Serialization → Search API → Inference Task

2) **Text Network Analysis**

   PDF → Convert to Text → Cleansing → Summary → Most relevant topics and their relations → visualize by Networkx→ Text Network Analysis

3) **Benford Histogram**

   PDF → Convert to Text → Cleansing → Crawl variable & numerical value from balance sheet → Cluster into Asset, Liability, Debts, Income, etc → Collect previous data → Calculate histogram

# Process

1.  Research of audit documents
2.  Scraping data
3.  Cleaning
4.  Text Analysis
5.  Visualization

# Team Roles

1.  Data Modeling/Mining : **AL**wi husada
2.  LEAD, FrontEnd Engineer : **BUD**i satrio
3.  Data Engineering : re**DY** andriyansah

# Tools

- **IndoBERT** untuk text embedding
- BERT extractive **Summarizer**
- Facebook AI Similarity Search (**FAISS**)
- **NeoVis** untuk Text Graph Analysis
- **Uvicorn** Web Server