



MINISTRY OF ECONOMY
DEPARTMENT OF STATISTICS MALAYSIA

Application of Language Models to Analyse Scanned Textual Compliance Reports*

Chun Onn Pam; Nur Aisyah binti Mohd Nasir; Siti Hajar binti Khairul Anwar
Central Bank of Malaysia
Payment Services Oversight Department

(*) The views and conclusions expressed herein are exclusively those of the author(s) and do not necessarily reflect the position of the Central Bank of Malaysia or of the Board members.

**11th MALAYSIA
STATISTICS CONFERENCE**
"Data and Artificial Intelligence: Empowering the Future"

**19th September
2024**

Organized by:



Application of language models enhances the utilisation of text-based reports

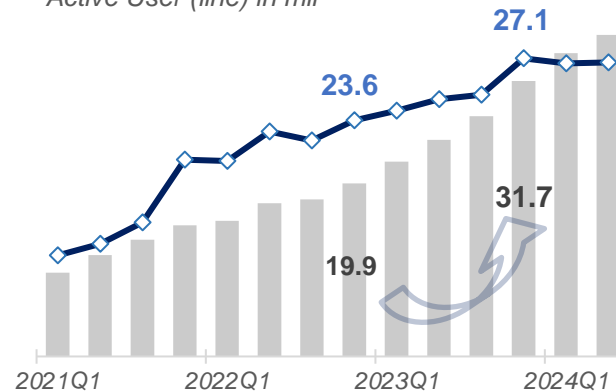


Supervisory authorities leverage on **Suptech** to support data-driven and timely monitoring for more effective supervision. Non-bank payment supervisors in BNM supervise over **50 e-money issuers (EMI)** and **260 money services businesses (MSB)**.



The **fund-safeguarding and/or capital adequacy** of the EMIs/MSBs are closely monitored, to ensure the protection of customers' fund and the sustainability of regulated businesses. EMIs and MSBs are expected to operate effectively, as transaction value and number of customers have increased in recent years.

E-Money Turnover (bar) in RM bil
Active User (line) in mil



Growing number of large size MSBs

2022: 83

2023: 109



MSBs and EMIs submit **scanned text-based compliance reports** prepared by external auditors to BNM, consisting of:

- Agreed-Upon Procedures (AUP)
- Fund Management Report (FMR)



Reports are **large** in number, **lengthy** and in **various formats** (imaged-based PDF, table-based report slides)- straining supervisors' resources to analyse this information efficiently.

This paper focuses on the application of:

- Optical Character Recognition (OCR)** to convert the scanned textual reports into machine-readable text; and
- language models (LMs)** and to provide a concise overview of the report contents

This compliance tool is deployed in Q1 2024, currently utilised by supervisors to assess non-compliances detected from the reports.

Extracting and summarising textual reports to monitor compliance with regulatory requirement

Assessment by auditors

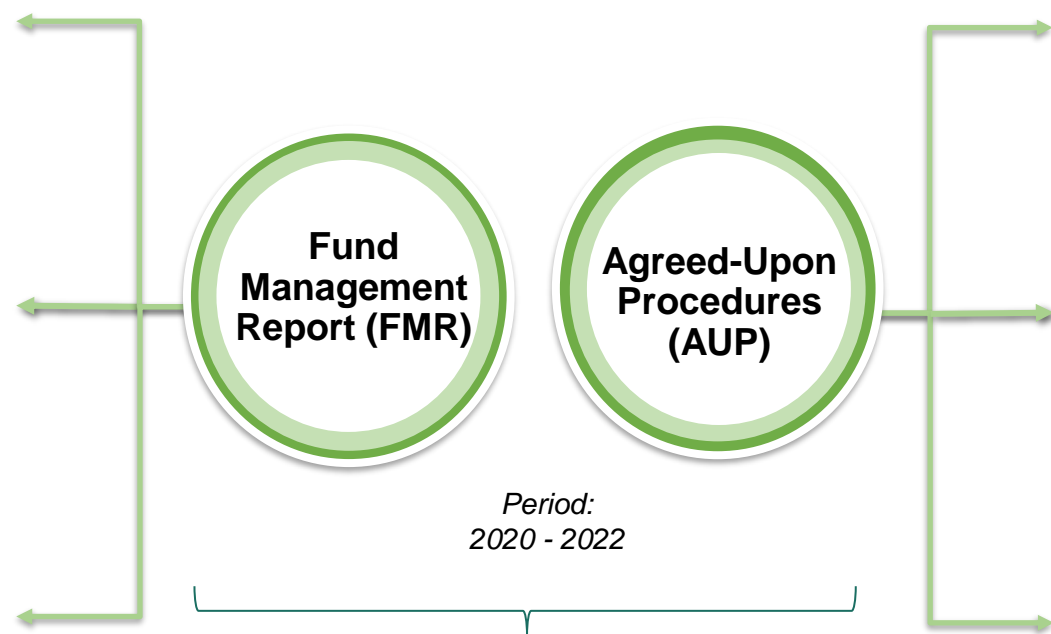
- Separation of funds to Trust account and the usage of funds.
- Trust Account balance and the effectiveness to top up in a timely manner.

Reports

- 123 reports received.
- **10 pages** (min: 1 page, max: 40 pages)

Sample Reports

Area of focus	Method of assessment	Audit techniques used	Observations	Result of audit
1. The separation of funds collected from users from the company's general account of revenue (under a dedicated account) or a separate account to be deposited in the public trust account.	1. Check the bank statements for the account from 01/01/2020 to 31/12/2022 for justification. 2. Request the bank statement from the account.	1. Company has a bank account and 2 operating accounts responsible for the revenue received from users. 2. 10 samples are submitted from the account with complete audit.	The controls are adequate.	Audit on the method of assessment and justification on the basis of the evidence of users, the audit evidence provided sufficient to provide a basis for our opinion.
2. The balance of the trust with a company bank (dedicated account) or (TRUST ACCOUNT) at the end of the month, to ensure the company has a balance greater than the total sum of the company's outstanding amount.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. Request the bank statement from the account.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. 10 samples are submitted from the account with complete audit.	The controls are adequate.	Audit on the method of assessment and justification on the basis of the evidence of users, the audit evidence provided sufficient to provide a basis for our opinion.
3. The balance of the company's bank account at the end of the month to ensure the company has a balance greater than the total sum of the company's outstanding amount.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. Request the bank statement from the account.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. 10 samples are submitted from the account with complete audit.	The controls are adequate.	Audit on the method of assessment and justification on the basis of the evidence of users, the audit evidence provided sufficient to provide a basis for our opinion.
4. The usage of the trust funds for the intended purpose of the account of the company, to ensure the funds are used for the intended purpose.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. Request the bank statement from the account.	1. Prepare regular daily for payment, balance and outstanding amount for the month to ensure the balance is greater than the liability (TRUST ACCOUNT) and its appropriate audit. 2. 10 samples are submitted from the account with complete audit.	The controls are adequate.	Audit on the method of assessment and justification on the basis of the evidence of users, the audit evidence provided sufficient to provide a basis for our opinion.



Format include table form, words, with imaged-based PDF (scanned PDF)

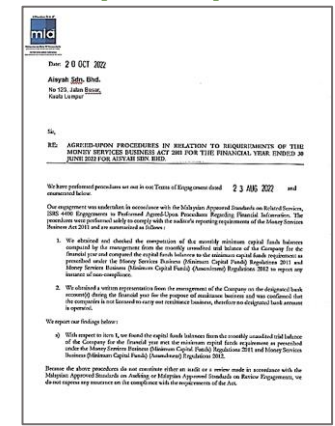
Assessment by auditors

- minimum capital fund, designated account and outstanding customers' monies.

Reports

- 482 reports received.
- **2 pages** (min: 2 pages, max: 14 pages)

Sample Report



Employing OCR to extract image-based reports, and language models for content summarisation



Key steps



Text Extraction

Using **optical character recognition (OCR)** for text extraction from scanned documents.



Text Summarisation

Using **language models (LMs)** to summarise content of the report for its findings.



Performance Evaluation

Evaluate Language Models' performance by comparing the similarity of meaning between the original text and the summarised text.

Detailed steps

...Python library, models and metrics used

Pre-processing

- ✓ PyMuPDF
- ✓ Pillow

Text Extraction

- ✓ PyTesseract
- ✓ Img2table
- ✓ PaddleOCR

Application of 3 LMs were compared for summarisation:



T5



BART



GPT-2

- × Inaccurate summarisation due to limited inputs allowed.

Performance metrics for sentences similarity were evaluated:



ROUGE - L



BLEU



BERTScore

- ✓ RegEx is applied to check on the compliance; or
- ✓ NLP classification model is applied for reports lacking keywords to determine compliance status.

...challenges faced, or highlights

- × Table inside a table or complex table.
- × Small/unrecognizable characters.
- × Inconsistent text format.

Text Extraction from table: Img2Table & PaddleOCR outperforms PyTesseract



Text Extraction (OCR)

PyTesseract, Img2table and PaddleOCR are used to accommodate different report formats.

- ✓ Accurate extraction of text paragraph from scanned document.
- ✓ Manage to extract tables in proper table format.
- ✓ Able to analyse data offline on local devices without sending through API-based engine.

Scanned Compliance Report using different OCR models:

(a) Scanned Compliance Report

No.	Procedures	Findings
1	Obtained and checked the computation of the monthly minimum capital funds balances computed by management of the Company (as disclosed in Appendix 1) from the monthly unaudited trial balance of the Company for the financial year, and compared the capital funds balances to the minimum capital funds requirement as prescribed under the Money Services Business (Minimum Capital Funds) Regulations 2011, Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2012 and Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2013 to report any instance of non-compliance.	With respect to the procedures, there were no exceptions noted.

(b) Extracted text data from table using *PyTesseract*

No.	Procedures	Findings
4	Obtained and checked the computation of the monthly minimum capital funds balances computed by management of the Company (as disclosed in Appendix 1) from the monthly unaudited trial balance of the Company for the financial year, and compared the capital funds balances to the minimum capital funds requirement as prescribed under the Money Services Business (Minimum Capital Funds) Regulations 2011, Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2012 and Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2013 to report any instance of non-compliance.	With respect to the procedures, there were no exceptions noted.

(c) Extracted text data from table using *Img2Table with PaddleOCR*

No	Procedures	Findings
1	Obtained and checked the computation of the monthly minimum capital funds balances computed by management of the Company (as disclosed in Appendix 1) from the monthly unaudited trial balance of the Company for the financial year, and compared the capital funds balances to the minimum capital funds requirement as prescribed under the Money Services Business (Minimum Capital Funds) Regulations 2011, Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2012 and Money Services Business (Minimum Capital Funds) (Amendment) Regulations 2013 to report any instance of non-compliance.	With respect to the procedures, there were no exceptions noted.

Text Summarisation: BART performed the best among three LMs



Performance Evaluation

Evaluation on the 3 LMs is measured using ROUGE - L, BLEU and BERTScore.

BART achieves highest similarity score across all metrics.

- ✓ Generate concise and accurate summary.
- ✓ Content of the summaries are similar to the original text.

BART result:

(a) Extracted Text from Compliance Report


1. We obtained and checked the computation of the monthly minimum capital funds balances computed by the management from the monthly unaudited trial balance of the Company for the financial year and compared the capital funds balances to the minimum capital funds requirement as prescribed under the Money Services Business (Minimum Capital Funds) Regulation 2011 and Money Services Business (Minimum Capital Funds) (Amendment) Regulation 2012 and Money Services Business (Minimum Capital Funds) (Amendment) Regulation 2013 to report any instance of noncompliance. There was no instance of noncompliance to be reported for the financial year ended 31 August 2022.

(b) Summarised Text

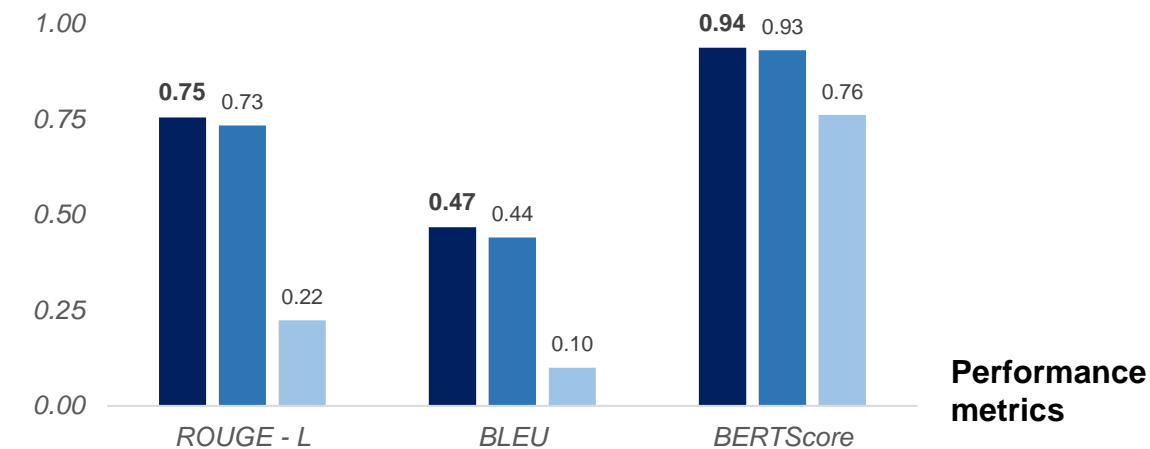
We obtained and checked the computation of the monthly minimum capital funds balances computed by the management from the monthly unaudited trial balance of the Company. There was no instance of noncompliance to be reported for the financial year ended 31 August 2022

Performance Evaluation on Text Summarisation Model

■ BART ■ T5 ■ GPT-2

 BART achieves highest similarity score across all metrics

Similarity Score



Text analysis output is visualised in dashboard for supervisors' monitoring



Enhanced supervisor efficiency: Reduce need for manual reading of reports.

Facilitates decision-making: Provides quick summaries of past submissions for a holistic view of compliance history of specific entity.

Information extracted from methodology

Visualised in

OCR for text extraction

Regex for word search

The keywords such as '**Exception noted**', '**Satisfactory**', '**met**', and '**not met**' and others are located for compliance checking.

LLM for text summarisation and classification

- ✓ Name of auditors.
- ✓ Content of the reports.
- ✓ Number of regulatees meeting and/or breaching the requirement.
- ✓ Summary of requirement breach, classified by specific act / compliance.

- AUP Monitoring Dashboard*
- FMR Monitoring Dashboard*

Example outcome: Several regulatees which did not meet capital adequacy requirement were imposed monetary penalty.

* The dashboards are not displayed due to privacy and confidentiality.

Application of language models enhances the utilisation of text-based reports

Summarisation accuracy matters

- BART outperforms other language models in text summarisation.

Flexible in approach

- Combining different OCR techniques for different reporting structures enhance data extraction and improve compliance checking.
- Language models (LMs) serve as an effective tool to help detect compliance, for reports lacking keywords to determine regulatory adherence.



Proprietary OCR Tools

- Paid services like Adobe Acrobat and Azure Intelligent Document Processing offer superior performance.

Multimodal Large Language Models

- Direct image pre-processing for NLP tasks like text summarisation and classification.

Report Standardisation

- Enhancing data extraction by imposing standardised reporting structure.
- Standardised and clearer policy assessment for auditors on compliance.

Thank you

**11th MALAYSIA
STATISTICS CONFERENCE**
"Data and Artificial Intelligence: Empowering the Future"

**19th September
2024**

Organized by:



	RegEx (Regular Expression)	Language Models (T5, BART, and GPT-2)	Performance Metrics (ROUGE – L , BLEU, BERTScore)
<i>What is it?</i>	<ul style="list-style-type: none"> Regex is a tool for defining patterns to search for certain characters or words inside text-based content Allows to process text and make locating information much easier Aside from literal characters (like 'abc'), Regex can also search other special characters (*,+,?,@ and so on) which may be relevant in the report 	<ul style="list-style-type: none"> LMs are pre-trained on a vast amount of text data to generate human-like responses or perform language-related tasks, including summarisation or classification An on-premises LM AI model allow LM to operate within our own infrastructure, rather than in the cloud ensuring no confidential data shared to other entity. LMs are transformer¹-based models capable of performing language tasks – <ul style="list-style-type: none"> T5 or Text-to-Text Transfer Transformer (Encoder-Decoder), developed by Google, converts all tasks into a text format, making it versatile for various text-related tasks. BART (Encoder-Decoder), developed by Meta, corrupts text with noise and learns to reconstruct the original text, making it effective for text generation and summarisation. GPT-2 (Decoder-Only), developed by OpenAI, generates relevant and coherent text by predicting the next word sequence. 	<ul style="list-style-type: none"> Performance metrics are used to calculate the similarity score to evaluate the performance of the summarisation models ROUGE – L formula: $ROUGE-L, F1 - score = \frac{2 \times (precision \times recall)}{precision + recall}$ BLEU formula: $BLEU = brevity_penalty \times e^{\sum_{n=1}^N w_n \log(p_n)}$ BERTScore formula: $Similarity(A, B) = \cos(\theta) = \frac{A \cdot B}{\ A\ \ B\ }$
<i>How we use it?</i>	<ul style="list-style-type: none"> We locate the keywords such as 'Exception noted', 'Satisfactory', 'met', and 'not met' and others for compliance checking. 	<ul style="list-style-type: none"> We use 3 LMs to summarize text and identify the one with the highest score as our champion LM. We use BART to classify the non-regex cases. (<i>Reports that do not have the specific word that regex can detect</i>) 	<ul style="list-style-type: none"> We calculate the average of the similarity for each performance metrics for all reports

¹ Transformers use an attention mechanism to understand the context and dependencies between words in a sentence. This mechanism allows the model to focus on different parts of the input sequence when generating the output sequence.