# Constructing and Benchmarking: a Labeled Email Dataset for Text-Based Phishing and Spam Detection Framework

Rebeka Tóth[1], Tamas Bisztray[1], Richard A. Dubniczky[2]

[1]University of Oslo, Oslo, Norway

[2]Eötvös Loránd University, Budapest, Hungary

rebekat@uio.no, tamasbi@uio.no, richard@dubniczky.com

Abstract:

Phishing and spam emails remain a major cybersecurity threat, with attackers increasingly leveraging Large Language Models (LLMs) to craft highly deceptive content. This study presents a comprehensive email dataset containing phishing, spam, and legitimate messages, explicitly distinguishing between human- and LLM-generated content. Each email is annotated with its category, emotional appeal (e.g., urgency, fear, authority), and underlying motivation (e.g., link-following, credential theft, financial fraud). We benchmark multiple LLMs on their ability to identify these emotional and motivational cues and select the most reliable model to annotate the full dataset. To evaluate classification robustness, emails were also rephrased using several LLMs while preserving meaning and intent. A state-of-the-art LLM was then assessed on its performance across both original and rephrased emails using expert-labeled ground truth. The results highlight strong phishing detection capabilities but reveal persistent challenges in distinguishing spam from legitimate emails. Our dataset and evaluation framework contribute to improving AI-assisted email security systems. To support open science, all code, templates, and resources are available on our project site.

## 1 INTRODUCTION

Phishing and spam emails continue to pose a significant threat to cybersecurity, targeting individuals and organizations with deceptive messages designed to steal sensitive information, compromise systems, or facilitate fraudulent activities. According to the ENISA Landscape 2025 report, phishing continues to dominate initial access techniques across major incident categories, with a steady year-over-year increase in both volume and sophistication [for Cybersecurity (ENISA), 2025]. As email-based attacks evolve, cybercriminals increasingly leverage advanced technologies, including LLMs, to craft highly convincing and tailored phishing and spam messages capable of bypassing traditional detection mechanisms. Studies show that AI-generated phishing emails can achieve click-through rates as high as 54%, compared to just 12% for generic phishing attempts [Heiding et al., 2024a]. This rapid increase in sophistication highlights the need for more robust and adaptive email security solutions.

Traditional email filtering techniques rely on rule-based systems, heuristics, or machine learning classifiers trained on features such as sender reputation, keyword frequency, or metadata. Although effective to some extent, these approaches often struggle to adapt to modern phishing strategies, particularly when attackers employ LLMs to create human-like messages that evade conventional detection. The 2025 State of the Phish report notes that 99% of phishing emails now incorporate some form of social engineering, making detection based solely on static features increasingly unreliable [Proofpoint, 2025]. Complementary studies also show that spam messages accounted for over 45.8% of global email traffic in 2023 [Petrosyan, 2024], while traditional spam filters may miss up to 25% of sophisticated phishing attempts [Josten and Weis, 2024].

Recent advances in natural language processing (NLP) and deep learning offer promising alternatives. LLMs have demonstrated impressive capabilities in text comprehension, classification, and contextual reasoning; however, their performance in real-world email threat detection remains insufficiently explored [Tarapiah and Others, 2025]. Understanding their strengths and limitations—especially in scenarios involving emotional manipulation, paraphrasing, or adversarial content—is essential for improving automated cybersecurity systems.

In this study, we develop and label a comprehensive dataset containing phishing, spam, and legitimate emails, distinguishing between human- and LLM-generated content. Each email is annotated with its type (phishing, spam, legitimate), underlying motivation (e.g., opening an attachment, following a link, subscription fraud), and emotional appeal (e.g., urgency, greed, curiosity). To assess LLM capabilities, we utilized multiple models on their ability to label emotional and motivational cues as well as rephrase emails using only the message body (including full link addresses and attachment names with extensions), subject line, and sender information, comparing the results against expert-validated ground-truth labels. Through this we create a robust labeled and detailed dataset for benchmark or training LLMs to evaluation, and identify key strengths, weaknesses, and potential failure cases in LLM-based email threat detection.

The study addresses the following research questions:

- RQ1: How effectively can LLMs identify emotional strategies and patterns used within email content?

- RQ2: How effectively can LLMs identify motivational strategies and patterns used within email content?

- RQ3: How effectively can a state-of-the-art LLM perform email classification of phishing, spam, and legitimate emails across both human-written and LLM-generated content?

The primary aim of this study is to develop a comprehensive and extensible labeled dataset for research on email threat detection. While we evaluate several state-of-the-art LLMs, the evaluation primarily serves to validate the dataset creation process and demonstrate how the data can support future research. The key contributions of this work are:

- Dataset construction and expansion: We construct a new multi-source email dataset that integrates recent real-world emails, curated public datasets, and controlled synthetic samples. Every email is consistently structured and labeled with its type, emotional appeal, attacker motivation, linguistic source (human or LLM), and year of origin. The resulting dataset is designed to be extensible and suitable for future risk-modeling and NLP research on email security. We also provide a framework for further expansion of the dataset by making all used materials open source.

- Emotional and motivational annotation: We annotate every email using an LLM to capture emotional tone and attacker motivation. A randomly selected subset of 100 emails is manually reviewed by a human expert to validate annotation reliability. In addition, we conduct a small comparative benchmark across several LLMs to determine which model provides the most consistent and accurate annotation. This evaluation allows us to select the most suitable model for labeling the full dataset while also offering insight into current LLM capabilities for emotion and motivation recognition in email contexts. The resulting multi-label emotional metadata captures the manipulative strategies commonly used in social engineering attacks.

- Rephrased variants: To evaluate robustness, we generate LLM-rephrased variants of all emails, preserving intent and emotional characteristics while altering linguistic expression.

Overall, this work contributes a high-quality, labeled, and extensible dataset together with a systematic framework for analyzing emotional manipulation, attacker intent, and classification robustness in email security. By emphasizing dataset construction and providing all processing scripts, rephrasing pipelines, and analysis tools, we aim to support open, reproducible research and enable continued dataset growth. All templates, code, and supplementary materials are publicly available at https://github.com/DataPhish/PhishingSpamDataSet.

## 2 RELATED LITERATURE

### 2.1 Spam Filtering

Modern spam filtering systems rely on a combination of rule-based heuristics, machine learning models, and reputation-based methods. Prominent commercial and open-source tools include Gmail's built-in spam filter and Apache SpamAssassin [Apache Software Foundation, 2025].

Several studies have evaluated the effectiveness of machine learning techniques for spam detection. Tusher et al. provide a comprehensive review of classical and deep learning classifiers—including Naïve Bayes, Support Vector Machines, Decision Trees, Random Forests, and neural networks—highlighting their performance characteristics and challenges such as concept drift, adversarial inputs, and real-time processing constraints [Tusher et al., 2024].

Mardiansyah et al. examine the capabilities of ChatGPT-4 and Google Gemini for spam classification using the SpamAssassin dataset. Their results show that ChatGPT-4 achieves balanced precision and recall, while Gemini obtains higher recall but at the cost of more false positives, indicating different trade-offs depending on deployment needs [Mardiansyah and Surya, 2024].

### 2.2 Phishing Detection

Phishing detection remains an active research area encompassing rule-based approaches, machine learning, and more recently LLM-based techniques. Eilertsen et al. introduce an LLM-powered framework for intent-based classification of phishing emails, focusing on identifying attacker objectives rather than binary detection [Eilertsen et al., 2025]. While their approach enriches threat intelligence analysis, it does not address emotional manipulation, paraphrasing robustness, or metadata cues.

Ige et al. provide a broad survey of Bayesian, non-Bayesian, and deep learning approaches to phishing detection, outlining their strengths, limitations, and typical deployment challenges [Ige et al., 2024]. As phishing attacks become increasingly AI-enhanced, studies have also examined the offensive capabilities of LLMs. Heiding et al. demonstrate that fully AI-generated phishing emails can match human experts in effectiveness and significantly outperform generic phishing attempts, highlighting the growing threat of AI-enabled social engineering [Heiding et al., 2024a].

More recent work by Afane et al. investigates the ability of modern detectors—including Gmail's spam filter, SpamAssassin, Proofpoint, and classical machine learning classifiers—to identify both traditional and LLM-rephrased phishing emails. Detection accuracy drops substantially across all systems when emails are paraphrased by an LLM, exposing critical weaknesses in current defenses [Afane et al., 2024].

Recent studies also, have evaluated traditional ML and deep-learning models for phishing detection, such as the comprehensive benchmark by Alhuzali et al. (2025), which compared fourteen algorithms across multiple public datasets and demonstrated high accuracy on static email frames. However, these approaches rely on fixed datasets and do not incorporate emotional or motivational cues or evaluate robustness against paraphrased or LLM-generated variants. Our work addresses these gaps by constructing a richer dataset and analyzing how LLMs behave under semantic variation, emotional manipulation, and rephrasing [Zhang et al., 2025].

### 2.3 Emotion Analysis

Emotional manipulation is a foundational component of many phishing and spam campaigns. Prior work has examined the psychological tactics used in deceptive emails, showing that attackers frequently employ fear, urgency, authority, curiosity, altruism, and financial incentives to influence user behavior [Wang and Lutchkus, 2023]. While existing research has explored sentiment and emotion detection in text, relatively few studies analyze emotional strategies specifically within phishing and spam emails or examine how LLMs handle such cues. Our work extends this area by incorporating detailed emotional and motivational annotations and evaluating how well LLMs leverage these signals during classification. Prior research demonstrates that GPT-4 and other frontier LLMs achieve near–state-of-the-art performance in zero- and few-shot emotion classification, often matching or exceeding supervised baselines across multiple datasets [Gilardi et al., 2023, Niu et al., 2024]. Benchmarks further indicate that GPT-4 performs strongly across diverse sentiment and affect recognition tasks, showing robustness for nuanced textual labeling. In phishing-specific contexts, recent studies have shown that GPT-4 and related models outperform smaller LLMs in detecting malicious intent and categorizing deceptive tactics [Afane
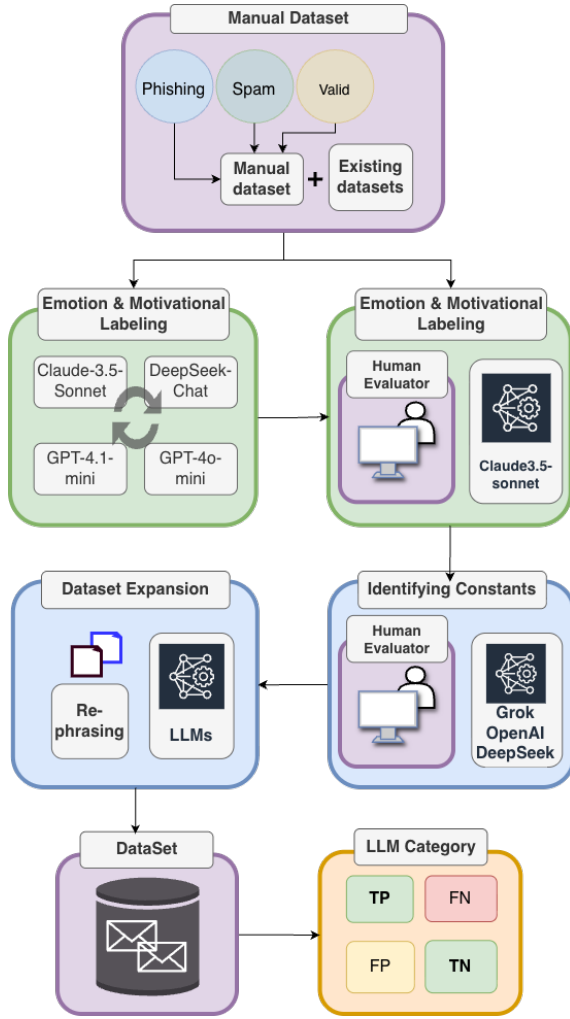
Figure 1: Methodology for dataset building and evaluating.

et al., 2024].

# 3 METHODOLOGY

Figure 1 shows a visual overview of our methodology. The research began by collecting recent phishing emails including both emails from enterprise phishing simulations used for cyber awareness trainings and phishing emails received from outside sources, spam emails, and legitimate emails. After creating a manual dataset of about 3,000 emails in each category from both work and personal accounts, a few entries from already existing datasets were added to expand the diversity. After merging the newly collected emails with existing datasets, we conducted an emotion-motivation labeling phase using multiple LLMs.

Each model was benchmarked on its ability to detect emotional cues and attacker motivations, allowing us to compare accuracy, consistency, and agreement with human annotations. Once the most reliable model was identified, it was used to label the full dataset. To validate the quality of this labeling, a subset of 100 emails was manually reviewed by two experts—a psychology specialist and a cyber-awareness analyst—to ensure alignment between LLM-derived and human—interpreted emotional signals. To create an even more diverse dataset each email has been given to multiple LLM agents to be rephrased, while keeping the context, intent and the emotion of the original email. After creating, labeling and expanding the entire dataset, it was given to an LLM agent that was excluded in the rephrasing stage to test one of the state-of-the-art best LLM model's efficacy in detecting and separating phishing and spam from legitimate emails.

## 3.1 Ethics

The research followed all Norwegian national and University of Oslo rules when collecting the emails for the dataset. Participants were asked to provide the body, title, and sender of the emails and were asked to remove all personal information. All submitted emails were checked and any additional private information left in the text was removed to ensure full anonymity. Participants who provided emails were given a consent form. The study was carried out in accordance with the research ethics guidelines of the University of Oslo. Participants were fully informed about the procedures involved and were free to withdraw at any point.

## 3.2 Dataset

### 3.2.1 Manually Created Dataset

To create the dataset, we collected emails across various platforms and sources including personal email accounts, work email accounts, and emails from existing databases, in three different categories: phishing (including enterprise phishing simulations), spam, and legitimate emails.

There are existing datasets like CAES_08, Nigerian Frauds, Enron or Ling on Kaggle [Alam, 2024], or designated spam datasets like the Guenter Spam Archive [Guenter, 2025] and the Warranted SPAM Archive [Martin, 2005]. These datasets were created using a similar approach

to ours, where information such as sender, receiver, date, subject, body, and label were collected. However, our dataset contains more up-to-date emails and additionally includes the full URL addresses, the names and extensions of file attachments, and distinctions between spam and phishing. To diversify our dataset, we included a subset of emails from the Phishing Email Dataset on Kaggle [Alam, 2024] and the Warranted SPAM Archive [Martin, 2005]. While existing datasets were used to expand our dataset and analyze patterns across years, our dataset further includes labeling such as emotional appeal and motivational intent.

As shown in the dataset at our project site: https://github.com/DataPhish/PhishingSpamDataSet, a each email was labeled with a set of metadata fields capturing its content, origin, and emotional and motivational characteristics:

- Subject: The email subject line.

- Body: The full text of the email body.

- Sender: The email address of the sender.

- URLs: Full URLs extracted from the message.

- Attachments: File names and extensions (if present).

- Motivation: Attacker intention (e.g., follow the link, open attachment, reply, financial fraud, credential harvesting).

- Emotion: Primary emotional cues (e.g., urgency, greed, altruism, fear, curiosity).

- Type: Phishing, spam, or legitimate.

- Created by: Whether the email was human-written or LLM-generated.

- Source: Source of the email (e.g., personal or corproate inbox, LLM rephrase, CAES_08, Warranted SPAM Archive).

- Year: The year the email was sent or collected.

### 3.2.2 Emotional and Motivational Analysis

Emotional and motivational labeling in this study builds upon prior research that investigated how phishing emails evoke psychological responses influencing user behavior from Tóth et al. That work identified key emotional indicators—such as urgency, fear, altruism, greed, or perceived scarcity—through a combination of qualitative analysis and a validation questionnaire in which participants rated the emotions elicited by real phishing examples. These categories informed

our annotation framework, as they represent common manipulation strategies used by attackers to prompt user actions, such as clicking malicious links or opening attachments [Tóth et al., 2024].

Using this validated set of emotional and motivational cues, a subset of emails in our dataset were analyzed and labeled. Emotions were interpreted based on textual markers (e.g., time pressure, threats, requests for help), while motivations captured the underlying attacker intent (e.g., inducing link-clicking, data submission, or file execution). To enlarge and diversify the dataset, LLM-generated paraphrased variants were also created from human-written emails, allowing the analysis to cover a broader linguistic space while maintaining the same underlying emotional and motivational structure. The same line of thought were given to LLM agents to label the entire dataset.

To assign emotional and motivational labels to every email, we evaluated multiple large language models using a controlled benchmarking setup. A subset of 100 emails was manually annotated by two human experts—a cyber awareness specialist with a psychology background, and a security researcher. This dual-review process created a high-quality ground truth and also allowed us to validate cases where LLMs identified additional emotional cues that were plausible but not explicitly marked in the initial expert pass.

Four models were selected for benchmarking: GPT-4o-mini, GPT-4.1-mini, Claude 3.5 Sonnet, and DeepSeek-Chat. Each model received the same structured prompt (Listing 1) and processed every email five times independently to evaluate output stability. This resulted in two types of evaluation:

- Internal Consistency: Whether a model produced the same labels across all five runs.

- Agreement with Human Labels: Measured using strict accuracy, close-enough accuracy, Jaccard similarity, precision, and recall.

Across all models, internal consistency was relatively high (78–80% for emotion and 83–88% for motivation). However, reliability and correctness varied substantially across models.

For emotion labeling, Claude 3.5 Sonnet achieved the strongest alignment with human annotations, obtaining a strict accuracy of 25%, a close-enough accuracy of 42%, and the highest average Jaccard similarity (0.60). GPT-4.1-mini followed with 21% strict accuracy and 30% close-enough accuracy (Jaccard 0.57). DeepSeek-Chat

and GPT-4o-mini performed notably weaker on both strict and approximate matching.

For motivation labeling, all models showed lower strict accuracy given the inherently ambiguous nature of email intent. However, close-enough accuracy remained between 53–61% across models. Claude displayed the highest internal consistency and strong recall, although its strict accuracy was low (1%) due to its tendency to infer additional plausible motivations. Human reviewers confirmed that in several cases these additional labels were indeed reasonable interpretations, indicating that strict accuracy alone does not fully capture model performance in this domain.

A summary of the benchmarking results is shown below:

- GPT-4o-mini: Emotion Jaccard 0.46, Close-enough 10%, Motivation close-enough 61%.

- GPT-4.1-mini: Emotion Jaccard 0.57, Close-enough 30%, Motivation close-enough 59%.

- Claude 3.5 Sonnet: Emotion Jaccard 0.60, Close-enough 42%, Motivation close-enough 61%.

- DeepSeek-Chat: Emotion Jaccard 0.45, Close-enough 18%, Motivation close-enough 53%.

Based on the combined criteria—alignment with human annotations, internal consistency, and expert review of ambiguous cases—Claude 3.5 Sonnet was selected as the final model to annotate the emotional and motivational labels for the entire dataset. Its outputs most closely matched human expert judgment while demonstrating stable behavior across repeated runs.

The exact prompt used for all models during emotional and motivational labeling is shown in Listing 1. In each of our prompts for the different methodology parts, the instruction style was optimized for task-specific objectives. Recent work shows that role-playing personas can influence model behavior in complex ways, sometimes improving reasoning, and in other times it would introduce bias or reliability [Kim et al., 2024]. Based on this, we have also conducted a small research on which type of prompt would return better results, role assigning or context assigning. We compared starting the prompt "I am a PhD Researcher and I need this for my Research" with "You are a careful assistant". Based on our tries, persona oriented prompts were more affective as it primes the model to adopt a consistent expert persona aligned with the cybersecurity labeling task to provide better outputs.

---

**Emotional and Motivational Analysis Prompt**

You are a phishing email analyst. Label each email with:
Emotions (1–3): Fear, Curiosity, Greed, Altruism, Urgency, Authority, Unic, Neutral
Motivations (1–2):Follow the link, Open attachment, Reply, Promotion/ad, Blackmail, Financial fraud, Data theft/credential harvesting, Unknown

Email:
Subject: {subject} Body: {body}

Respond ONLY in strict JSON: { "emotions": ["..."], "motivations": ["..."] }

Listing 1: Emotional and Motivational Analysis Prompt.

### 3.2.3 Rephrasing

Since AI-generated phishing is increasing [Heiding et al., 2024b], LLMs were also used to rephrase emails while preserving their intent and emotional tone [Deng et al., 2024].

To ensure that rephrased emails remained true to the original content, we first identified which elements must remain unchanged. A subset of 20 emails was manually reviewed to extract constants such as URLs, filenames, dates, numerical values, identifiers, and named entities. In addition to manual inspection, three different LLMs, GPT-4o, DeepSeek-Chat, and Grok, were asked to determine which components of an email should remain unmodified during rephrasing. The overlapping results from humans and models informed the final constants used in the rephrasing prompt.

Using these constants, the manually collected dataset (3,000 emails) was rephrased using two LLMs—GPT-4o-mini and DeepSeek-Chat—producing about 9,000 emails. A further rephrasing phase via OpenRouter expanded the dataset to approximately 12,000 emails using a wide selection of models (OpenAI, Google, Amazon, Meta, xAI, Qwen, Mistral, DeepSeek).

For every generated email, metadata such as Created by, Rephrase Model, and Rephrase Error were stored to ensure full reproducibility and support future verification and dataset extension.

The exact prompt used for all rephrasing operations is shown in Listing 2.

```
Email Rephrasing Prompt

You are a careful writing assistant. Task:
Rephrase the email's Subject and Body to
new wording but keep: - the SAME meaning,
intent, and emotion; - all entities, numbers,
dates, identifiers, filenames UNCHANGED; -
all URLs EXACTLY as is; - KEEP the origi-
nal language.
Do NOT add warnings. {emotion_line}
Original: Subject: {subject} Body: {body}
Return ONLY JSON: { "rephrased_subject":
"...", "rephrased_body": "..." }
```

Listing 2: Email Rephrasing Prompt.

## 3.3 Categorization

For the final classification step, Claude 3.5 Son-
net was used to categorize each email into one of
three classes: Phishing, Spam, or Valid. This
step relied on a neutral prompting setup in
which the model received only the textual com-
ponents of the email (sender, subject, body),
along with extracted filenames with their ex-
tensions and URLs when applicable. All per-
sonal information (names, addresses, organiza-
tions) was anonymized prior to processing unless
public datasets were used. We selected Claude
3.5 Sonnet because it demonstrated the highest
reliability in our earlier benchmarking of emo-
tional and motivational labeling, and, critically,
it was never used during the rephrasing stage.
This avoids stylistic imprinting or data leakage
and keeps the final classification evaluation com-
pletely neutral.

In addition to assigning a category label, the
model was instructed to provide a short explana-
tion (maximum two sentences) identifying key in-
dicators such as suspicious senders, unusual sub-
jects, file attachments, or potentially malicious
links. The model's response format was strictly
constrained to JSON to support reproducibility
and automated evaluation.

The exact categorization prompt is shown in
Figure 2.

```
Category Prompt

Classify this email as one of: Spam, Valid,
Phishing.
Personal information (names, addresses, etc.)
has been replaced with Anonymous.
Give a short explanation with key indicators
(sender, subject, file, links).
Keep the explanation short (maximum 3
sentences).


Sender: {sender}
Subject: {subject}
Body: {body}
File: {file}
Links: {urls}


Respond ONLY in JSON:
{
"category": "Phishing | Spam | Valid",
"explanation": "short reason here"
}
```

## 3.4 Evaluation Metrics

To evaluate the performance of the language
model classifier, we compare Claude's predicted
labels against the human-labeled ground truth in
the merged dataset. The evaluation is performed
under two parts: (1) strict three-class classifica-
tion and (2) relaxed binary classification.

Strict Classification (Phishing / Spam / Valid).
In the strict setting, the model must correctly
assign one of the three classes. Performance is
quantified using standard metrics derived from
the confusion matrix:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$\text{Precision} = \frac{TP}{TP + FP},$$
$$\text{Recall} = \frac{TP}{TP + FN},$$
$$\text{F1} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$

These metrics are reported for each email
source group (original emails, DeepSeek
rephrases, GPT-4o rephrases, and RandomAPI

rephrases), enabling a comparison of performance across human-written and LLM-generated variants.

**Relaxed Classification (Unwanted vs. Valid).** For many security applications, the practical question is simply whether an email is desired or potentially harmful. To reflect this operational focus, we merge Phishing and Spam into a single Unwanted class:

$$\text{Unwanted} = \{\text{Phishing}, \text{Spam}\}.$$

Accuracy and F1 are then computed using the resulting binary labels. This relaxed evaluation reduces penalty for fine-grained misclassification between harmful subclasses and measures overall detection capability.

**Emotional and Motivational Labeling.** For emotion and motivation annotations, the model predictions are compared with human-reviewed labels using the Jaccard similarity coefficient:

$$\text{Jaccard Similarity} = \frac{|L_{\text{LLM}} \cap L_{\text{human}}|}{|L_{\text{LLM}} \cup L_{\text{human}}|}, \qquad (2)$$

where $L_{\text{LLM}}$ and $L_{\text{human}}$ denote the sets of emotion or motivation tags assigned to each email.

Together, the strict and relaxed classification metrics assess the model's ability to assign correct threat categories, while the Jaccard similarity quantifies its ability to reproduce human judgment of emotional and motivational cues. Evaluating all rephrasing sources ensures that results reflect performance not only on natural emails but also on diverse linguistic variants introduced by different LLMs. This provides a comprehensive measurement of robustness in automated email threat detection.

## 4 DISCUSSION & RESULTS

### 4.1 LLMs' Efficacy in Text-Based Email Filtering

Claude 3.5 Sonnet was used as the final classifier across all versions of the dataset. To assess robustness, each original email was paired with several rephrased variants generated by three independent pipelines: (1) DeepSeek-Chat, (2)

Table 1: Strict three-class evaluation of Claude 3.5 Sonnet on original and rephrased emails. Acc. = overall accuracy; F1_P/S/V = F1 for Phishing/Spam/Valid.

| Group | Acc. (%) | F1_P | F1_S | F1_V |
|---|---|---|---|---|
| Original | 66.89 | 0.937 | 0.208 | 0.639 |
| DeepSeek | 66.34 | 0.936 | 0.197 | 0.631 |
| GPT4o | 66.93 | 0.937 | 0.225 | 0.638 |
| Random | 66.95 | 0.931 | 0.207 | 0.632 |

Table 2: Relaxed binary evaluation (Unwanted = Phishing + Spam vs. Valid) for Claude 3.5 Sonnet.

| Group | Acc. (%) | F1_U | F1_V |
|---|---|---|---|
| Original | 69.53 | 0.736 | 0.639 |
| DeepSeek | 68.99 | 0.731 | 0.634 |
| GPT4o | 69.31 | 0.734 | 0.638 |
| Random | 69.60 | 0.741 | 0.632 |

F1_U = F1 for Unwanted; F1_V = F1 for Valid.

GPT-4o, and (3) a diverse multi-model Open-Router pipeline (Gemini 1.5 Pro, Nova Pro, Llama 3.3, Grok, Mistral Medium, and others). These rephrased samples alter the linguistic surface while preserving semantic content and emotional tone, enabling an evaluation of how stable the classifier remains under paraphrasing.

Table 1 summarizes strict, three-class performance (Phishing, Spam, Valid). Across all four groups, accuracy is highly consistent at 66–67%. Claude performs very strongly on the Phishing class (F1 $\approx 0.93$), while performance on Spam remains weak (F1 $\approx 0.20$–0.23), with many spam messages reclassified as Valid. The Valid class achieves moderate performance (F1 $\approx 0.63$).

### 4.2 Relaxed Evaluation (Unwanted vs. Valid)

Since both phishing and spam constitute unwanted messages, we also evaluate a relaxed binary setting. Table 2 shows that accuracy increases only slightly, to 69–70% across all rephrasing conditions. The Unwanted class reaches F1 values of 0.73–0.74, while the Valid class stays near 0.63–0.64.

The small difference between strict and relaxed accuracy indicates that most errors occur between Phishing and Spam—not between Unwanted and Valid. This suggests that Claude reliably detects that a message is "problematic," even when it struggles to distinguish precise threat types.

Table 3: Emotion labeling performance (100 emails).

| Model | Cons. | Close | Jaccard |
|---|---|---|---|
| GPT4omini | 79% | 10% | 0.46 |
| GPT4.1mini | 78% | 30% | 0.57 |
| Claude3.5Sonnet | 80% | 42% | 0.60 |
| DeepSeekChat | 79% | 18% | 0.45 |

Table 4: Motivation labeling performance (100 emails).

| Model | Cons. | Close | Jaccard |
|---|---|---|---|
| GPT4omini | 83% | 61% | 0.33 |
| GPT4.1mini | 84% | 59% | 0.35 |
| Claude3.5Sonnet | 88% | 61% | 0.23 |
| DeepSeekChat | 84% | 53% | 0.20 |

## 4.3  Impact of Rephrasing

A key goal of the study was to determine whether LLM-based paraphrasing affects classifier performance. Results across the three rephrasing pipelines show that strict and relaxed accuracy on rephrased emails is nearly identical to that on the original dataset, with deviations within about 0.6 percentage points (maximum 0.55 pp strict and 0.54 pp relaxed). Moreover, differences between DeepSeek-, GPT-4o-, and RandomAPI-generated paraphrases were negligible, suggesting that surface-level linguistic variation has limited impact on Claude's classification decisions.

These findings indicate that Claude 3.5 Sonnet is robust under significant textual variation, maintaining stable performance even when emails are rewritten by heterogeneous LLMs. This is particularly encouraging for security settings, as attackers increasingly use paraphrasing—automated or manual—to evade detection.

### 4.3.1  Emotional and Motivational Labeling Benchmark

To choose a model for annotating emotional and motivational cues across the full dataset, we evaluated four LLMs (GPT-4o-mini, GPT-4.1-mini, Claude 3.5 Sonnet, DeepSeek-Chat) on a 100-email subset. Each model labeled every email five times, allowing us to measure (i) internal consistency across runs and (ii) agreement with human expert labels.

Table 3 summarizes the emotion-labeling results. All models showed high internal consistency (78–80%), but Claude 3.5 Sonnet achieved the strongest alignment with human annotations, with the highest close-enough accuracy (42%) and Jaccard similarity (0.60).

Motivation results (Table 4) show similar trends. Although strict accuracy was low across all models, close-enough accuracy remained between 53–61%, with Claude again achieving the highest consistency (88%).

Overall, Claude 3.5 demonstrates strong and stable phishing detection capabilities, consistently achieving around 66–67% strict accuracy

and 69% relaxed accuracy across all dataset variants. While distinguishing spam from valid emails remains a challenge, the model is highly reliable at separating benign from harmful content when evaluated in the more practical Unwanted vs. Valid setting. Crucially, classification performance remained nearly unchanged across all rephrased email groups, confirming the robustness of Claude 3.5 Sonnet against paraphrasing-based variations and demonstrating the reliability of the proposed dataset and evaluation framework.

## 5  LIMITATIONS & FUTURE RESEARCH

### 5.1  Limitations

This study has several limitations that must be acknowledged. First, the dataset is constrained by the limited sources from which the emails were collected. Many existing public phishing and spam datasets are outdated and often lack critical fields such as embedded URLs, file attachments, or sender information. These missing attributes reduce realism and the complexity of the classification task.

Second, the categorization of Spam emails presents inherent challenges. Unlike phishing, which has clearer malicious intent, the definition of spam is often subjective and context-dependent, varying between users and environments. This complicates ground-truth labeling and may influence evaluation outcomes.

### 5.2  Future Research

Future research coudl be focusing on enhancing both dataset realism and model robustness. A key direction involves comparing pre-LLM phishing emails with those generated or rephrased by modern LLMs to better understand how linguistic fluency and subtle persuasion techniques have evolved. Building on this, new adversarial

datasets could be created where hidden or obfuscated cues—such as embedded code, manipulated links, or metadata features invisible to human readers—are safely simulated to test model resilience.

Additionally, personalized and context-aware LLM-based email assistants may dynamically adapt to user-specific communication patterns, improving discrimination between legitimate and malicious messages. Training specialized models using both traditional machine learning and deep learning on enriched datasets could further strengthen detection capabilities. Expanding and modernizing the dataset with more diverse and up-to-date phishing, spam, and legitimate emails—including full metadata—could also be essential for improving generalizability.

Together, these improvments aim to advance the development of reliable, context-aware, and adaptive defenses against evolving phishing and spam threats.

# 6 CONCLUSION

In this study, we introduced a new, richly annotated dataset of phishing, spam, and legitimate emails, including emotional and motivational labels, rephrased variants, and model-generated samples. We benchmarked multiple LLMs for emotional and motivational analysis, selected the best-performing model, and conducted a large-scale classification experiment across original and rephrased emails. Below, we summarize the answers to our research questions.

- RQ1: How effectively can LLMs identify emotional strategies and patterns used within email content?

  Answer: LLMs can identify emotional cues with moderate reliability, but performance varies substantially across models. In our controlled benchmark, Claude 3.5 Sonnet achieved the highest alignment with human annotations, reaching a Jaccard similarity of 0.60 and a close-enough accuracy of 42%. Other models (e.g., GPT-4o-mini, DeepSeek-Chat) showed lower agreement. This indicates that while LLMs can capture emotional signals such as urgency, fear, and authority, human-level accuracy has not yet been reached.

- RQ2: How effectively can LLMs identify motivational strategies and patterns used within

email content?

Answer: Motivation labeling proved more challenging than emotion detection. Strict accuracy was low across all models due to the inherent ambiguity of attacker intent, but close-enough accuracy remained between 53–61% for the top-performing LLMs. Claude 3.5 Sonnet showed the highest consistency (88%) and strong recall, making it the most reliable model for large-scale annotation. Overall, LLMs can detect high-level motivations (e.g., link-clicking, credential theft), but fine-grained distinctions remain difficult.

- RQ3: How effectively can a state-of-the-art LLM perform email classification of phishing, spam, and legitimate emails across both human-written and LLM-generated content?

  Answer: Using Claude 3.5 Sonnet, we observed consistent classification performance across original and rephrased emails. Strict accuracy remained stable at approximately 66–67% across all groups, with high recall for phishing, but persistent confusion between spam and legitimate emails. Under the relaxed binary setting (Unwanted vs. Valid), performance increased to 69–70%, demonstrating that the model can reliably separate harmful from benign emails even under paraphrasing. This indicates robust detection of problematic emails, though fine differentiation between spam and phishing remains a challenge.

Overall, our findings show that contemporary LLMs have strong potential for automated email analysis, particularly in detecting harmful messages and identifying broad emotional and motivational patterns. However, challenges remain in achieving precise intent classification and fully matching human-level emotional interpretation. The dataset and evaluation framework developed in this work provide a foundation for future research into LLM-based email security, robustness, and explainability.

The dataset introduced in this study—combining human-written and LLM-generated emails with emotional, motivational, and rephrased variants—offers a compact yet comprehensive resource for studying text-based email threats. By releasing the dataset and tools openly, we aim to support reproducible research and enable more resilient and adaptable email-security models.

# REFERENCES

Afane, K., Wei, W., Mao, Y., Farooq, J., and Chen, J. (2024). Next-generation phishing: How llm agents empower cyber attackers.

Alam, N. A. (2024). Phishing email dataset. Accessed: March 25, 2025.

Apache Software Foundation (2025). SpamAssassin. https://spamassassin.apache.org/. [Online; accessed 31-Mar-2025].

Deng, Y., Zhang, W., Chen, Z., and Gu, Q. (2024). Rephrase and respond: Let large language models ask better questions for themselves.

Eilertsen, E., Mavroeidis, V., and Grov, G. (2025). Llm-powered intent-based categorization of phishing emails.

for Cybersecurity (ENISA), E. U. A. (2025). Enisa threat landscape 2025. Technical report, European Union Agency for Cybersecurity.

Gilardi, F., Alizadeh, M., and Kubli, M. (2023). Chatgpt outperforms crowd-workers for text-annotation tasks. Proceedings of the National Academy of Sciences, 120(30).

Guenter, B. (1998–2025). SPAM Archive. An archive of spam emails collected for research purposes, maintained by Bruce Guenter. Accessed 2025-06-10.

Heiding, F., Lermen, S., Kao, A., Schneier, B., and Vishwanath, A. (2024a). Evaluating large language models' capability to launch fully automated spear phishing campaigns: Validated on human subjects.

Heiding, F., Schneier, B., and Vishwanath, A. (2024b). Ai will increase the quantity — and quality — of phishing scams. Harvard Business Review.

Ige, T., Kiekintveld, C., Piplai, A., Waggler, A., Kolade, O., and Matti, B. H. (2024). An investigation into the performances of the current state-of-the-art naive bayes, non-bayesian and deep learning based classifier for phishing detection: A survey.

Josten, M. and Weis, T. (2024). Investigating the effectiveness of bayesian spam filters in detecting llm-modified spam emails. arXiv preprint arXiv:2408.14293.

Kim, J., Yang, N., and Jung, K. (2024). Persona is a double-edged sword: Mitigating the negative impact of role-playing prompts in zero-shot reasoning tasks.

Mardiansyah, K. and Surya, W. (2024). Comparative analysis of chatgpt-4 and google gemini for spam detection on the spamassassin public mail corpus. Preprint.

Martin, E. (2005). Warranted spam dataset. Accessed: 2025-05-14.

Niu, M. et al. (2024). Unveiling the emotion annotation capabilities of llms. In Interspeech 2024.

Petrosyan, A. (2024). Spam: Share of global e-mail traffic monthly 2014-2023. Statista.

Proofpoint (2025). 2025 state of the phish report. Technical report, Proofpoint, Inc.

Tarapiah, S. and Others (2025). Evaluating the effectiveness of large language models versus machine learning in identifying phishing email attempts. Algorithms, 18(10):599.

Tóth, R., Limonova, O., Voldokhin, S., and Belorusec, A. (2024). Impact of emotions on user behavior toward phishing emails. In Norsk IKT-konferanse for forskning og utdanning. Accessed: March 28, 2025.

Tusher, E. H., Ismail, M. A., Rahman, M. A., Alenezi, A. H., and Uddin, M. (2024). Email spam: A comprehensive review of optimize detection methods, challenges, and open research problems. IEEE Access, 12:143627–143657.

Wang, P. and Lutchkus, P. (2023). Psychological tactics of phishing emails. Issues in Information Systems, 24(2):71–83.

Zhang, J., Wu, P., London, J., and Tenney, D. (2025). Benchmarking and evaluating large language models in phishing detection for small and midsize enterprises: A comprehensive analysis. IEEE Access, 13:28335–28352.