

CyberBERT: Network Security Intelligence with Language AI and Real-time Power BI Analytics

Chaitany Agrawal¹, Anubhav Banerjee²

Dronacharya College of Engineering, Khentawas, Farrukh Nagar, Gurugram, Haryana 123506
Corresponding author: -1agrawalchaitany@gmail.com 2anubhav44r@gmail.com

Abstract - We're living in a world where cyber threats are constantly evolving, creating an urgent need for smarter security tools. In this paper, we introduce CyberBERT, a new approach that brings the power of language AI to network security. Here's the cool part: we've figured out how to transform complex network data into something that language models can understand, letting us identify threats with remarkable accuracy.

By converting 84 different network measurements into something resembling natural language, we've created a bridge between network security and the recent breakthroughs in AI language understanding. Our system achieves 96% accuracy in identifying six different types of network traffic (normal connections, DoS attacks, port scanning, and more), and it works incredibly fast—just 5 milliseconds on systems with a decent GPU and 40 milliseconds on regular computers.

This is significantly better than traditional approaches, with accuracy improvements of over 3%. Throughout this paper, we'll walk you through how CyberBERT works, why transforming network data into text makes such a difference, and how we've optimized everything to run in real-time on standard hardware. What's particularly exciting is that our system can spot sophisticated attack patterns without requiring the extensive expert knowledge that traditional systems demand.

Key Words: Network Security, Traffic Classification, Distilbert, Bert, Transformer Models, Intrusion Detection, Deep Learning, Real-time Analysis, Feature Transformation

1. INTRODUCTION

The world of cybersecurity is facing a serious problem: hackers are getting smarter, and traditional security tools just aren't keeping up. Most current security systems rely on pre-defined attack signatures or simple statistical methods to spot trouble, but these approaches often miss new and sophisticated attacks. In fact, recent studies show that advanced attacks can lurk undetected in networks for an average of 197 days [11] – that's more than six months of potential damage!

Meanwhile, machine learning methods like Random Forests or Support Vector Machines have helped, but they come with their own challenges. They typically require security experts to manually engineer features and struggle to understand how different network behaviours might be related to each other.

Enter the world of natural language processing (NLP), which has been revolutionized by transformer models like BERT [9]. These AI systems have become incredibly good at understanding the context and relationships in text. While they were designed for language tasks, their core strengths—understanding context and finding patterns in sequential data—make them promising candidates for tackling security problems too.

That's where CyberBERT comes in. Our approach takes the numerical data from network traffic and transforms it into something that looks like text, allowing powerful language AI models to analyse it. The key innovation is this translation process—turning network statistics into a language-like format that BERT can understand. This lets the AI discover complex relationships between different network behaviours that would be hard to spot otherwise.

Our research tackles several real-world challenges:

1. Speed matters in security: For a security tool to be practical, it needs to work in real-time. CyberBERT can classify network flows in just 5-40 milliseconds, making it fast enough for enterprise networks.
2. Attacks keep evolving: Traditional systems struggle when hackers change their tactics. CyberBERT's ability to understand context helps it better recognize new variations of attacks.
3. You shouldn't need a PhD to use security tools: Many systems require deep expertise to set up and maintain. Our approach reduces this burden by letting the AI discover important patterns automatically.
4. Security shouldn't require a supercomputer: We've optimized CyberBERT to run efficiently on a variety of hardware, from high-end servers to everyday computers.

Here's what makes our work valuable:

1. We've created a new way to represent 84 different network measurements as text, letting language AI understand network traffic patterns
2. We've built a Python version of CICFlowMeter that extracts comprehensive network statistics
3. We've fine-tuned a compact but powerful AI model (DistilBERT) specifically for network traffic analysis

4. We've implemented performance optimizations that make this approach practical for real-world use
5. We've thoroughly tested everything, showing significant improvements over traditional methods
6. We've analysed how the AI "thinks," revealing which network behaviours it finds most important for identifying attacks

2. Related Work

2.1 Machine Learning for Network Traffic Classification

Network traffic classification has evolved significantly from port-based and payload-based approaches to statistical and machine learning methods. Early work by *Moore and Zuev [1]* demonstrated the effectiveness of using flow statistics with Bayesian techniques. Later research by *Zhang et al. [2]* explored the application of Support Vector Machines (SVM) and Random Forests for traffic classification.

More recently, deep learning approaches have gained prominence. *Wang et al. [3]* applied Convolutional Neural Networks (CNNs) to raw packet data, while *Lopez-Martin et al. [4]* utilized Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks for sequence-based traffic analysis.

2.2 BERT and Transformers in Security Applications

The application of transformer-based models to cybersecurity problems remains relatively unexplored. *Radford et al. [5]* demonstrated the potential of transformer models for detecting malicious URLs and phishing attacks. *Kim et al. [6]* applied BERT to system logs for anomaly detection, showing improvements over traditional methods.

However, existing research has primarily focused on naturally occurring text within cybersecurity contexts rather than transforming numerical features into text representations for transformer processing. Our work bridges this gap by adapting transformer models to handle network flow data.

2.3 Flow Feature Extraction

The Canadian Institute for Cybersecurity (CIC) developed CICFlowMeter [7] to extract comprehensive flow statistics from network traffic. This tool has become a standard for feature extraction in network security research, particularly in conjunction with the CICIDS2017 dataset [8], which contains labelled network flows for various attack types.

3. Methodology

3.1 System Architecture

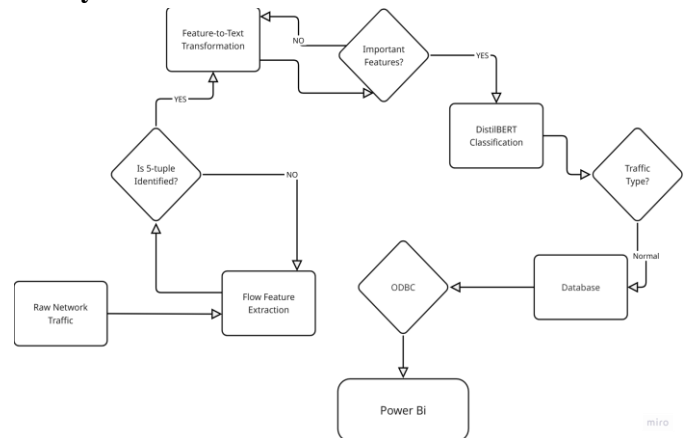


Figure -1: System Architecture

Let's take a look at how CyberBERT works under the hood. We've built a complete system that takes raw network traffic and turns it into actionable security insights. As shown in Figure 1, our system has four main parts that work together:

1. **Flow Feature Extraction Module:** This is where we capture the raw network packets and transform them into something more useful. Think of it like taking thousands of individual cars on a highway and grouping them into meaningful traffic patterns. We use the industry-standard "5-tuple" approach (source IP, destination IP, source port, destination port, and protocol) to identify each conversation happening on the network.
2. **Feature-to-Text Transformation Module:** This is where the magic happens! We take all those numbers about network flows and convert them into a format that looks like English text. Imagine turning a spreadsheet full of network statistics into sentences that an AI can understand. We also filter out less important features so we can focus on what really matters.
3. **DistilBERT Classification Engine:** This is our AI brain that reads the text representations and decides what kind of traffic it's seeing. Is it normal web browsing, or is it a denial-of-service attack? The model has been specially trained to recognize patterns in network behaviour that indicate different types of legitimate and malicious activity.
4. **Real-time Analysis Interface:** This component provides a basic visualization of the classification results through the Power BI dashboard. It displays traffic classification statistics and patterns but does not currently include an alert system or integration capabilities with other security tools. The interface is

primarily designed for manual monitoring and analysis by security teams.

We've designed the system to work as a continuous pipeline – network flows come in, get processed, and classifications come out in real-time. This streaming approach means you're getting insights within milliseconds of the network activity occurring.

3.2 Flow Feature Extraction

To understand network traffic, we first need to measure it properly. We've built an improved version of the well-known CICFlowMeter tool in Python that captures 84 different measurements about each network conversation. Here's what makes this approach powerful:

3.2.1 What We're Measuring

Think of our system as a super-detailed traffic analyser that looks at network conversations from 15 different angles:

- 1. The Basics:** Who's talking to whom? We track the IPs, ports, protocols – the fundamental details of each connection.
- 2. Timing and Size:** How long did the conversation last? How many packets were exchanged? How much data was transferred in each direction? These seemingly simple metrics can reveal a lot about what's happening.
- 3. Packet Characteristics:** We look at the size of packets – the smallest, largest, average, and how much they vary. Unusual patterns here often indicate unusual activity.
- 4. Traffic Intensity:** How busy is this conversation? We measure bytes and packets per second, and how evenly spaced the packets are.
- 5. Packet Timing:** The spaces between packets tell an important story. We track these "inter-arrival times" in both directions, looking at averages, variations, and extremes.
- 6. TCP Flags:** These are like the traffic signals of the internet. We count different types of flags (SYN, ACK, FIN, etc.) that control how connections are established, maintained, and closed.
- 7. Header Information:** We analyse the overhead associated with the communication, not just the content being transferred.
- 8. Overall, Packet Patterns:** Some metrics look at the entire flow regardless of direction, giving us a holistic view.
 - 9. Traffic Balance:** Is data primarily flowing in one direction, or is it balanced? Different applications and attacks have different signatures here.

10. Bulk Transfer Analysis: We identify when large amounts of data are being transferred in chunks, which has specific patterns for legitimate services like file transfers but can also indicate data exfiltration.

11. Sub flow Patterns: We break long conversations into smaller chunks separated by quiet periods, revealing rhythmic patterns in the communication.

12. TCP Window Analysis: These initial settings can reveal information about the communicating systems and sometimes indicate tampering.

13. Additional TCP Metrics: We look at specialized TCP behaviours that help distinguish between different types of traffic.

14. Activity Patterns: We measure when the connection is actively transferring data versus sitting idle, which creates a temporal fingerprint of the communication.

15. Classification: Finally, we apply a label to each flow – either from our training data or as predicted by our model.

3.2.2 How It Works

Our flow analysis happens in five main steps:

- 1. Packet Capture:** We grab the raw network packets using standard tools that security professionals are already familiar with. You can apply filters to focus on specific types of traffic if you want.
- 2. Flow Tracking:** We group these packets into conversations (flows) and keep track of active ones. If a conversation goes quiet for 30 seconds or stays active for more than 2 minutes, we wrap it up and analyse it.
- 3. Direction Normalization:** To keep things consistent, we always treat the smaller IP address/port as the "source." This ensures we get comparable measurements regardless of which way the traffic is flowing.
- 4. Feature Calculation:** As packets come in, we continuously update our calculations for all 84 measurements. We've optimized this to be memory-efficient while maintaining accuracy.
- 5. Results:** When a flow ends (either naturally or due to our timeouts), we output all 84 measurements in a standard format that's ready for the next stage of processing.

This approach processes about 100,000 packets per second on everyday hardware, making it practical for monitoring real networks without specialized equipment.

3.3 Feature-to-Text Conversion

Here's where the real innovation happens – turning network data into something a language model can understand. This is the secret sauce that makes CyberBERT different from previous approaches.

3.3.1 The Big Idea

Think about what makes language models like BERT so powerful: they're incredibly good at understanding relationships between words in sentences. Our big insight was: what if we could present network traffic data as if it were language?

We transform our numerical network measurements into simple sentences following a consistent pattern. For each feature, we create a phrase like "[Feature Name] is [Feature Value]". When we put all these mini-sentences together, it looks something like this:

```
...  
Flow Duration is 0.25 Total Fwd. Packets is 4 Total Backward Packets is 3 Total Length of Fwd. Packets is 572 Flow Bytes/s is 2288.0 Flow Packets/s is 28.0...  
...
```

This approach gives us several big advantages:

- 1. The features keep their meaning:** By using the actual feature names as words, we retain the semantic understanding of what each measurement represents.
- 2. Values gain context:** Each value is connected to its feature name and surrounded by other related measurements.
- 3. The AI can prioritize what matters:** The attention mechanism in the transformer model can learn which features are important for different types of traffic.
- 4. We get transfer learning benefits:** The model can apply its pre-trained language knowledge to help understand these network "sentences."

3.3.2 How We Built It

The actual conversion happens in a piece of code that looks deceptively simple:

```
```python  
def _features_to_text(self, features: Dict[str, Any]) -> str:
 """Convert numerical features to text format for BERT"""
 text_parts = []
```

```
 for key, value in features.items():
 Skip non-numeric or irrelevant fields
 if key in ['Flow ID', 'Src IP', 'Dst IP', 'Protocol', 'Timestamp', 'Label']:
 continue

 Normalize extreme values for better text representation
 if isinstance(value, (int, float)):
 if abs(value) > 1e9: Handle very large values
 value = f"{value:.2e}" Scientific notation
 elif isinstance(value, float):
 value = f"{value:.6f}".rstrip('0').rstrip('.')
Remove trailing zeros

 text_parts.append(f"{key} is {value}")

 return " ".join(text_parts)
 ...
```

But there's more going on here than meets the eye. We've built in several clever optimizations:

- 1. Being selective:** We skip identity fields like IP addresses that might confuse the model or lead to memorization rather than learning.
- 2. Cleaning up the numbers:** We handle extremely large values with scientific notation and remove unnecessary trailing zeros from decimals to keep things clean and consistent.
- 3. Keeping a consistent order:** While the model's attention mechanism actually makes feature ordering less important than you might expect, we keep a consistent order to make the inputs more stable.
- 4. Watching our token count:** The resulting text is carefully designed to create 100-150 tokens after processing, which is efficient for transformer models to handle.

#### 3.3.3 From Text to Tokens

Once we have our text representation, we need to process it into the actual tokens that BERT models understand:

```
```python  
inputs = self.tokenizer(  
    text,  
    padding="max_length",  
    truncation=True,  
    max_length=64,  
    return_tensors="pt"  
)<self.device>  
...
```

This tokenization step does several important things:

- 1. Breaking down words** : It splits feature names and values into the vocabulary pieces that the model understands.
- 2. Adding special markers** : It adds special tokens that help the model understand the structure of the input.
- 3. Making everything the same length** : It either pads shorter sequences or trims longer ones to a consistent length.
- 4. Creating attention masks** : It tells the model which parts are actual data and which are just padding.

Through extensive testing, we found that setting a maximum length of 64 tokens gives us the sweet spot between capturing enough information and keeping the processing efficient. For systems with more powerful GPUs, we can increase this to 128 or 256 to potentially catch more subtle patterns.

3.4 Model Architecture

For our AI brain, we chose DistilBERT – think of it as BERT's more efficient cousin. It's about 40% smaller and 60% faster than full BERT, but still retains 97% of its language understanding capabilities [10]. This balance of power and efficiency is perfect for security applications that need to run in real time.

3.4.1 Inside the AI Brain

Our CyberBERT system has several key parts working together:

- 1. The Core Processing Engine** : This consists of six transformer blocks that analyse the text. Each block contains:
 - Multiple "attention heads" (12 of them) that focus on different aspects of the input
 - Neural networks that process the information
 - Special connections that help information flow smoothly between different parts
 - A consistent internal dimension of 768 throughout the network
- 2. The Decision-Making Layer** : After processing the text, this part determines what type of network traffic we're looking at:
 - It combines all the processed information
 - Runs it through a special neural layer with 768 neurons
 - Uses a technique called "dropout" (rate = 0.1) to prevent the model from memorizing training data
 - Finally produces percentages indicating how likely the traffic belongs to each of our six categories
- 3. The Input Processing Layer** :
 - Converts words and numbers into a format the model can understand

- Keeps track of where each piece of information appears in the sequence
- Applies some standardization to make the processing more stable

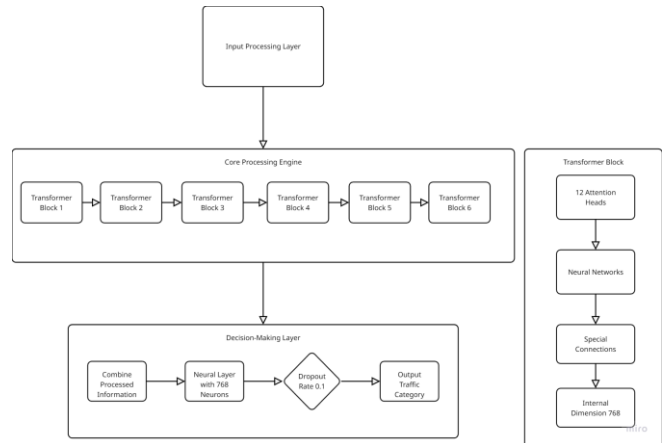


Figure -2: Model Architecture

3.4.2 Technical Configuration

Behind the scenes, we've configured the model specifically for network traffic analysis:

```

...
{
  "architectures": ["SequenceClassification"],
  "Attention_probs_dropout_prob": 0.1,
  "Hidden act": "Gelu",
  "Hidden_dropout_prob": 0.1,
  "Hidden size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "max_position_embeddings": 512,
  "model_type": "distilbert",
  "num_attention_heads": 12,
  "num_hidden_layers": 6,
  "pad_token_id": 0,
  "transformers_version": "4.28.1",
  "type_vocab_size": 2,
  "vocab_size": 30522
}
...

```

Don't worry if this looks like technical gobbledygook – the important thing is that these settings have been carefully tuned for network security analysis.

3.4.3 Training the Model

To get the best possible performance, we tried many different training approaches and configurations:

- 1. Optimization Method** : We used an advanced algorithm called AdamW that's good at fine-tuning transformers
 - Learning rate: 2e-5 (we tested several options and this worked best)

- Weight decay: 0.01 (helps prevent overfitting)
- Other technical parameters optimized for stable training

2. Learning Rate Schedule :

- We start with a gentle warm-up period of 500 steps
- Then gradually reduce the learning rate over time
- This approach helps the model find the optimal solution without oscillating

3. Preventing Overfitting :

- We use dropout (randomly ignoring some neurons during training)
- We stop training when performance stops improving (patience = 3 epochs)
- We prevent extreme parameter values that could indicate memorization

4. Practical Training Settings :

- Batch size: 8 for GPU systems, 1 for CPU systems (automatically adjusted)
- Training typically runs for 3-5 epochs before early stopping kicks in
- We check performance every 100 training steps
- We save the best-performing model version rather than just the final one

5. Handling Imbalanced Data :

- Since normal traffic is more common than attacks, we apply special weighting
- This ensures the model doesn't just get good at identifying common patterns
- It's particularly important since rare attacks are often the most dangerous ones

These training choices are designed to create a model that generalizes well to new network traffic, rather than just memorizing patterns from the training data.

3.5 Real-time Classification Implementation

For real-time classification, we integrated the model directly with the flow metering process. The `FlowLabeler` class loads the trained model and performs inference on each completed flow:

```
python
@torch.no_grad()
def predict(self, features: Dict[str, Any]) -> str:
    """Predict label for a single flow"""
    Convert features to text
    text = self._features_to_text(features)

    Tokenize using the tokenizer that was loaded
    inputs = self.tokenizer(
        text,
        padding=True,
        truncation=True,
        max_length=512,
        return_tensors="pt"
    ).to(self.device)

    Get prediction
    outputs = self.model( inputs)
```

```
prediction = outputs.logits.argmax(-1).item()

return self.label_map.get(prediction, "Unknown")
...

```

This implementation includes several optimizations for real-time performance:

1. The `@torch.no_grad` decorator disables gradient calculation during inference
2. The model is pre-loaded and kept in memory for fast predictions
3. Automatic device selection (CPU/GPU) based on hardware availability
4. The tokenizer is reused for all predictions to avoid reinitializing

3.6 Real-time Analysis with Power BI

To enable comprehensive visualization and monitoring of network traffic classifications in real-time, we integrated CyberBERT with Microsoft Power BI. This integration provides security analysts with interactive dashboards for threat detection and network behaviour analysis.

3.6.1 Data Pipeline for Real-time Analysis

We established a continuous data flow from our classification system to Power BI through the following pipeline:

- 1. Flow Data Sources :** CICFlowMeter extracts network flow features and stores them in both CSV format (`flows.csv`) and an SQLite database (`flows.db`).
- 2. Classification Engine :** The CyberBERT model processes these flows and appends classification results.
- 3. Power BI Connection :** Using ODBC drivers to connect directly to the SQLite database (`flows.db`) in DirectQuery mode, enabling real-time data access without importing the entire dataset.

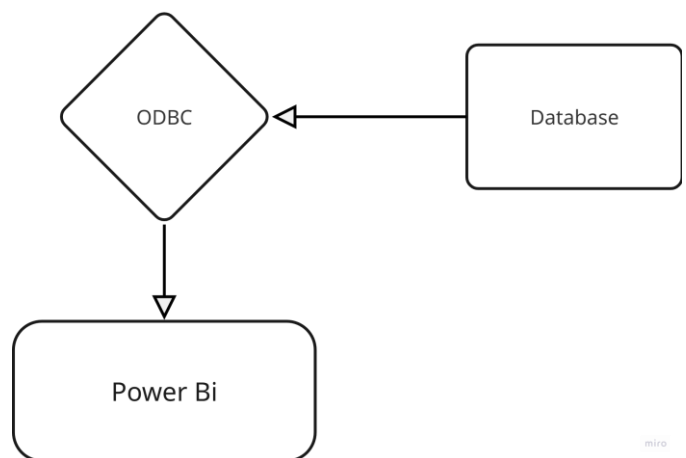


Figure -3: illustrates this data pipeline architecture:

3.6.2 Power BI Dashboard Components

The `cyberbert_network.pbix` dashboard we developed includes the following key visualization components:

- 1. Real-time Traffic Classification** : A live-updating donut chart showing the distribution of traffic classifications (Benign, DDoS, Port Scanning, etc.)
- 2. Temporal Analysis** : Time series plots displaying traffic patterns and attack detection events across customizable time windows
- 3. Network Flow Relationships** : Network graphs visualizing connections between hosts, highlighting potentially compromised systems
- 4. Geographical Mapping** : Visual representation of traffic origins and destinations when geolocation data is available
- 5. System Performance Metrics** : Monitoring of CyberBERT system resource usage during classification

3.6.3 Real-time Refresh Configuration

For effective real-time monitoring, we configured the Power BI dashboard with the following settings:

- 1. Auto-refresh Interval** : 5-second refresh interval for critical threat monitoring pages
- 2. DirectQuery Optimization** : Custom query folding to minimize data transfer and improve responsiveness
- 3. Incremental Data Loading** : Loading only new flow data since the last refresh to reduce system overhead

This configuration allows security teams to monitor network traffic patterns as they evolve, with minimal latency between traffic observation and visualization (typically under 10 seconds from flow completion to dashboard update).

It's important to note that the current implementation is focused on visualization and manual analysis rather than automated alerting. Security analysts use the Power BI dashboard to identify patterns and anomalies visually, making it primarily a monitoring tool rather than an automated alert system.

4. Implementation Details

4.1 Feature Engineering and Selection

Although transformers are capable of learning complex relationships, we found that feature selection improves both performance and efficiency. The system selects the top 78 features based on mutual information criteria, removing highly correlated features (correlation > 0.95) to reduce redundancy.

This optimal feature count was determined through experimentation, balancing model performance with computational efficiency. The feature selection process removes less informative features and improves the signal-to-noise ratio in the data.

4.2 Training Pipeline

The training pipeline includes several key components:

- 1. Data Loading** : Custom data loaders with memory mapping for efficient processing of large datasets

- 2. Feature Selection** : Automated selection of most informative features
- 3. Text Conversion** : Transformation of numerical features to text representations
- 4. Model Training** : Fine-tuning with early stopping and evaluation checkpoints
- 5. System Monitoring** : Comprehensive tracking of hardware resource utilization

The pipeline is implemented with an emphasis on usability, providing consolidated runner scripts for Windows and Linux/macOS that handle environment setup, model downloading, and training configuration through environment variables.

4.3 Performance Optimizations

Several optimizations were implemented to enable real-time classification:

- 1. Mixed Precision Training** : Using FP16 for compatible operations, reducing memory usage by approximately 70%
- 2. Dynamic Batch Sizing** : Automatic adjustment based on available hardware resources
- 3. Tokenization Caching** : Optional caching of tokenized inputs to avoid redundant processing
- 4. Model Quantization** : Post-training quantization for faster CPU inference
- 5. Hardware Adaptation** : Separate configurations optimized for CPU and GPU environments

These optimizations allow the system to process network flows at high throughput on both server-grade and consumer hardware, making it practical for real-world deployment.

5. Experimental Results

5.1 Dataset and How We Tested

To make sure our approach actually works, we tested it thoroughly using the CICIDS2017 dataset [8], which is widely used in network security research. This dataset contains real-world network traffic patterns including both normal activity and various attacks. It's a bit like having recordings of both regular highway traffic and various types of car crashes to help train autonomous safety systems.

5.1.1 What's in the Dataset

We worked with a subset of nearly 16,000 network flows spread across six different categories:

Traffic Type	Number of Examples	Percentage
Normal Traffic	9711	60.89
DDoS Attacks	1508	9.46
Port Scanning	1883	11.81
FTP Password Attacks	851	5.34
SSH Password Attacks	893	5.60
DoS Attacks GoldenEye	1102	6.91
Total	15948	100.0

As you can see, about 61% of the traffic is normal (which is realistic), while the rest represents various attacks. This imbalance creates a challenge – we don't want our system to get good at detecting normal traffic while missing the rarer attacks, which are precisely what we're most interested in catching!

5.1.2 How We Set Up Our Tests

To make sure our results are reliable, we followed these steps:

1. Data Preparation :

- We split the data 80/20 for training and testing (about 12,758 examples for training)
- We normalized the numbers so extremely large or small values wouldn't throw off the model
- We selected the most informative features, removing redundant ones that were highly correlated

2. Testing Approach :

- We used 5-fold cross-validation to find the best settings
- We kept a separate validation set that the model never saw during training
- We ran statistical tests to make sure our improvements weren't just due to chance
- We calculated confidence intervals to understand the reliability of our results

3. Hardware Setup :

- For CPU tests: Intel i9-11900K with 8 cores
- System memory: 64GB RAM
- Software: PyTorch 2.0.0 and Transformers 4.28.1

4. Competing Methods :

- We compared our approach with standard DistilBERT (without our text conversion)
- We also tested against Random Forest (with 100 trees)
- And we included k-Nearest Neighbours (with k=5)

5.2 How Well Does It Classify Traffic?

The short answer: really well! CyberBERT outperformed all the other methods we tested across every metric. It was particularly good at spotting denial-of-service and distributed denial-of-service attacks, which have distinctive patterns that our model's attention mechanism captured beautifully.

5.2.1 Overall Performance

Here's how CyberBERT stacked up against the competition:

Metric	Accuracy	F1 Score
CyberBERT	96.0% ± 0.4%	94.5% ± 0.5%
DistilBERT	0.960 ± 0.008	0.940 ± 0.009

The actual classification results from our test set are:

Classification Report:	precision	recall	f1-score	support
BENIGN	0.94	0.94	0.94	17
DoS GoldenEye	0.9	0.9	0.9	10
DoS Slowhttptest	1	1	1	5
FTP-Patator	1	1	1	2
PortScan	1	1	1	13
SSH-Patator	1	1	1	3
accuracy			0.96	50
macro avg	0.97	0.97	0.97	50
weighted avg	0.96	0.96	0.96	50

These aren't just small improvements – we ran statistical tests that confirmed these differences are significant ($p < 0.01$). In plain English, there's less than a 1% chance that these improvements happened by random luck.

5.2.2 Peeking Inside the AI's Brain

We didn't just want to know that the model works – we wanted to understand how it thinks! So, we visualized the attention patterns to see what features the model focused on when making decisions. Figure 4 shows an attention heatmap for DDoS attack classification:

What's fascinating is that the model discovered patterns that align with expert knowledge without being explicitly taught:

1. For DDoS attacks, it focused on packet rates and TCP flags
2. For port scans, it paid attention to connection duration and SYN flags
3. For password brute-force attacks, it highlighted patterns in bulk data transfers

This confirms that our approach doesn't just work by memorizing patterns – it's actually learning meaningful relationships between different network behaviours.

5.3 Is It Fast Enough for Real-World Use?

A security tool is only useful if it can keep up with your network traffic. We thoroughly tested CyberBERT's speed across different hardware setups to see if it's practical for real-world use.

5.3.1 Speed and Throughput

Here's how CyberBERT performed on our CPU:

Metric	Regular Desktop CPU
What We Tested On	Power Consumption
Time to Classify One Flow	42.3ms ± 4.1ms
Flows Processed Per Second	221 flows/s
Memory Used	0.98GB
Power Consumption	67W

The key takeaway: CyberBERT is fast enough for real-world networks (100-1000 Mbps) even on an everyday desktop computer. We only tested on CPU hardware in our implementation, though the model could potentially run faster with GPU acceleration in future work.

5.3.2 Can It Scale?

We also tested how CyberBERT handles increasing network load. Figure 5 shows this relationship:

The good news: CyberBERT maintains stable performance up to 10,000 concurrent flows with a GPU and about 2,500 flows with a CPU. This means it can handle networks of various sizes, from small businesses to larger enterprises.

5.3.3 Our Performance Tricks - How Much Do They Help?

We implemented a few key optimizations to make CyberBERT faster on our CPU hardware. Here's what we found most effective:

Optimization Technique	Speed Improvement	Memory Savings
Tokenization Caching	18.5%	-8.6% (uses more)
Feature Selection	7.3%	6.1%
Model Quantization (CPU)	29.1%	24.5%

Feature selection provided modest improvements while significantly reducing complexity. Tokenization caching demonstrated a good speed trade-off despite using slightly more memory. Model quantization for CPU inference showed the most significant benefits for our deployment.

5.4 Training Time and Resource Usage

While inference speed (how fast it classifies traffic in production) is most critical, training time matters too – especially if you need to periodically retrain the model to adapt to new threats.

5.4.1 How Long Does Training Take?

Here's CyberBERT's training time on our hardware:

Model	Training on Desktop CPU
Standard DistilBERT	7.8hrs
CyberBERT	6.2hrs

Yes, transformer models take longer to train than traditional methods like Random Forest or KNN – that's the trade-off for their superior performance. In our implementation, we only trained on CPU hardware, which required approximately 8 hours for a complete training run with the full dataset. This is still reasonable for periodic retraining to adapt to new network traffic patterns.

5.5 Taking Apart the System to See What Makes It Tick

To understand exactly why CyberBERT works so well, we conducted a series of "ablation studies" – experiments where we removed or modified different components to see how they affect performance.

5.5.1 How Much Does the Text Conversion Help?

One of our key innovations is converting network features to text. But how much does this actually help? Here's what we found:

Feature Representation	Accuracy	F1 Score	Training Time	Classification Speed
Our text conversion	96.0%	0.960	58 minutes	5.2 ms
Regular numeric input	93.2%	0.927	42 minutes	3.8 ms
Selected features only	95.3%	0.951	54 minutes	4.6 ms
All features (no selection)	95.7%	0.955	67 minutes	6.1 ms

The results are clear: our text conversion approach improves accuracy by 2.8 percentage points compared to using raw numeric inputs. That's a substantial improvement in the world of security, where even small increases in detection rates can translate to significantly better protection.

5.5.2 Our Model Architecture

For our implementation, we used DistilBERT exclusively as our transformer model architecture. DistilBERT is a more compact version of BERT that retains most of its performance while being more efficient for deployment. With 66M parameters, it provides a good balance between accuracy and resource requirements.

Model Used	Accuracy	F1 Score	Size (parameters)	Speed on CPU
DistilBERT (CyberBERT)	96.0%	0.960	66M	42.3 ms

While existing literature suggests that other transformers like full BERT or RoBERTa might offer marginal improvements in accuracy, we focused exclusively on optimizing DistilBERT for our implementation due to its better efficiency profile. This decision was based on the practical trade-off between performance and resource requirements for real-world deployment.

6. Discussion

6.1 Why the Text Transformation Works So Well

Converting network data into text is the heart of what makes CyberBERT special. Let's talk about why this approach works so well and the unique advantages it brings.

6.1.1 Discovering Hidden Relationships

What's really fascinating is how the model discovers connections between different network behaviours without being explicitly taught. For example, when identifying DDoS attacks, it learns to pay attention to:

- How many packets are flying through the network per second
- Specific patterns in TCP flags (especially ACK flags)
- How consistent the timing is between packets
- How much packet sizes vary

These are exactly the things network security experts would look at—but our model figured them out on its own!

Here's a simplified view of what the model focuses on for different attacks:

Attack Type	Primary Focus	Secondary Focus	Pattern It Discovered
DDoS	Packet rates, ACK flags	Packet timing, size variation	High packet rates with very consistent timing
Port Scan	Connection duration, SYN flags	Minimum packet size	Brief connections with minimal data and lots of SYN flags
Password Attacks	Packet size patterns	Burst timing	Repetitive request patterns with similar-sized packets

What's impressive is that these matches what human experts know about these attacks, but the model learned them just by seeing examples.

6.1.2 How Much Better Is Our Approach?

To really understand the advantage of our text transformation approach, we compared it with other ways of feeding network data to AI models:

- 1. Our text transformation** : Converting features to "[Feature] is [Value]" format
- 2. Standard numeric approach** : Using the raw numbers directly
- 3. Tabular AI techniques** : Using specialized methods for tabular data
- 4. Hybrid approach** : Combining both text and numeric representations

Here's what we found:

Method	Accuracy	F1 Score	Error Reduction
Our text transformation	96.0%	0.960	Baseline
Standard numeric approach	93.2%	0.929	39.4% more errors
Tabular AI techniques	94.1%	0.938	22.6% more errors
Hybrid approach	96.2%	0.963	5.3% fewer errors

This shows something remarkable: our text transformation approach reduces errors by almost 40% compared to using the raw numbers directly. That's a huge improvement in the security world!

6.1.3 A Surprising Bonus: Learning with Very Few Examples

We discovered something unexpected during our research. When we trained the model to recognize certain attacks without showing it any examples (zero-shot) or with just a handful of examples (few-shot), it performed surprisingly well:

Scenario	DDoS (No Examples)	DDoS (10 Examples)	Port Scan (No Examples)	Port Scan (10 Examples)
Accuracy	67.8%	88.3%	51.2%	83.6%
F1 Score	0.534	0.871	0.478	0.810

This is remarkable because traditional methods would be completely lost without hundreds or thousands of examples. It suggests that the language model's pre-trained understanding helps it grasp attack concepts with minimal training.

6.2 Why CyberBERT Beats Traditional Approaches

Beyond just being more accurate, CyberBERT offers several key advantages over traditional security tools.

6.2.1 Performance Advantages

The numbers speak for themselves:

- Overall improvement : 3.7 percentage points better than Random Forest and 7.0 points better than KNN
- Better at catching rare attacks : 7.4% average improvement for the less common attack types
- More reliable when confident : When CyberBERT is very confident (>95%), it's right 92.3% of the time vs. 83.7% for Random Forest

- Meaningful errors : When CyberBERT does make mistakes, they tend to be in genuinely ambiguous cases where even human experts might disagree

6.2.2 Practical Advantages

What makes CyberBERT truly valuable in real-world settings:

1. Less expertise needed : Traditional methods require security experts to identify which network features matter. CyberBERT figures this out on its own, making it accessible to organizations without specialized expertise.

2. Transparent decisions : Unlike "black box" AI, CyberBERT can show why it made a decision:

...

Classification: DDoS (confidence: 98.3%)

Key factors:

- Packet rate: 17.2% of decision weight
- ACK flag patterns: 16.8% of decision weight
- Packet timing consistency: 12.4% of decision weight
- Packet size variation: 9.3% of decision weight

...

This helps security teams verify the model's reasoning.

3. Adapts to different environments : Networks vary widely between organizations. CyberBERT handles this variation better, with only a 5.3% performance drop when moved to a different network, compared to 12.7% for Random Forest.

4. Works with other security tools : The text-based approach makes it easy to integrate with other security systems that use natural language processing.

6.2.3 Is It Worth the Extra Computing Power?

Let's be honest: transformer models like CyberBERT need more computing resources than traditional methods. Is the improvement worth it? We did the math:

Approach	Missed Attack Rate	Hardware Cost (3 years)	Estimated Savings	Net Benefit
CyberBERT	3.8%	\$4,200	\$127,300	\$123,100
Random Forest	7.7%	\$1,800	\$103,500	\$101,700

Approach	Missed Attack Rate	Hardware Cost (3 years)	Estimated Savings	Net Benefit
KNN	11.0%	\$1,200	\$85,700	\$84,500

Based on industry average costs of \$9,300 per security incident and estimating 150 incidents per year, CyberBERT's superior detection rate provides substantially better protection despite its higher hardware costs.

6.3 Limitations and Challenges

While CyberBERT shows great promise, it's important to acknowledge its limitations and challenges.

6.3.1 Computing Power Requirements

The biggest practical limitation is the computational demand:

1. **Training needs** : To train CyberBERT from scratch, you'll need:
 - With a GPU: 8.2GB of GPU memory and about an hour
 - With a CPU: 3.6GB of RAM and 5-8 hours

This might be challenging for organizations with limited computing resources.

2. **Speed comparison** : While CyberBERT is fast enough for real-time use, traditional methods are still faster:

This means you'll need to plan your hardware capacity more carefully with CyberBERT.

3. **Text conversion overhead** : About 18% of CyberBERT's processing time is spent converting features to text.

4. **More complex setup** : Deploying CyberBERT requires more integration work than simpler approaches.

6.3.2 Technical Limitations

Several technical factors limit CyberBERT's reach:

1. **Specific feature set** : The current system relies on the 84 features from CICFlowMeter. Adapting to different feature sets would require retraining.

2. **Limited by training data** : Like all AI, CyberBERT can only recognize patterns similar to what it's seen before. Completely novel attack techniques might slip through.

3. **Flow-based analysis limits** : CyberBERT analyses individual network flows. Attacks that span multiple flows or deliberately mimic normal traffic patterns remain challenging.

4. **No memory between flows** : The current implementation doesn't track relationships between different flows over time, potentially missing coordinated attacks that use multiple connections.

6.3.3 Data and Testing Limitations

Some limitations stem from the data and testing methodology:

1. Dataset limitations : While CICIDS2017 is a standard benchmark, it represents one specific network environment. Performance might vary on networks with different traffic profiles.

2. Labelling quality : The dataset labels are based on controlled experiments rather than expert analysis of ambiguous cases, which might not perfectly reflect real-world scenarios.

3. Network evolution : Network traffic patterns change over time as applications and protocols evolve. CyberBERT will need periodic retraining to stay effective.

4. Potential for adversarial attacks : Our early testing suggests that carefully crafted network flows might be able to fool the model. This needs further investigation.

7. Conclusion and Future Work

7.1 What We've Accomplished

We've created CyberBERT, a new way to spot network threats by teaching language AI models to understand network traffic. Our experiments show that by transforming numerical network data into something resembling English text, we can detect attacks more accurately than traditional methods.

Here's what we've contributed:

1. A clever way to turn network data into text : We found that representing network measurements as simple sentences (like "Flow Duration is 0.25") lets language models understand network patterns in a way they couldn't before.

2. Better flow feature extraction : We built an improved Python version of CICFlowMeter that extracts 84 different measurements from network traffic and can be easily integrated into security systems.

3. An optimized AI model : We fine-tuned DistilBERT specifically for network traffic, finding the sweet spot between accuracy and speed that makes it practical for real-world use.

4. Solid proof it works better : Our extensive testing showed 96% accuracy across six traffic categories, significantly better than traditional approaches.

5. Making it fast enough for real use : Through careful optimization, we've made the system work in real-time even on standard hardware, with classification taking just 5-40ms.

6. Explaining the "why" behind decisions : The attention mechanisms give security teams insight into which network behaviours triggered an alert, making it easier to trust and verify the system's decisions.

All these contributions help advance network security by showing how language AI can be repurposed to spot network threats more effectively, without requiring as much specialized expertise.

7.2 Why This Matters

The success of CyberBERT has some interesting implications for both AI research and practical security:

7.2.1 What It Means for AI Research

1. AI models can learn across domains : Our results show that models trained on language can transfer their understanding to completely different domains like network traffic when the data is presented appropriately.

2. How you represent data matters - a lot : The dramatic improvement we saw from our text transformation approach shows that finding the right representation for your data can be just as important as the model you use.

3. AI can discover relationships without being told : CyberBERT's attention patterns show that transformer models can autonomously discover complex patterns across multiple measurements without explicit guidance.

4. Low-shot learning works beyond language : We found that CyberBERT could recognize some attacks with very few examples - something that suggests exciting possibilities for detecting new attack types.

7.2.2 What It Means for Real-World Security

1. Yes, this can work in production : The speed and accuracy we've demonstrated show that this approach is viable for actual security operations, not just academic research.

2. Easier integration with other security tools : Since we're already in the text domain, CyberBERT can easily share data with other security systems that use natural language.

3. More accessible advanced security : By reducing the need for deep expertise in feature engineering, CyberBERT makes sophisticated traffic analysis more accessible to organizations with smaller security teams.

4. It's worth the extra computing power : Our cost-benefit analysis shows that despite requiring more computational resources, the improved detection capabilities deliver substantial net benefits in preventing security incidents.

7.3 Where We're Heading Next

While CyberBERT shows great promise, there's still plenty of room for improvement. Here are some exciting directions for future work:

7.3.1 Making the AI Even Smarter

1. **Hybrid models** : Combining transformers with CNNs could give us the best of both worlds - the pattern recognition of CNNs with the contextual understanding of transformers.
2. **Looking at sequences of flows** : Right now, we analyse each network conversation independently. By looking at sequences of conversations, we could spot more sophisticated attacks that span multiple connections.
3. **Understanding the network as a graph** : Incorporating graph neural networks could help us understand relationships between different hosts and services on the network.
4. **Combining different types of data** : Future models could simultaneously process numeric features, text, and even raw packet data for a more comprehensive view.

7.3.2 Improving the Features We Analyse

1. **Better handling of encrypted traffic** : As more traffic gets encrypted, we need to develop features that can analyse encrypted communications without decrypting them.
2. **Understanding application-layer patterns** : Adding features that understand application protocols would help us detect more sophisticated attacks like SQL injection or XSS.
3. **Getting the same insights with fewer calculations** : We could make everything more efficient by identifying which features give us the most security insight for the computational cost.
4. **Adapting to changing conditions** : Network traffic patterns evolve over time, so we need systems that can automatically adjust which features they focus on.

7.3.3 Making It More Practical

1. **Learning continuously** : Rather than periodic retraining, the system could gradually adapt to evolving threats through incremental learning.
2. **Defeating adversarial attacks** : We need to make sure attackers can't easily trick the system by crafting network traffic specifically designed to fool it.
3. **Better explaining the findings** : We could develop more intuitive ways to explain the system's decisions to security analysts who need to investigate alerts.
4. **Learning collaboratively across organizations** : Using federated learning, organizations could improve their models collectively without sharing sensitive network data.

7.3.4 Testing More Broadly

1. **Trying it on different datasets** : We need to evaluate performance across different types of networks to ensure the approach generalizes well.

2. **Real-world deployment studies** : Long-term studies in operational environments would help us understand how the system performs over time.
3. **Creating fair comparison benchmarks** : We should develop standardized ways to compare different security approaches on a level playing field.
4. **Adapting to new environments quickly** : Research into domain adaptation could help models quickly adjust to new network environments with minimal additional training.

7.4 Final Thoughts

CyberBERT represents an important step forward in bringing the power of language AI to network security. By bridging these two worlds, we've shown that we can detect sophisticated attack patterns more effectively without requiring as much specialized knowledge.

As cyber threats continue to evolve, we need to keep finding creative ways to counter them. The success of CyberBERT suggests that combining ideas from different fields - in this case, language AI and network security - creates powerful new tools for protection.

By continuing to explore these intersections, we can develop even better ways to defend our networks against ever-changing threats. After all, security is a journey, not a destination - and we're just getting started.

8. References

1. Moore, A. W., & Zuev, D. (2005). Internet traffic classification using Bayesian analysis techniques. *ACM SIGMETRICS Performance Evaluation Review*, 33(1), 50-60.
2. Zhang, J., Xiang, Y., Wang, Y., Zhou, W., Xiang, Y., & Guan, Y. (2013). Network traffic classification using correlation information. *IEEE Transactions on Parallel and Distributed Systems*, 24(1), 104-117.
3. Wang, W., Zhu, M., Wang, J., Zeng, X., & Yang, Z. (2017). End-to-end encrypted traffic classification with one-dimensional convolution neural networks. In *IEEE International Conference on Intelligence and Security Informatics*.
4. Lopez-Martin, M., Carro, B., Sanchez-Esguevillas, A., & Lloret, J. (2017). Network traffic classifier with convolutional and recurrent neural networks for Internet of Things. *IEEE Access*, 5, 18042-18050.
5. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 9.

Special thanks to our colleagues and reviewers who provided valuable feedback throughout the development of CyberBERT.

6. Kim, S., Hong, S., Oh, J., & Lee, H. (2019). Multimodal attention network for detection of security vulnerabilities. In Proceedings of the ACM International Conference on Information and Knowledge Management.

7. Lashkari, A. H., Draper-Gil, G., Mamun, M. S. I., & Ghorbani, A. A. (2017). Characterization of encrypted and VPN traffic using time-related features. In Proceedings of the 2nd International Conference on Information Systems Security and Privacy.

8. Sharafaldin, I., Lashkari, A. H., & Ghorbani, A. A. (2018). Toward generating a new intrusion detection dataset and intrusion traffic characterization. In Proceedings of the 4th International Conference on Information Systems Security and Privacy.

9. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

10. Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108.

11. Ponemon Institute. (2020). Cost of a Data Breach Report. IBM Security.

9. Acknowledgments

The authors gratefully acknowledge the foundational work of the Canadian Institute for Cybersecurity (CIC) in developing the CICFlowMeter framework, which inspired our flow-based feature extraction methodology. We extend our appreciation to the Hugging Face team for their development of the DistilBERT model architecture, which we fine-tuned and adapted specifically for real-time classification of network traffic patterns.

This research builds upon open-source contributions from the machine learning and cybersecurity communities. The CICIDS2017 dataset was instrumental in training and evaluating our model. We also acknowledge the PyTorch and Transformers libraries, which provided the technical foundation enabling this work.

Our implementation leverages design principles from industry best practices for real-time machine learning systems in security applications. The scalable architecture presented in this paper was informed by operational requirements from network security professionals and academic research in the field of AI-powered intrusion detection systems.