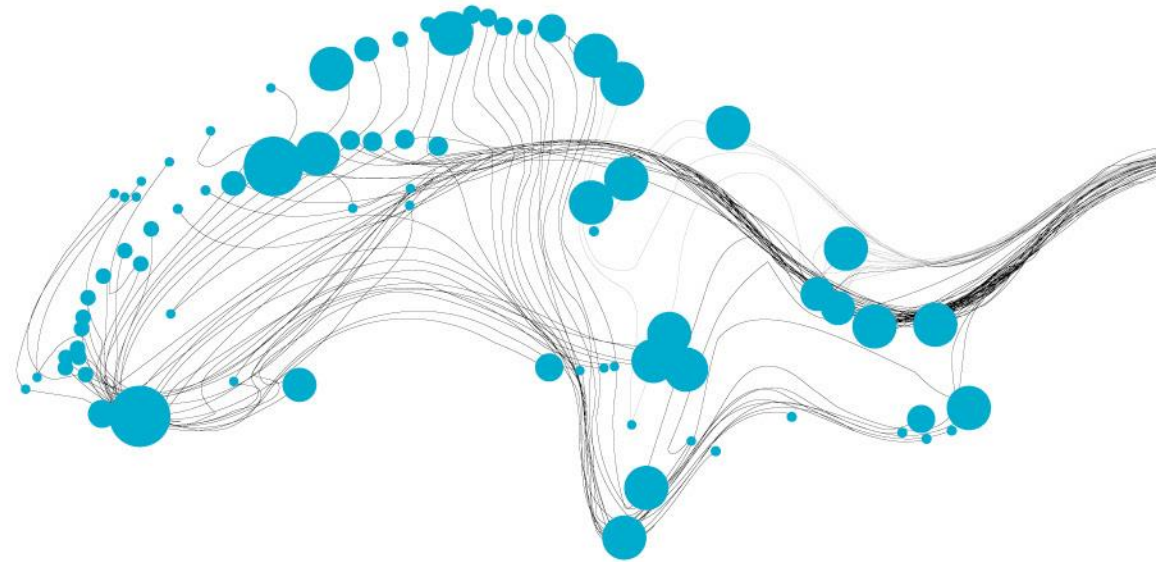# AUTOMATING THE CYBERSECURITY TRIAGE PROCESS

## A COMPARATIVE STUDY ON THE PERFORMANCE OF LARGE LANGUAGE MODELS

**PASCAL BAKKER**
**SUPERVISED BY JAIR SANTANNA**
**2024-07-05**

UNIVERSITY
OF TWENTE.

# HOSPITAL – TRIAGE

1. Immediate
2. Emergent
3. Urgent
4. Less urgent
5. Non-urgent

# HOSPITAL – TRIAGE
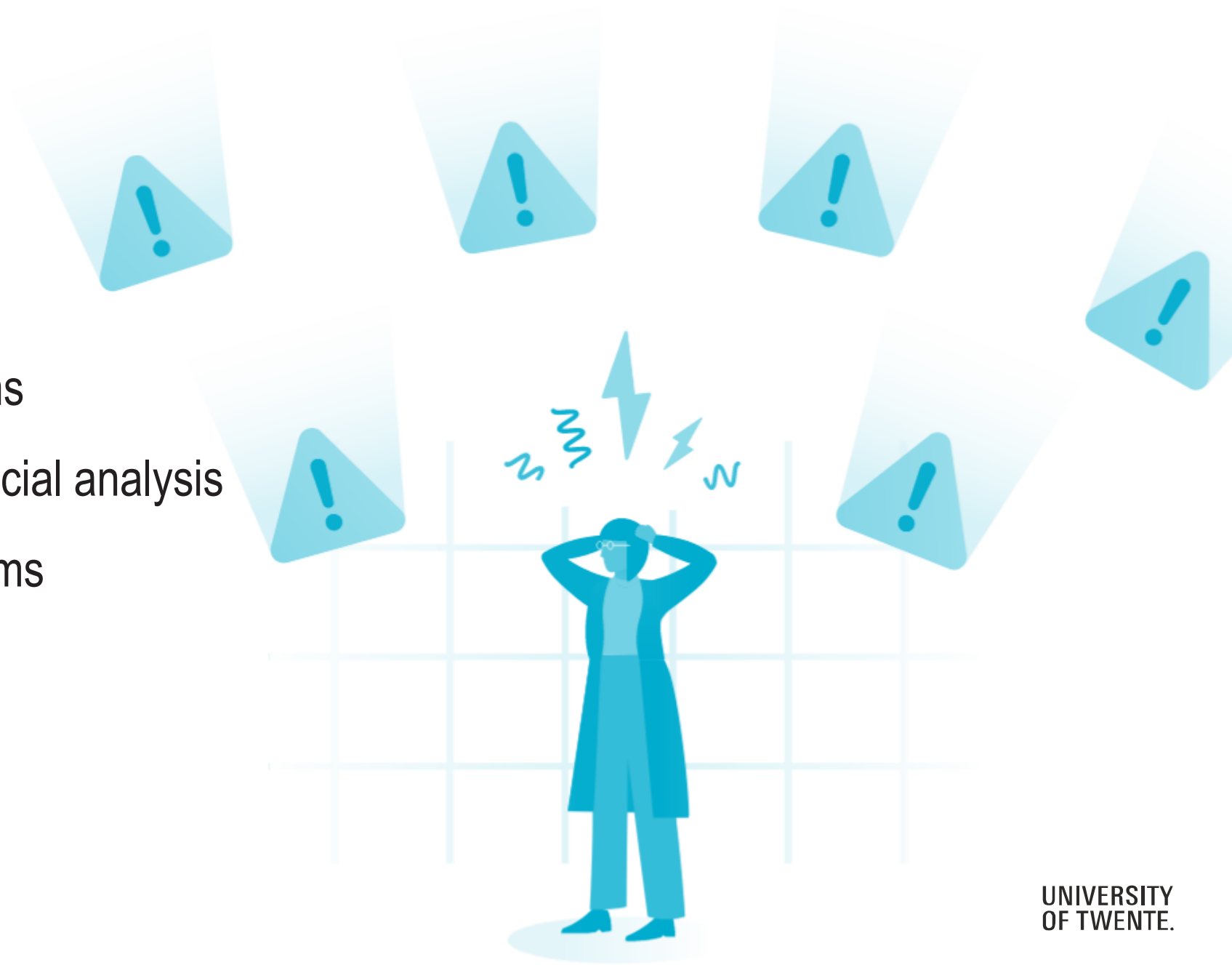
1. Immediate
2. Emergent
3. Urgent
4. Less urgent
5. Non-urgent

UNIVERSITY OF TWENTE.

# THE PROBLEM WITH TRIAGE

Immense numbers of alarms

- Too little time → Superficial analysis

- Fatigue → Missed alarms

- Burnout / turnover

- Human error

UNIVERSITY OF TWENTE.

# LARGE LANGUAGE MODELS (LLMS)

- Natural Language Processing

- Much data → General understanding

1. Generating natural language
2. Identifying contextual relationships
3. Recognizing complex patterns
4. Analyzing semantics

- Many existing applications

→ Automate triage

UNIVERSITY
OF TWENTE.

# RESEARCH QUESTIONS

How can <u>LLMs</u> be integrated into the existing incident response workflow to streamline the <u>triage process</u>?

What suitable <u>evaluation metrics</u> should be used to assess the performance of LLMs in cybersecurity triage?

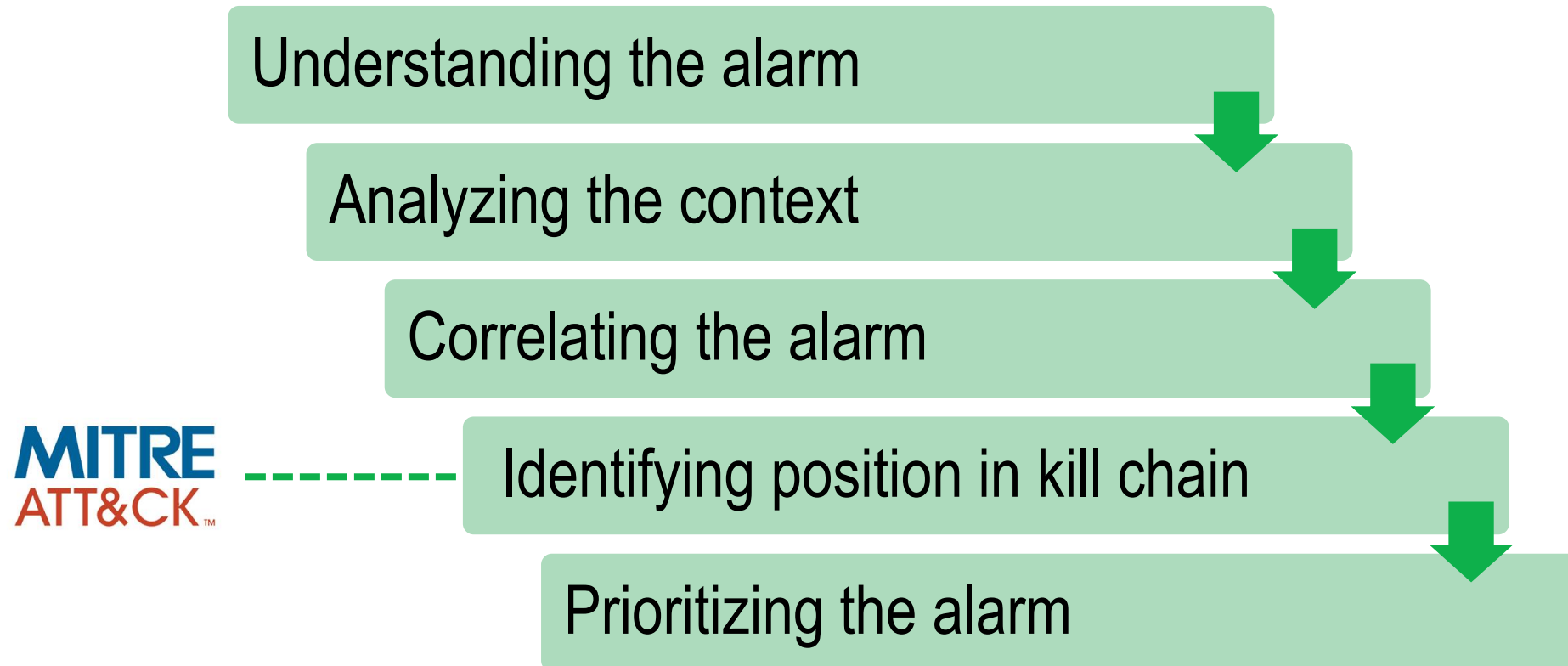How do different LLMs <u>compare in performance</u> when optimizing the cybersecurity triage process?

UNIVERSITY
OF TWENTE.

# DEFINING TRIAGE



Very little documentation

Interview

# DEFINING TRIAGE

Understanding the alarm

Analyzing the context

Correlating the alarm

Identifying position in kill chain

Prioritizing the alarm

UNIVERSITY
OF TWENTE.

# OPTIMIZING TRIAGE
## USING LLMS

**EXISTING OPTIMIZATIONS**

- Organize alarms in trees [12]

- Integrate thread intelligence [41]

- Follow steps of senior analysts [25, 48]

UNIVERSITY
OF TWENTE.

# OPTIMIZING TRIAGE
## USING LLMS

**EXISTING OPTIMIZATIONS**

- Organize alarms in trees [12]

- Integrate thread intelligence [41]

- Follow steps of senior analysts [25, 48]

Do not automate natural language tasks

UNIVERSITY
OF TWENTE.

# OPTIMIZING TRIAGE
## USING LLMS

1. Detecting cybersecurity announcement emails

2. Detecting relation between email and alarm

3. Finding correlation between alarms

4. Determine position in kill chain

5. Determine priority of alarm

UNIVERSITY
OF TWENTE.

# OPTIMIZING TRIAGE
## USING LLMS

1. <u>Detecting cybersecurity announcement emails</u>

2. Detecting relation between email and alarm

3. Finding correlation between alarms

4. Determine position in kill chain

5. Determine priority of alarm

UNIVERSITY
OF TWENTE.

# TASK 1
## EVALUATING LLMS

**ANNOUNCEMENT DETECTION**

1. Give email to LLM

2. Is this a cybersecurity announcement?

3. True / false

UNIVERSITY
OF TWENTE.

**EVALUATING LLMS**

## ANNOUNCEMENT DETECTION

1. Give email to LLM

2. Is this a cybersecurity announcement?

3. True / false

> **Customer X will add the user "sea_line" to the local administrators on the computer "FRLIM-IPC-0017".**

"True" → ✅

"False" → ❌

UNIVERSITY
OF TWENTE.

# TASK 2
## EVALUATING LLMS

**TACTIC DETECTION**

1. Give email to LLM

2. What MITRE ATT&CK <u>tactic</u> can <u>consequential alarms</u> have?

3. e.g. "exfiltration"

UNIVERSITY
OF TWENTE.

# TASK 2
## EVALUATING LLMS

**TACTIC DETECTION**

1. Give email to LLM

2. What MITRE ATT&CK <u>tactic</u> can <u>consequential alarms</u> have?

3. e.g. "exfiltration"

> **Customer X will add the user "sea_line" to the local administrators on the computer "FRLIM-IPC-0017".**

"privilege escalation" → ✅

"persistence" → ✅

"reconnaissance" → ❌

UNIVERSITY
OF TWENTE.

# EVALUATING AND COMPARING LLMS

**EVALUATION METRICS**

For different prompts:
- F1-score / accuracy
- Median time
- Error rate

UNIVERSITY
OF TWENTE.

# EVALUATING AND COMPARING LLMS

**EVALUATION METRICS**

For different prompts:

- F1-score / accuracy

- Median time

- Error rate

**DATASET**

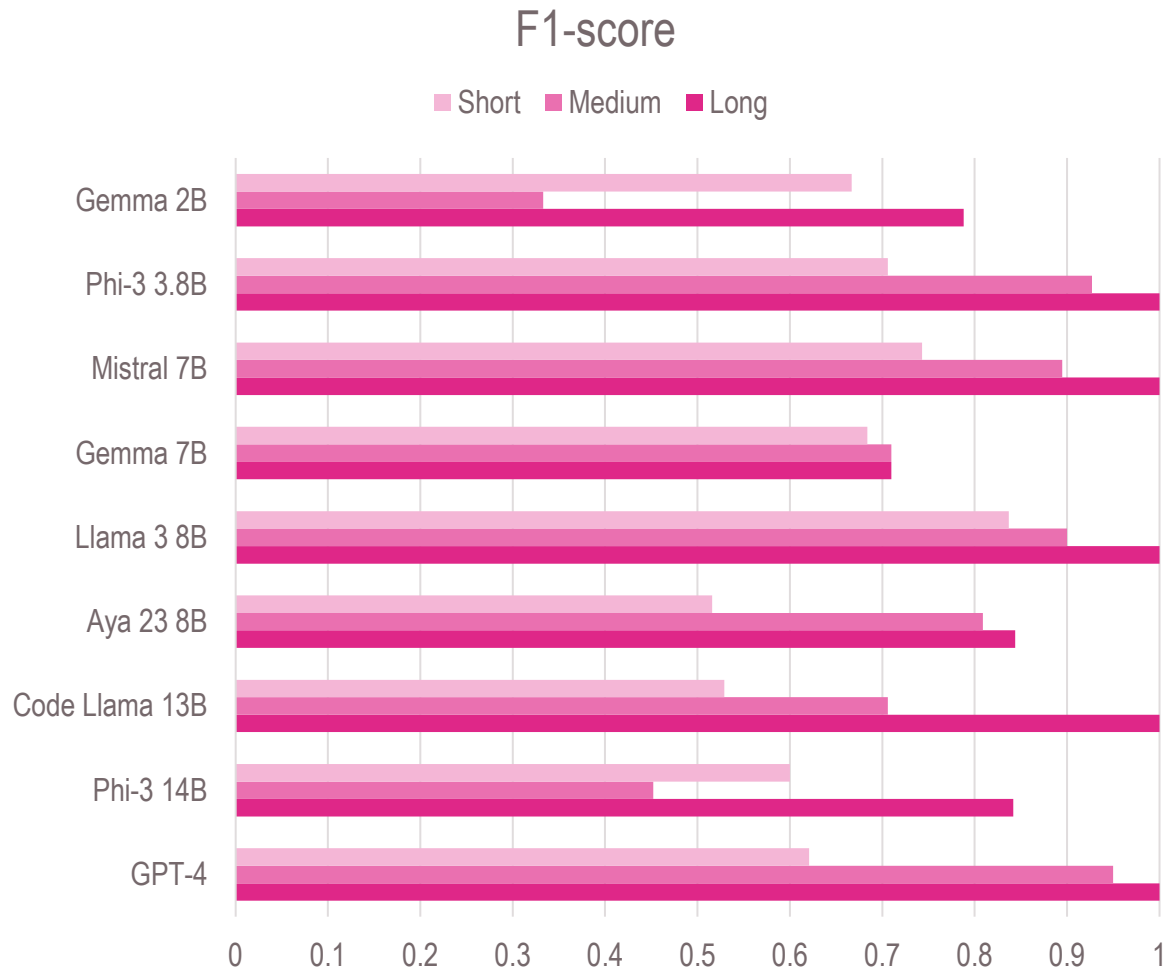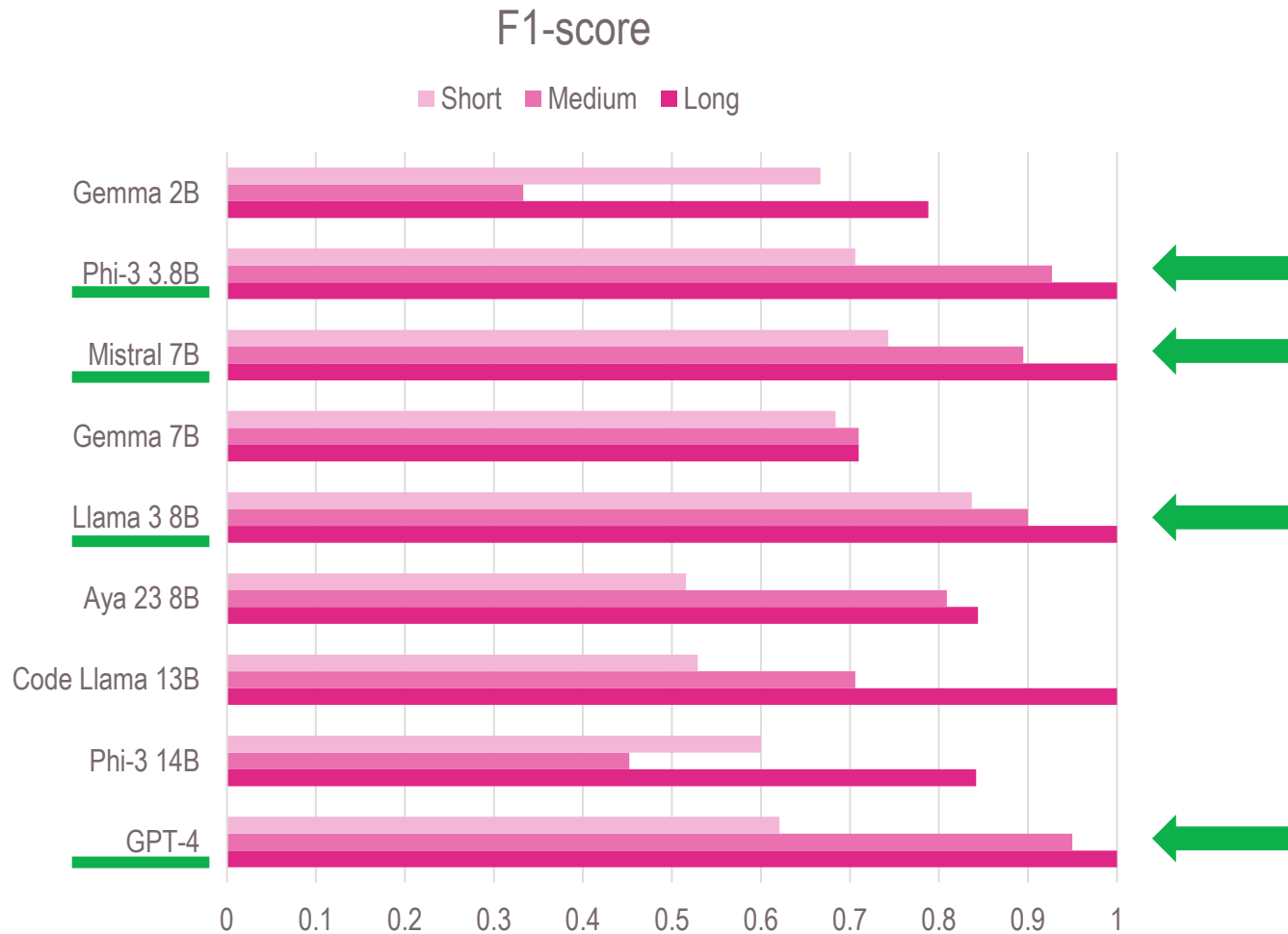- 40 labeled emails

- 20 announcements with labeled tactics

UNIVERSITY
OF TWENTE.

# EVALUATING AND COMPARING LLMS

**EVALUATION METRICS**

For different prompts:
- F1-score / accuracy
- Median time
- Error rate

**DATASET**
- 40 labeled emails
- 20 announcements with labeled tactics

**LARGE LANGUAGE MODELS**

(And variations)
- GPT-4
- Llama 3
- Mistral
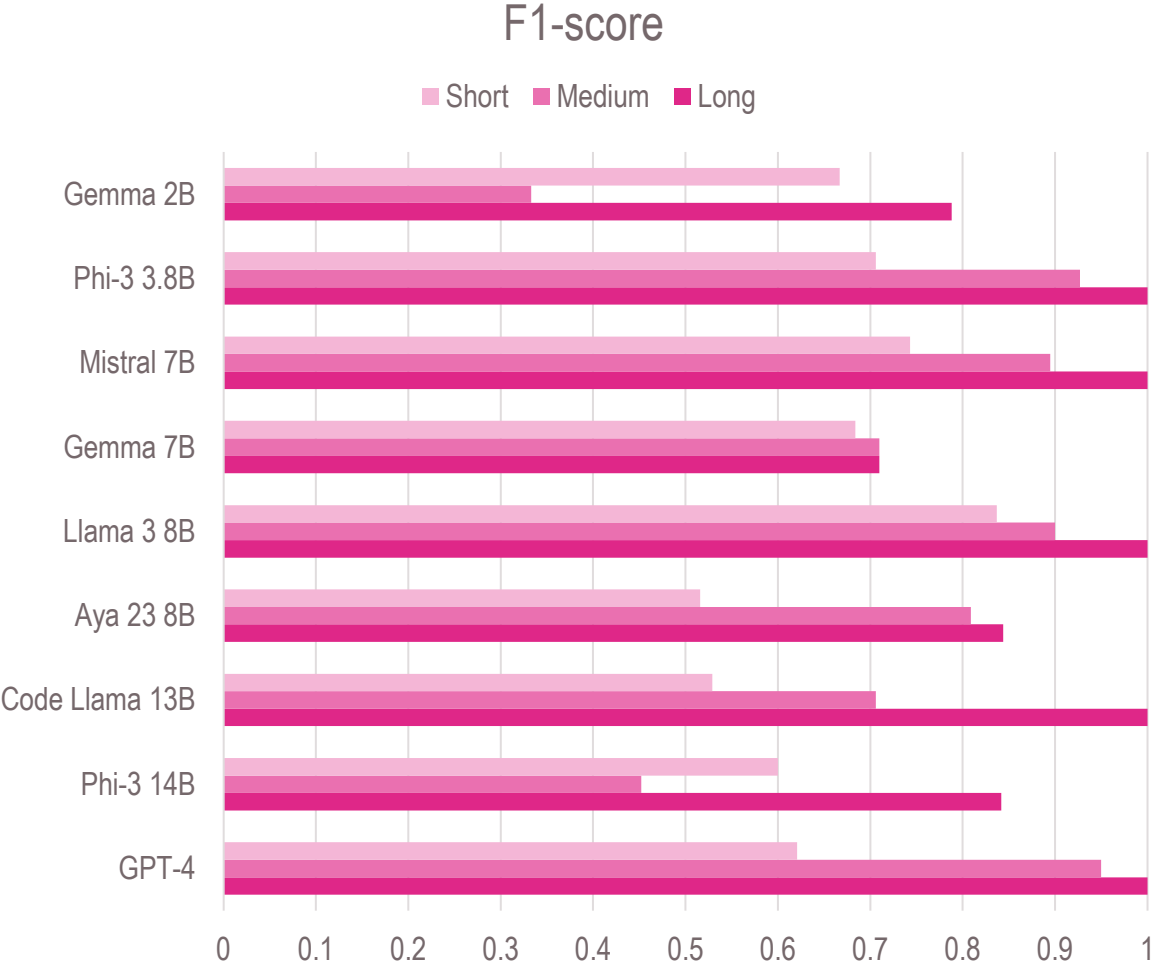- Phi-3
- Gemma
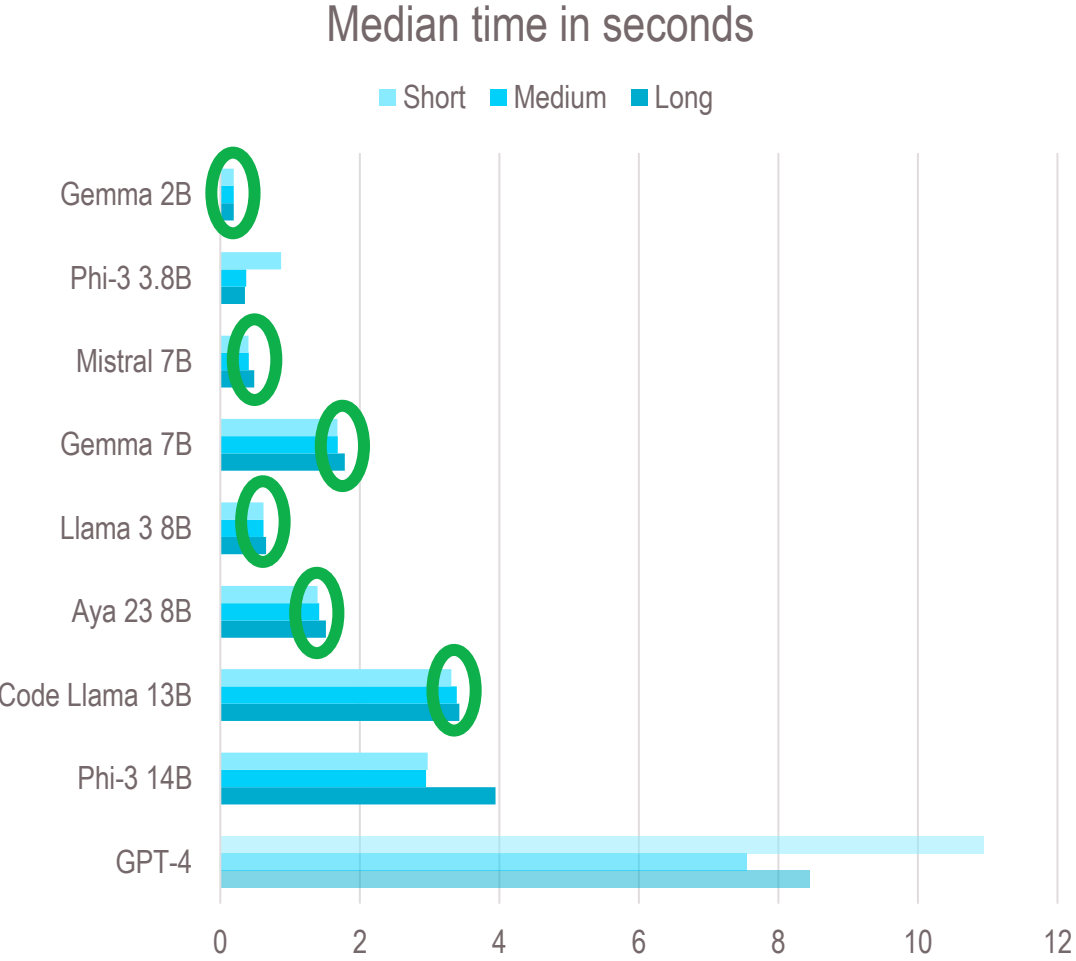- Aya 23
- Code Llama
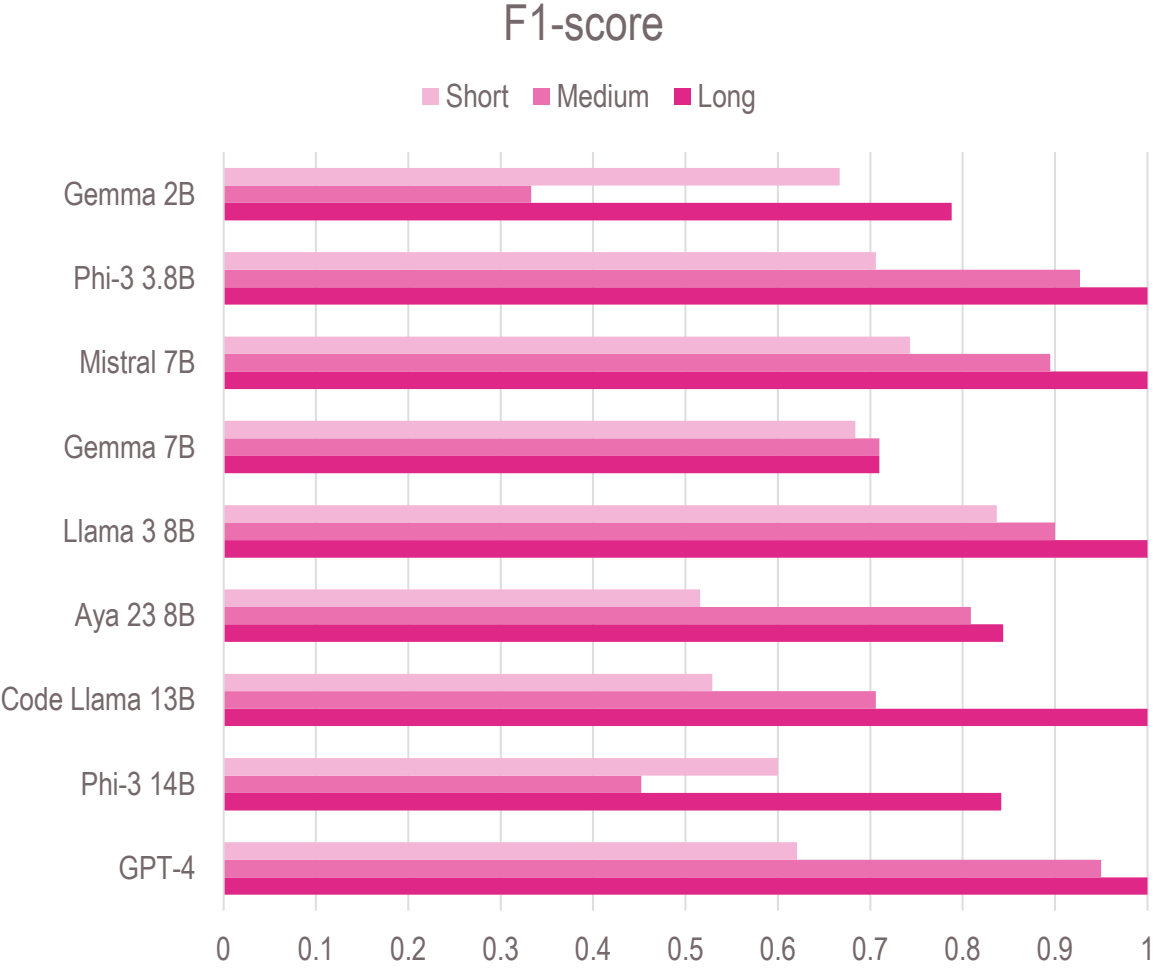
# RESULTS – ANNOUNCEMENT DETECTION



F1-score

Short    Medium    Long

Models (top to bottom): Gemma 2B, Phi-3 3.8B, Mistral 7B, Gemma 7B, Llama 3 8B, Aya 23 8B, Code Llama 13B, Phi-3 14B, GPT-4

X-axis: 0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1

UNIVERSITY
OF TWENTE.

# RESULTS – ANNOUNCEMENT DETECTION

F1-score

■ Short  ■ Medium  ■ Long

# RESULTS – ANNOUNCEMENT DETECTION



F1-score

Median time in seconds

UNIVERSITY
OF TWENTE.
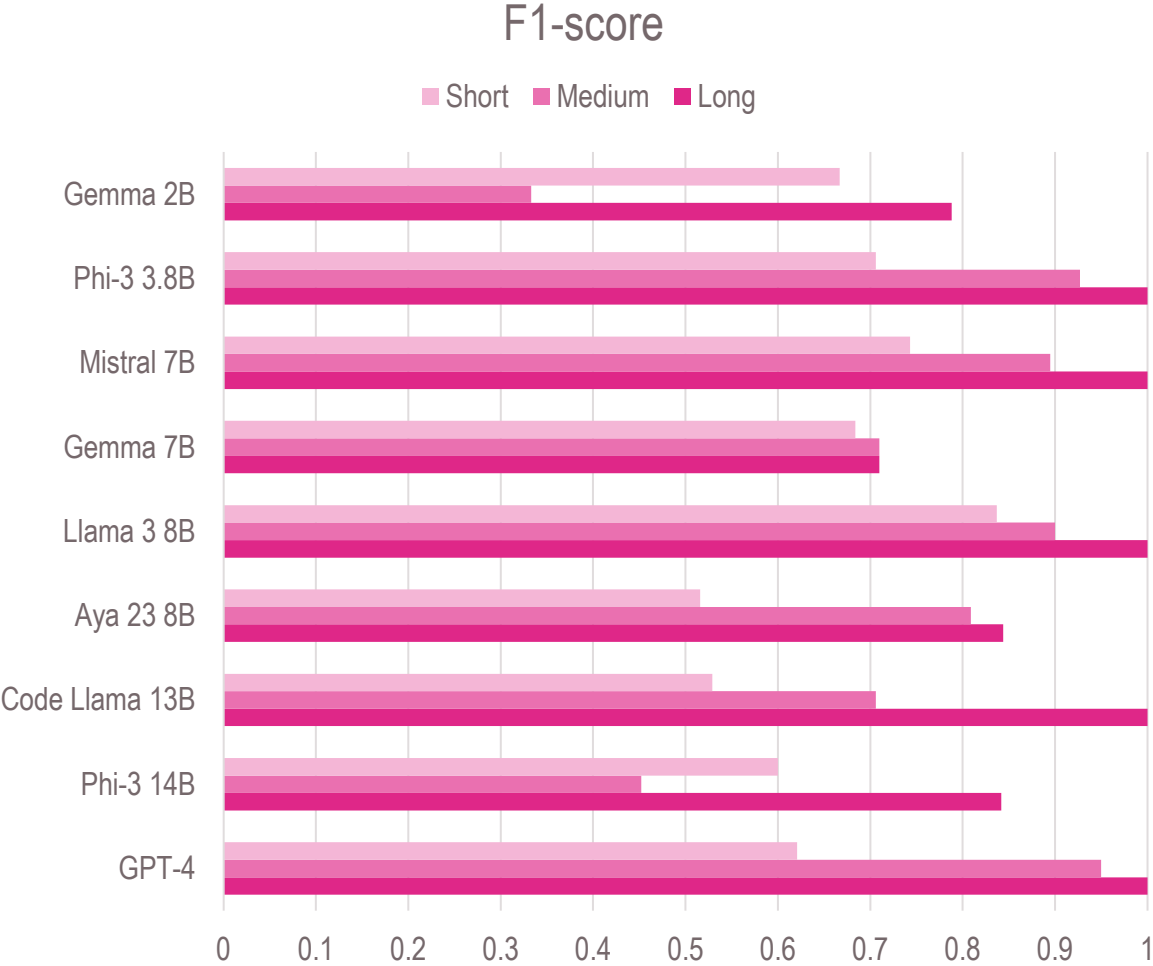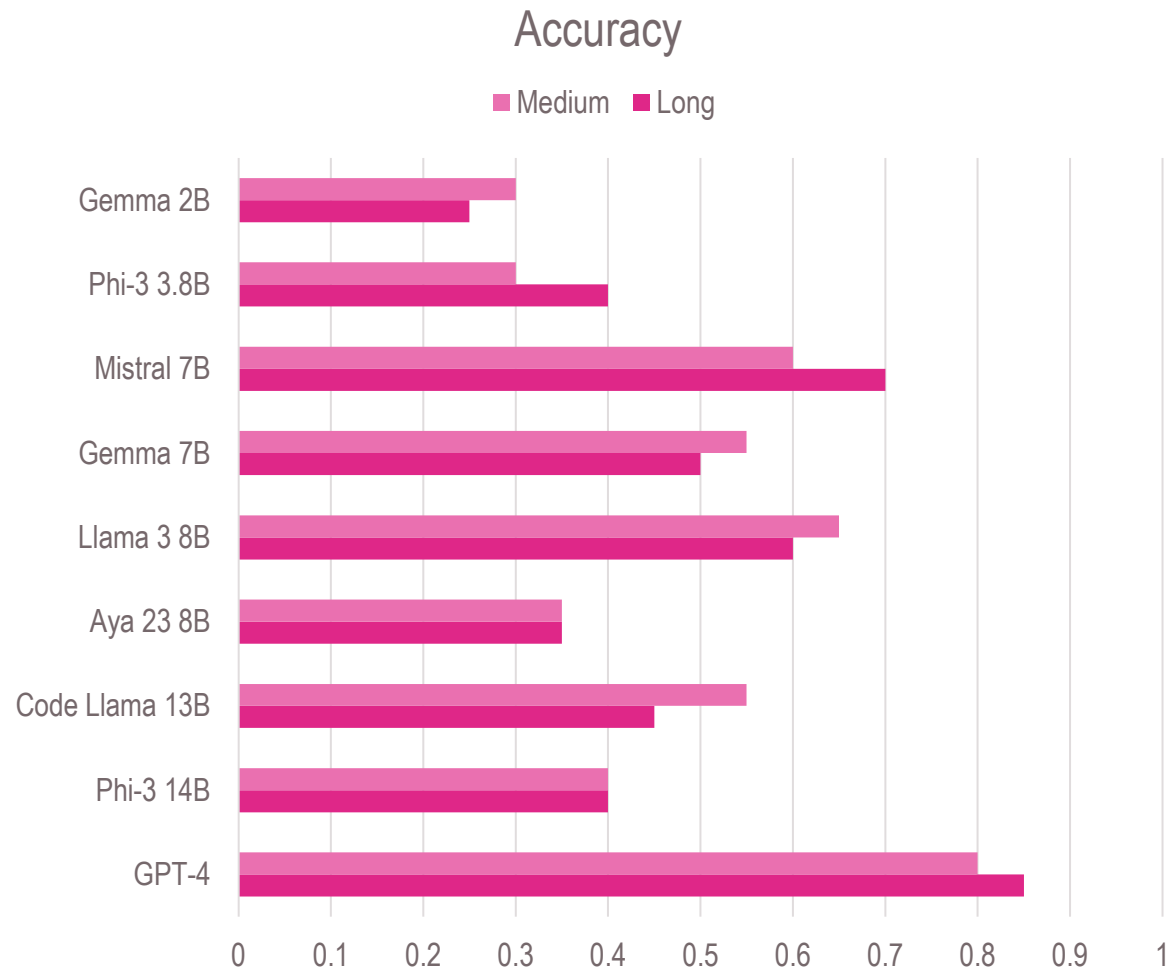
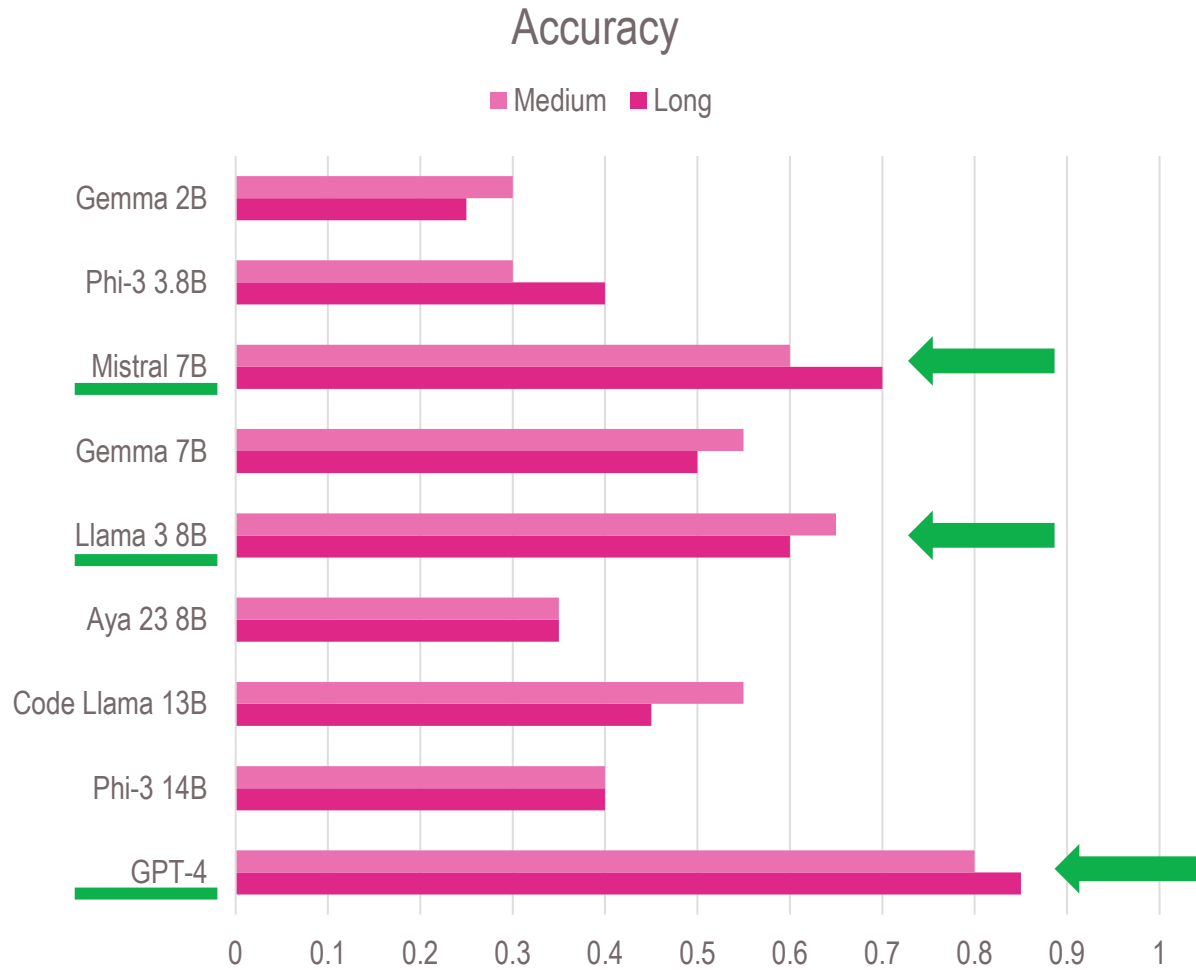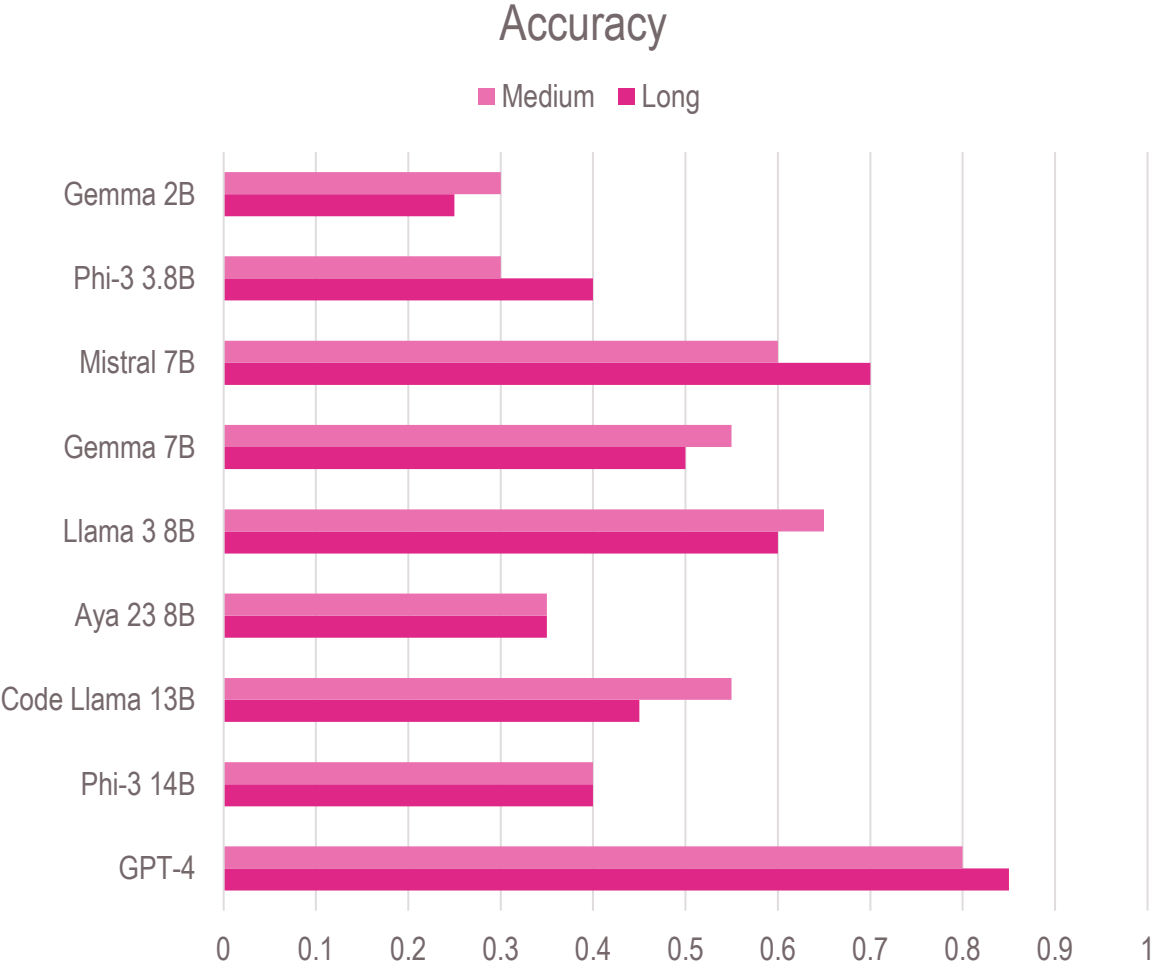# RESULTS – ANNOUNCEMENT DETECTION



13

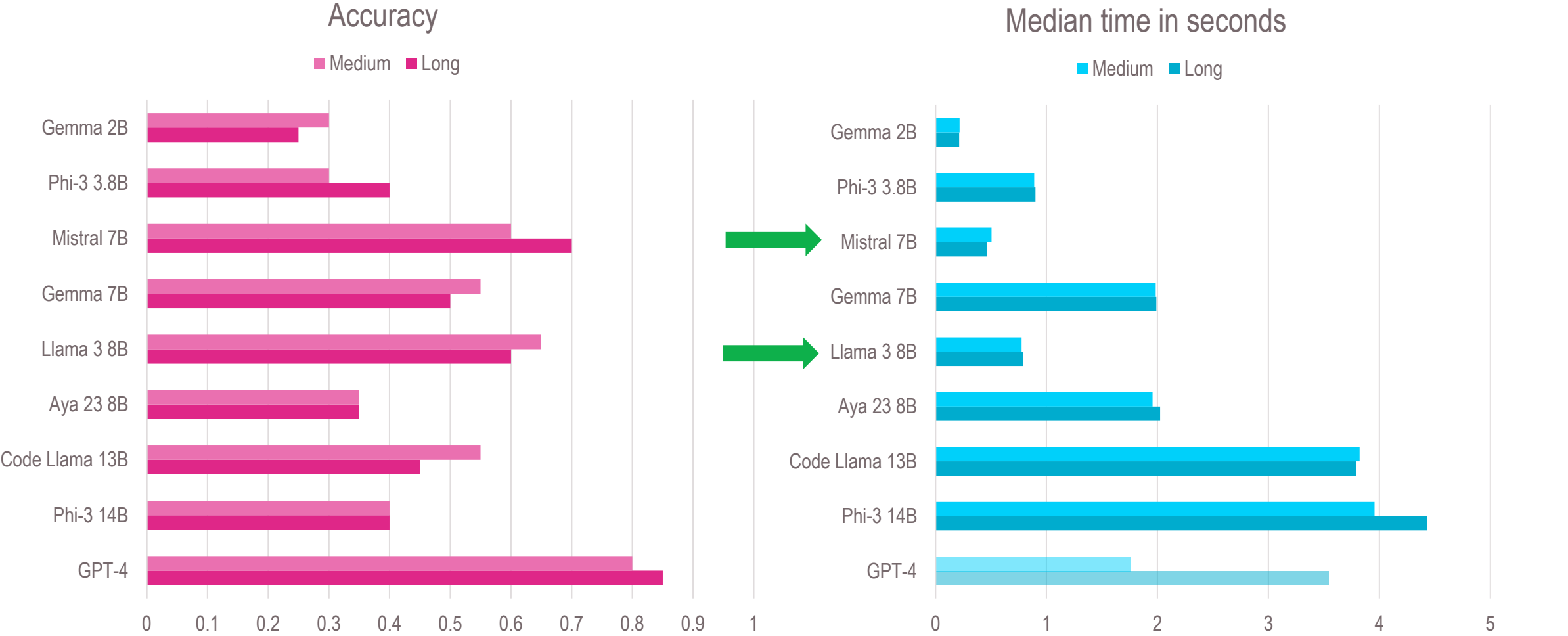# RESULTS – ANNOUNCEMENT DETECTION



13

# RESULTS – TACTIC DETECTION



Accuracy

■ Medium  ■ Long

UNIVERSITY
OF TWENTE.

# RESULTS – TACTIC DETECTION



Accuracy

■ Medium  ■ Long

UNIVERSITY OF TWENTE.

# RESULTS – TACTIC DETECTION



Accuracy

■ Medium  ■ Long

Median time in seconds

■ Medium  ■ Long

UNIVERSITY
OF TWENTE.

# RESULTS – TACTIC DETECTION

UNIVERSITY
OF TWENTE.

# RESULTS – COMPARISON



Announcement – F1-score

Tactic – Accuracy

UNIVERSITY
OF TWENTE.

# RESULTS – COMPARISON



Announcement – F1-score

Tactic – Accuracy

15

# RESEARCH QUESTIONS

How can <u>LLMs</u> be integrated into the existing incident response workflow to streamline the <u>triage process</u>?

What suitable <u>evaluation metrics</u> should be used to assess the performance of LLMs in cybersecurity triage?

How do different LLMs <u>compare in performance</u> when optimizing the cybersecurity triage process?

UNIVERSITY
OF TWENTE.

# RESEARCH QUESTIONS

How can <u>LLMs</u> be integrated into the existing incident response workflow to streamline the <u>triage process</u>?

What suitable <u>evaluation metrics</u> should be used to assess the performance of LLMs in cybersecurity triage?
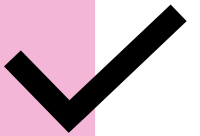
How do different LLMs <u>compare in performance</u> when optimizing the cybersecurity triage process?

UNIVERSITY
OF TWENTE.

# RESEARCH QUESTIONS

How can <u>LLMs</u> be integrated into the existing incident response workflow to streamline the <u>triage process</u>?

What suitable <u>evaluation metrics</u> should be used to assess the performance of LLMs in cybersecurity triage?
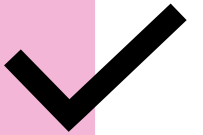
How do different LLMs <u>compare in performance</u> when optimizing the cybersecurity triage process?

UNIVERSITY
OF TWENTE.

# RESEARCH QUESTIONS

How can <u>LLMs</u> be integrated into the existing incident response workflow to streamline the <u>triage process</u>?

What suitable <u>evaluation metrics</u> should be used to assess the performance of LLMs in cybersecurity triage?
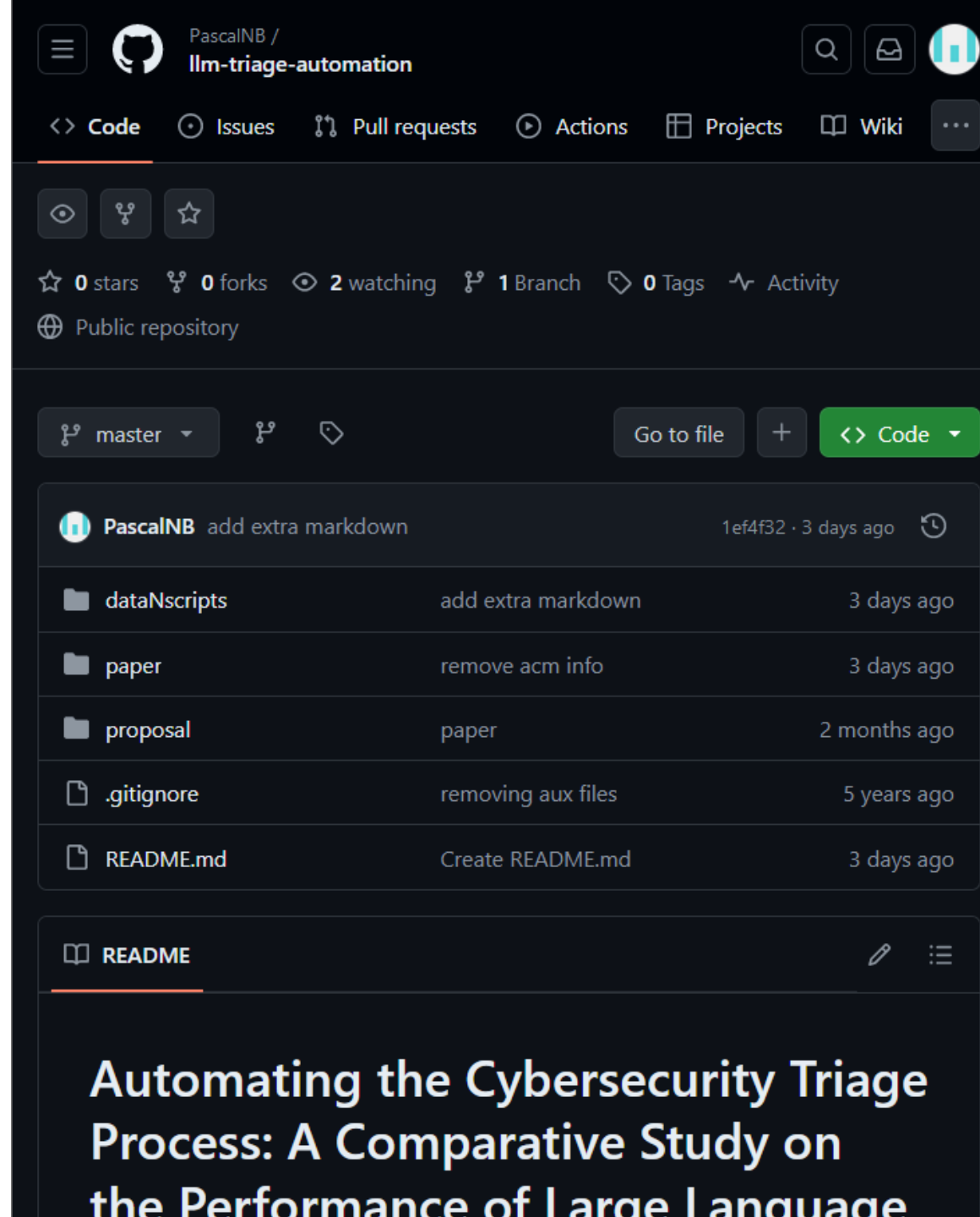
How do different LLMs <u>compare in performance</u> when optimizing the cybersecurity triage process?
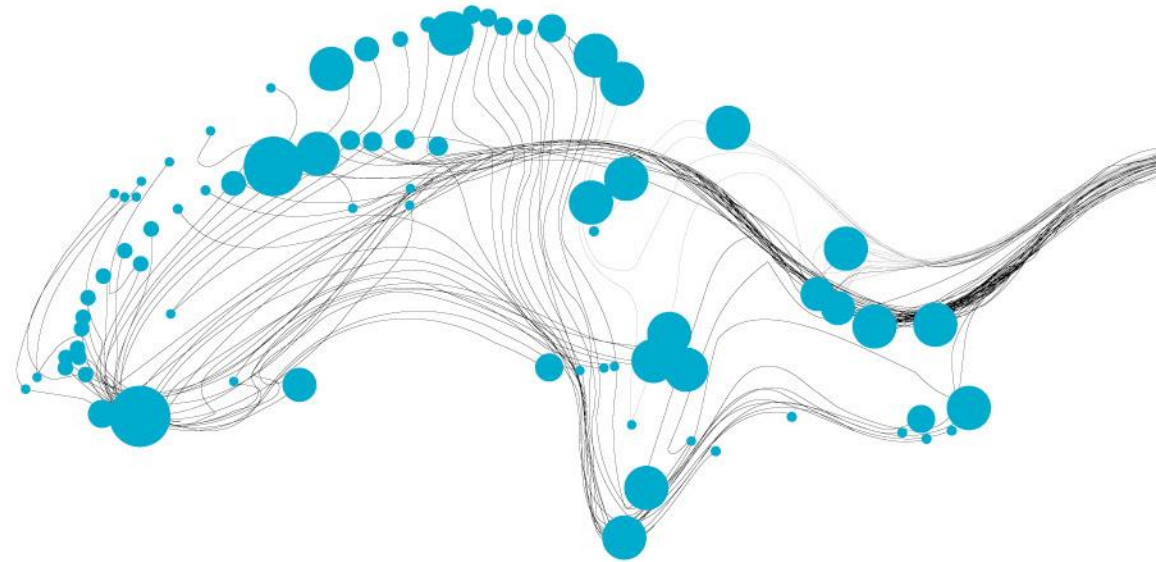
UNIVERSITY OF TWENTE.

# TAKEAWAYS

1. GPT-4 performed the best

   - Followed by Llama 3 and Mistral
   - Phi-3 3.8B good in simple tasks

2. Prompt size had no effect on time

3. Baseline towards further usage of LLMs

   - Defined key steps of triage
   - Identified optimizations
   - Evaluation framework
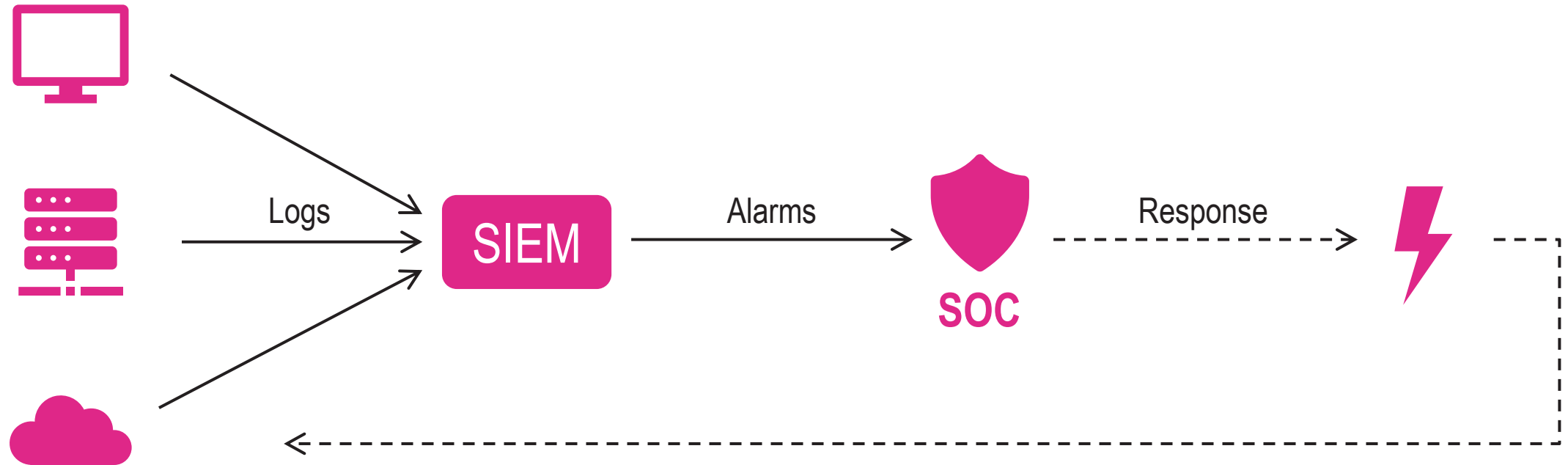   → Other tasks and models

# AUTOMATING THE CYBERSECURITY TRIAGE PROCESS

## A COMPARATIVE STUDY ON THE PERFORMANCE OF LARGE LANGUAGE MODELS

PASCAL BAKKER
SUPERVISED BY JAIR SANTANNA
2024-07-05

UNIVERSITY
OF TWENTE.

# THE INCIDENT RESPONSE WORKFLOW



UNIVERSITY
OF TWENTE.

# ERROR RATE

```
{
 "is_announcement": True
}
```

✅

```
{
 "is_annoonment": True
}
```

❌

UNIVERSITY
OF TWENTE.

# EVALUATION METRICS

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$