

54,000,000

DatAvengers Cybersecurity



Amaan Ansari

[Am-aan]

Lead Analyst

**MS in Business Analytics and
Information Management**

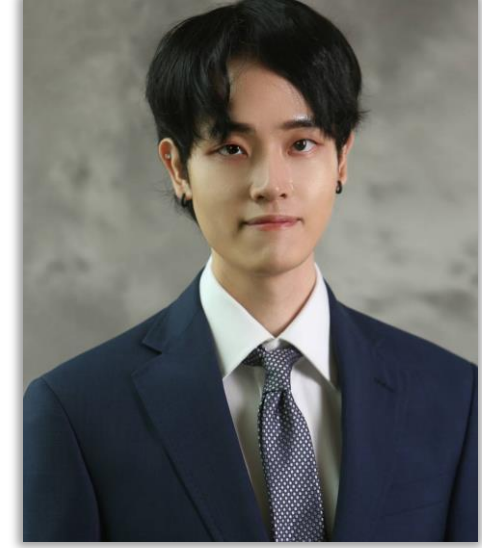


Jai Woo Lee

[Jay]

Data Analyst

**MS in Business Analytics and
Information Management**



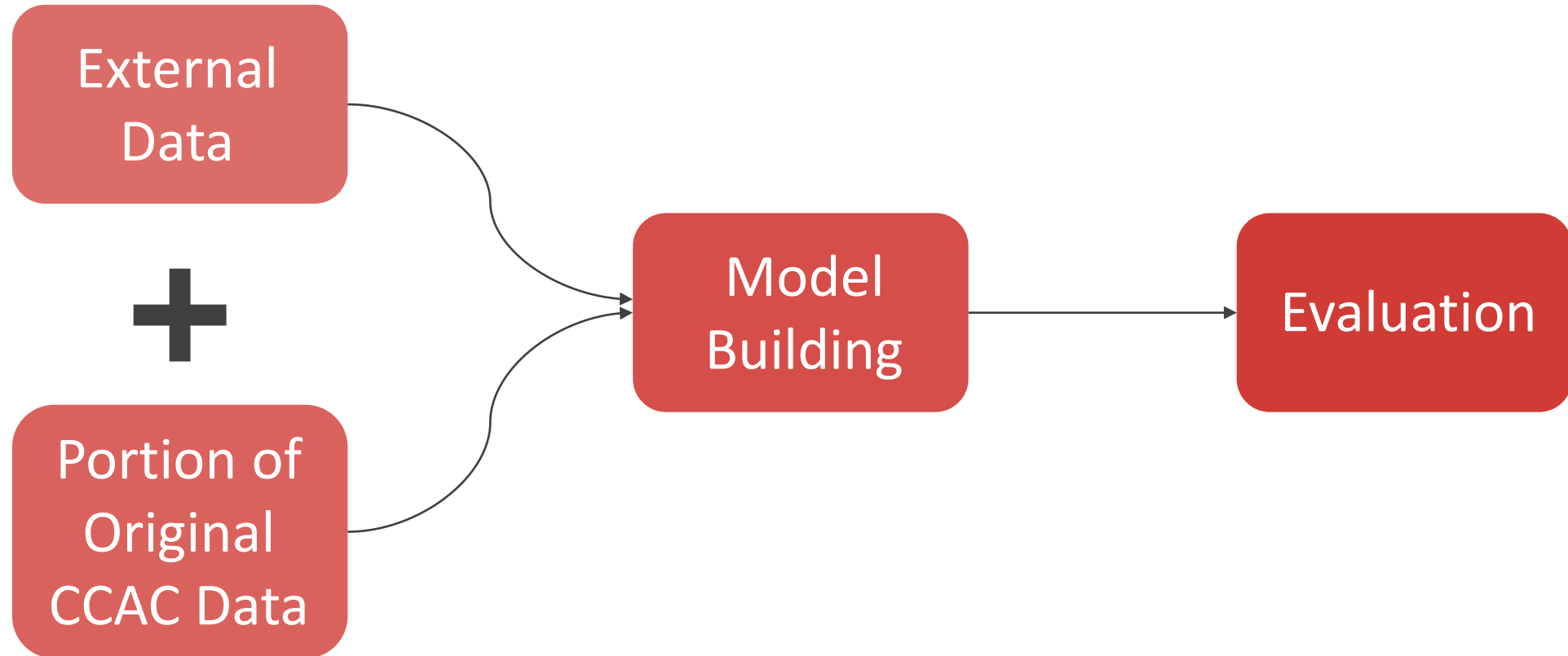
Paul Chen

[Paul]

Data Scientist

**MS in Business Analytics and
Information Management**

APPROACH



1
How?

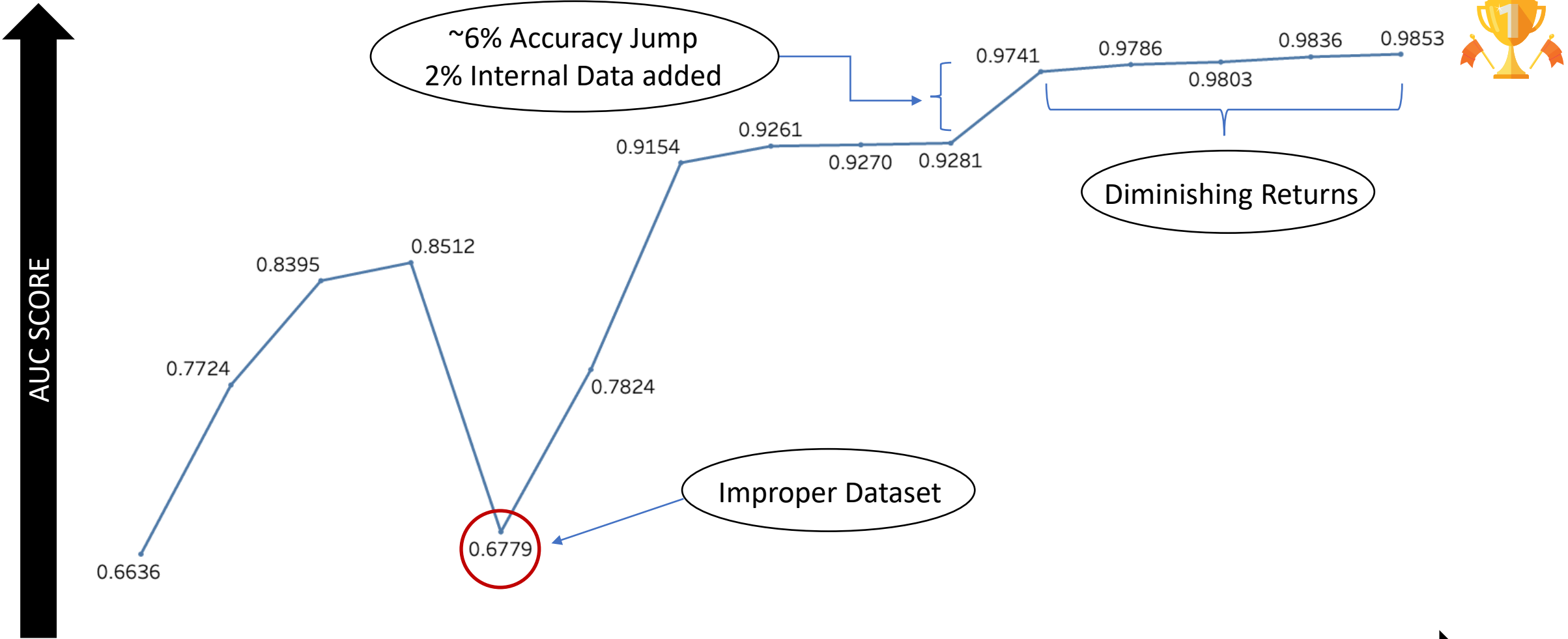
2
What?

3
Solution

4
Risks

5
Conclusion

ROADMAP



SVM+XGBoost



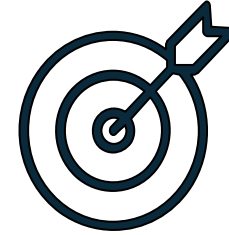
PROJECT TIMELINE

- 1 How?
- 2 What?
- 3 Solution
- 4 Risks
- 5 Conclusion

WHY INTERNAL DATA?



Internal Data Patterns



Higher Accuracy



Creates custom fit



URL Link Discounting

F.I.S.H PROTOTYPE



Filter

FISH's high accuracy model can filter most of the phishing email



Summary

FISH provides monthly summary of the latest phishing news and word cloud



Identify

The hackers are improving everyday so is FISH. FISH identifies the recent phishing key words and update our model everyday



Hook-free

FISH provides user friendly interface letting user decide what is phishing and the adjustable threshold design

1
How?

2
What?

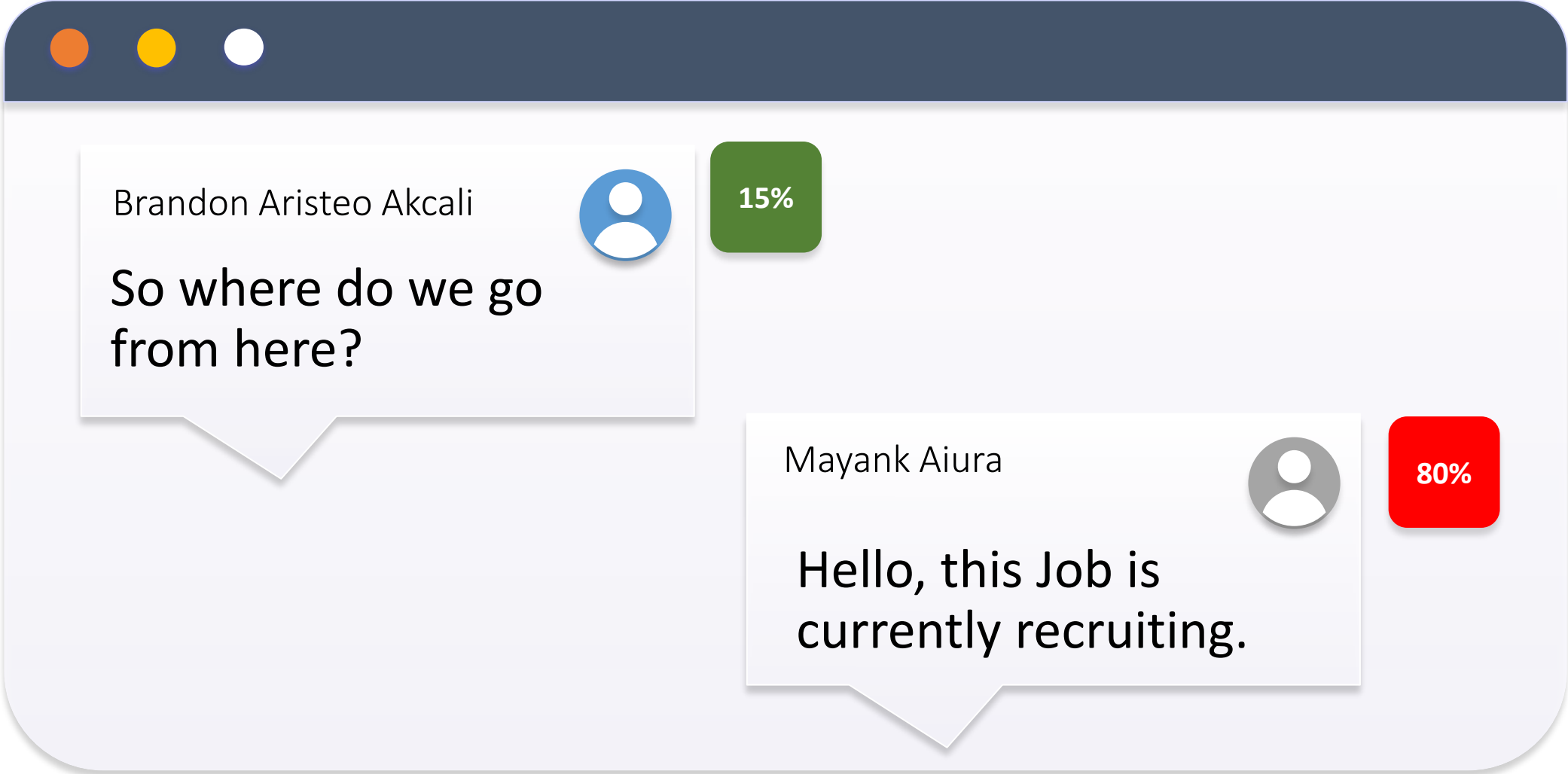
3
Solution

4
Risks

5
Conclusion

7

F.I.S.H PROTOTYPE



F.I.S.H PROTOTYPE

Low Risk

Less than 30%

Neutral

Between 30% to 60%

High Risk

Higher than 60%



The number can be tailored !

RISK ANALYSIS

1 Confusion Matrix

		Actual	
		Phishing	Non-Phishing
Predicted	Phishing	9700(TP)	60(FP)
	Non-Phishing	90(FN)	150(TN)

2 Cost Matrix

		Actual	
		Phishing	Non-Phishing
Predicted	Phishing	\$0	\$0
	Non-Phishing	\$224*	\$0

Opportunity Cost

- Work-related emails
- Important personal emails

Exposed to Phishing Emails

- Malware/ransomware
- Data breach

3 Mitigation Strategy

- Adjust a threshold

Potential Cost for Misclassification

\$20,160

* Average Loss in US (2021): \$14.8M

*In 2020, Number of phishing complaints: 241,342 | Loss: \$54 million in US (IC3 Report, 2020), \$54M / 241,342 = 224

CONCLUSION

Identification of Phishing

Word Cloud and Pattern Analysis
Ex. Upper Letters, Job Offers, Promotions, !!!!!!!!!!!!!!!

Strategic Data Selection

Incorporated internal data gradually and proved adding ~2% data generates the highest accuracy rate efficiently

98.53%

AUC Score based on Ensemble model

F.I.S.H

High accuracy Filter,
Identify new phishing key words,
Monthly Summary,
Adjustable threshold to provide Hook-free environment

Thank You