

A Three-way Decision Making Approach to Malware Analysis

Mohammad Nauman^a, Nouman Azam^a and JingTao Yao^b

^aNational University of Computer and Emerging Sciences, Peshawar, Pakistan

^bDepartment of Computer Science, University of Regina

[mohammad.nauman,nouman.azam]@nu.edu.pk
jtyao@cs.uregina.ca



November 17, 2015

Introduction

- The protection of digital devices from illegal use is an important issues for obvious reasons.
 - Protection of personal data.
 - Protection of enterprise data.
 - Protection of governmental data.
- **Malware Analysis:** The set of techniques and tools used to ensure protection of digital devices.
- A recent increase in targeted attacks.
 - In 2013, a 91% increase in targeted attack campaigns (Sem., 2015).
 - No less than 38% users have experienced mobile cyber crime in 2014 and 2015 (Sem., 2015).
- **Challenge:** Investigation of effective techniques for detecting malicious activity on digital devices.

Symantec Inc., 2015. Internet security threat report, volume 19. Accessed: Apr 12, 2015, http://www.symantec.com/security_response/publications/threatreport.jsp.

An Issue with Existing Malware Analysis Techniques

- **Existing approaches:** They are generally based on two-way classification of application behaviour.
 - An application behaviour is either classified as being malicious (harmful) or benign (not harmful).
- The two-way classification may not be effective in many cases.
 - A malicious application occasionally behaving like benign (for deceiving the analysis engine). This may lead to ambiguous information not sufficient for precise classification.
 - **Problem:** The two-way approaches are based on classifying every case. We may misclassify the cases with low quality of associated information.
- We propose and examine **three-way decision making approach for malware analysis**.
 - The rationale is to defer the classification decisions of cases that have low level of associated information.

Three-way Decision Approach to Malware Analysis

- Three-way decisions are based on three decisions of,
 - acceptance,
 - rejection,
 - deferment.
- The deferment decision option provides benefits in at least two aspects.
 - More flexible compared to two-way → immediate decisions versus deferment.
 - Provides a mechanism for explicitly identifying the cases with low quality of information.
- We consider rough sets based three-way approaches.

Rough Sets

- Sets from imperfect, imprecise and incomplete data may not be precisely defined.
 - Sets have to be approximated.
- Approximating a concept C with objects in U (Pawlak 1982).
 - Lower approximation given by $\underline{apr}(C) = \{x \in U \mid [x] \subseteq C\}$,
 - Upper approximation given by $\overline{apr}(C) = \{x \in U \mid [x] \cap C \neq \phi\}$.
- Three regions may be defined using these approximations.
 - $POS(C) = \underline{apr}(C)$,
 - $BND(C) = \overline{apr}(C) - \underline{apr}(C)$,
 - $NEG(C) = (\overline{apr}(C))^c$.

Pawlak, Z. (1982). Rough sets, *International Journal of Computer and information Sciences*, 11.

Probabilistic Rough Sets (PRS)

- Restrictness of the Pawlak model.
 - The degree of an overlap between $[x]$ and C is not considered.
 - Strict conditions for inclusion in positive and negative regions.
- Probabilistic rough sets (PRS) (Yao, 2008).
 - Considers the overlap between $[x]$ and C in the form of conditional probability.
 - Pair of thresholds (α, β) are used to define approximations (Yao, 2008).
 - $\underline{apr}_{(\alpha, \beta)}(C) = \bigcup\{x \in U \mid Pr(C|[x]) \geq \alpha\}$,
 - $\overline{apr}_{(\alpha, \beta)}(C) = \bigcup\{x \in U \mid Pr(C|[x]) > \beta\}$.
 - Probabilistic positive, negative and boundary regions,
 - $POS_{(\alpha, \beta)}(C) = \{x \in U \mid Pr(C|[x]) \geq \alpha\}$,
 - $NEG_{(\alpha, \beta)}(C) = \{x \in U \mid Pr(C|[x]) \leq \beta\}$,
 - $BND_{(\alpha, \beta)}(C) = \{x \in U \mid \beta < Pr(C|[x]) < \alpha\}$.

Yao, Y. Y., (2008). Probabilistic rough set approximations, IJAR, 49.

PRS Models and Approaches

- PRS Models.
 - Decision-theoretic rough sets (Yao & Wong, 1992).
 - Variable precision rough sets (Ziarko, 1992).
 - 0.5-probabilistic rough sets (Pawlak, 1988).
 - Information-theoretic rough sets (Deng & Yao, 2012).
 - Game-theoretic rough sets (Yao & Herbert, 2008).
- Approaches to determination of (α, β) thresholds.
 - Optimization viewpoint (Jia et al., 2011).
 - Multi-view model (Li & Zhou, 2011).
 - Method using probabilistic model criteria (Liu et al., 2011).

Jia, X. Y., Li, W. W., Shang, L., & Chen, J. J., (2011). An optimization viewpoint of DTRS. In: (RSKT'11).

Li, H.X., & Zhou, X.Z., (2011). Risk decision making based on DTRS... IJCIS 4.

Liu, D., Li, T.R., & Ruan, D., (2011). Probabilistic model criteria with DTRS. Information Science 181.

Deng, X. F., & Yao, Y. Y., (2012). An information-theoretic interpretation of thresholds in PRS. In: (RSCTC'12).

Yao, J.T., & Herbert, J.P., (2008). A game-theoretic perspective on rough sets. JCUPT, 20(3).

Pawlak, Z., Wong, S. K. M., & Ziarko, W., (1988). Rough sets: probabilistic versus IJMMS, 29(1).

Ziarko, W., (1993). Variable precision rough set model. Journal of Computer and System Sciences, 46(1).

Yao, Y. Y., & Wong, S. K. M., (1992). A decision theoretic framework for approximating concepts. IJMM

A Key in PRS: Determination of Thresholds

- User's or expert's opinion about the thresholds may involve several error and trails (Herbert, 2010).
 - Moreover, one can not set thresholds once and for all.
- Needing a scientific method to determine thresholds (Herbert, 2010).
- The GTRS approach for threshold determination.
 - Provides threshold determination mechanism based on a game involving an often contradictive criteria (or properties) (Yao & Herbert, 2008; Herbert & Yao, 2011; Azam & Yao, 2014; Zhang and Yao, 2012).
- The ITRS approach for threshold determination.
 - Determining an effective configuration of thresholds by minimizing the overall uncertainty of probabilistic rough sets regions using the measure of Shannon Entropy (Deng & Yao, 2012; Deng & Yao, 2014).

Herbert, J.P., & Yao, J.T., (2011). Game-theoretic rough sets. *Fundamenta Informaticae*, 108(3-4).

Yao, J.T., & Herbert, J.P., (2008). A game-theoretic perspective on rough sets. *JCUPT*, 20(3).

Azam, N., & Yao J. T. (2014). Analyzing uncertainties of PRS regions with GTRS. *IJAR*, 55(1).

Yao, J. T., & Azam, N. (2014a). Three-way Decision Making in WMDSS with GTRS. *IEEE TFS*.

Zhang, Y., & Yao J. T. (2012). Rule measure tradeoff using GTRS. In: *BI'12*.

Deng, X. F., & Yao, Y. Y., (2012). An information-theoretic interpretation of thresholds in PRS. In: (RSC) UNIVERSITY OF REGINA

Deng, X. F., & Yao, Y. Y., (2014). A multifaceted analysis of probabilistic three-way decisions. *FI 132(3)*

The GTRS Approach for Threshold Determination

- A specific (α, β) pair represents a particular rough set model.
 - $(\alpha, \beta) = (1, 0)$ = Pawlak model and $(\alpha = \beta)$ = probabilistic two-way model.
- Choosing a best or better rough set model based on some properties, such as, accuracy and generality.
- However, some of these properties may have a conflict.
 - Increasing one may decrease the other.
 - Accuracy versus generality in the Pawlak model and probabilistic model.
- Examining such properties to produce a pair of (α, β) in a game setting.

A Typical Game in Game Theory

- Game theory is a core subject in decision sciences.
- The basic game components include.
 - Players.
 - Strategies.
 - Payoffs.
- A classical example in Game Theory: The prisoners dilemma.

		p_2	
		confess	don't confess
p_1	confess	p_1 serves 10 years, p_2 serves 10 years	p_1 serves 0 year, p_2 serves 20 years
	don't confess	p_1 serves 20 years, p_2 serves 0 year	p_1 serves 1 year, p_2 serves 1 year

A Typical Game in GTRS: Accuracy Versus Generality

- **Players:** Accuracy versus Generality.
- **Strategies:** Three types of strategies were formulated for each player.
 - s_1 (α_{\downarrow} , decrease of α),
 - s_2 (β_{\uparrow} , increase of β),
 - s_3 ($\alpha_{\downarrow}\beta_{\uparrow}$, decrease of α and increase of β).
- **Payoffs:** They are based on the measure of accuracy and generality.
 - $u_A(s_m, s_n) = \text{Accuracy}(\alpha, \beta)$.
 - $u_G(s_m, s_n) = \text{Generality}(\alpha, \beta)$.

The Game in the Payoff Table

		G		
		$s_1 = \alpha_{\downarrow}$	$s_2 = \beta_{\uparrow}$	$s_3 = \alpha_{\downarrow}\beta_{\uparrow}$
A	$s_1 = \alpha_{\downarrow}$	$u_A(s_1, s_1), u_G(s_1, s_1)$	$u_A(s_1, s_2), u_G(s_1, s_2)$	$u_A(s_1, s_3), u_G(s_1, s_3)$
	$s_2 = \beta_{\uparrow}$	$u_A(s_2, s_1), u_G(s_2, s_1)$	$u_A(s_2, s_2), u_G(s_2, s_2)$	$u_A(s_2, s_3), u_G(s_2, s_3)$
	$s_3 = \alpha_{\downarrow}\beta_{\uparrow}$	$u_A(s_3, s_1), u_G(s_3, s_1)$	$u_A(s_3, s_2), u_G(s_3, s_2)$	$u_A(s_3, s_3), u_G(s_3, s_3)$

The ITRS Approach for Threshold Determination

- The PRS regions have a degree of uncertainty (Deng & Yao, 2014).
 - Acceptance/rejection decisions are made with uncertainty.
- **ITRS approach:** Configuring the thresholds in order to optimize the overall uncertainty of the PRS.
 - The (α, β) thresholds control the uncertainty of the regions.

Calculating the Uncertainty in the Probabilistic Regions

- Measuring uncertainty in probabilistic regions.
 - Considering a partition based on a concept C , $\pi_C = \{C, C^c\}$.
 - Another partition based on the (α, β) thresholds,
 - $\pi_{(\alpha, \beta)} = \{\text{POS}_{(\alpha, \beta)}(C), \text{NEG}_{(\alpha, \beta)}(C), \text{BND}_{(\alpha, \beta)}(C)\}$.
 - Uncertainty in $\pi_C = \{C, C^c\}$ with respect to the three regions.
 - Using Shannon entropy (Deng & Yao, 2012).
 - E.g., the uncertainty in π_C due to positive, negative and boundary regions are,

$$\begin{aligned} \Delta_P(\alpha, \beta) &= H(\pi_C | \text{POS}_{(\alpha, \beta)}(C)) &= -P(C | \text{POS}_{(\alpha, \beta)}(C)) \log P(C | \text{POS}_{(\alpha, \beta)}(C)) \\ & & \quad -P(C^c | \text{POS}_{(\alpha, \beta)}(C)) \log P(C^c | \text{POS}_{(\alpha, \beta)}(C)), \\ \Delta_N(\alpha, \beta) &= H(\pi_C | \text{NEG}_{(\alpha, \beta)}(C)) &= -P(C | \text{NEG}_{(\alpha, \beta)}(C)) \log P(C | \text{NEG}_{(\alpha, \beta)}(C)) \\ & & \quad -P(C^c | \text{NEG}_{(\alpha, \beta)}(C)) \log P(C^c | \text{NEG}_{(\alpha, \beta)}(C)), \\ \Delta_B(\alpha, \beta) &= H(\pi_C | \text{BND}_{(\alpha, \beta)}(C)) &= -P(C | \text{BND}_{(\alpha, \beta)}(C)) \log P(C | \text{BND}_{(\alpha, \beta)}(C)) \\ & & \quad -P(C^c | \text{BND}_{(\alpha, \beta)}(C)) \log P(C^c | \text{BND}_{(\alpha, \beta)}(C)), \end{aligned}$$

Deng, X. F., & Yao, Y. Y., (2012). An information-theoretic interpretation of thresholds in PRS. In: (RSCTC'12)



Overall Uncertainty of the PRS

- The overall uncertainty is determined as the weighted average.

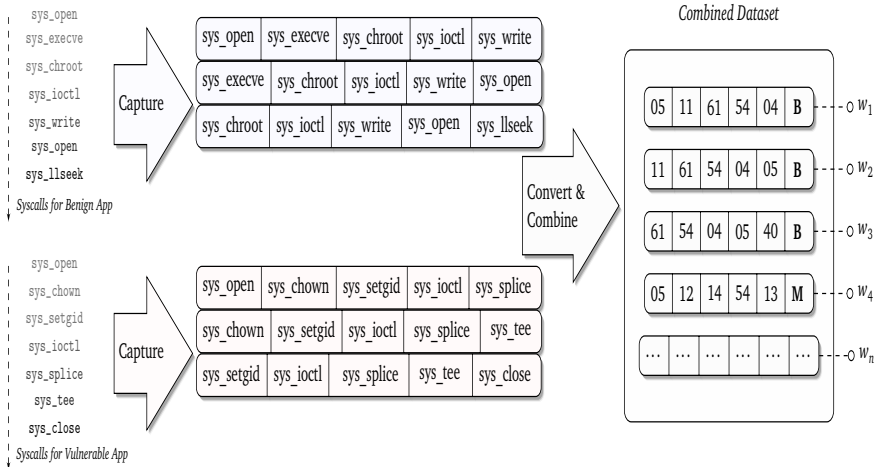
$$\begin{aligned}\Delta(\alpha, \beta) &= P(\text{POS}_{(\alpha, \beta)}(C)) * \Delta_P(\alpha, \beta) + \\ &P(\text{NEG}_{(\alpha, \beta)}(C)) * \Delta_N(\alpha, \beta) + \\ &P(\text{BND}_{(\alpha, \beta)}(C)) * \Delta_B(\alpha, \beta)\end{aligned}$$

- Configuring the thresholds to decrease the uncertainty of a particular region may increase the uncertainty of some other region.
- Consider optimization of the above equation based on (α, β) thresholds.

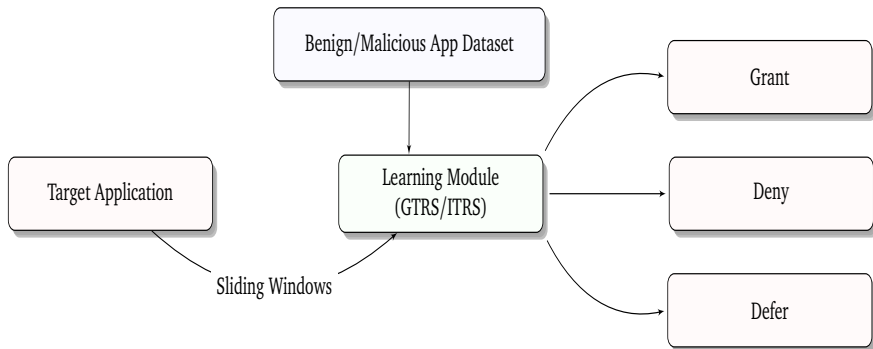
An Architecture for Malware Analysis with Three-way Decisions

- We propose an architecture for malware analysis with three-way decisions.
 - The architecture is based on capturing and analysing the system call sequences of applications.
 - These system call sequences are converted to chunks using sliding window.
 - Each of these chunks are used as a row (object) of an information table.
 - Three-way decision models are then trained on the information table and three-way decisions are obtained.

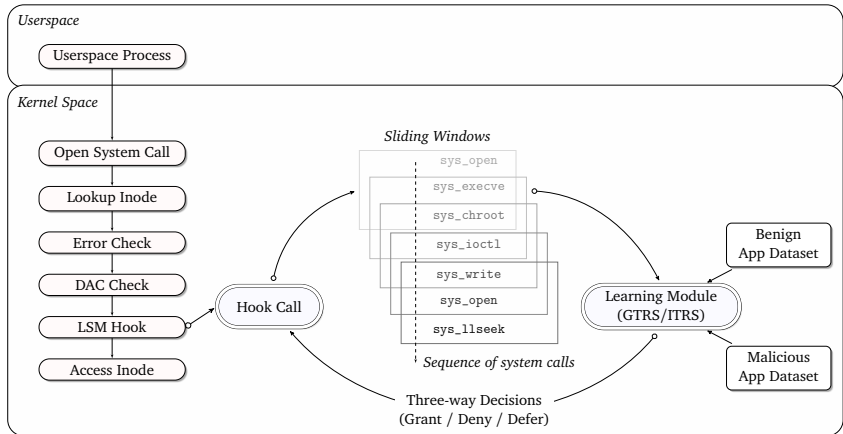
Sliding Windows for Capturing System Call Sequences



High Level View of the Proposed Architecture



Detailed View of the Proposed Architecture



System Call Sequences with Class Labels

Window	s_1	s_2	s_3	s_4	s_5	Behavior	Window	s_1	s_2	s_3	s_4	s_5	Behavior
w_1	106	106	106	106	106	M	w_2	125	5	5	3	90	B
w_3	125	5	5	3	90	M	w_4	106	5	90	6	5	B
w_5	106	106	106	106	106	M	w_6	125	5	5	3	90	M
w_7	3	90	90	90	6	B	w_8	3	90	90	90	6	M
w_9	106	5	90	6	5	M	w_{10}	125	5	5	3	90	M
w_{11}	125	5	5	3	90	M	w_{12}	5	108	3	19	6	B
w_{13}	108	3	19	6	33	B	w_{14}	108	3	19	6	33	M
w_{15}	3	90	90	90	6	M	w_{16}	3	90	90	90	6	M
w_{17}	106	5	90	6	5	M	w_{18}	3	6	5	108	3	B
w_{19}	5	108	3	19	6	B	w_{20}	5	108	3	19	6	M
w_{21}	125	5	5	3	90	B	w_{22}	45	45	5	108	45	B
w_{23}	45	45	5	108	45	M	w_{24}	3	6	5	108	3	B
w_{25}	3	6	5	108	3	B	w_{26}	3	6	5	108	3	M
w_{27}	125	5	5	3	90	B	w_{28}	6	5	108	3	19	B
w_{29}	3	6	5	108	3	M	w_{30}	45	45	5	108	45	B
w_{31}	5	108	3	19	6	B	w_{32}	6	5	108	3	19	B

- Considering the above table as information table.
- The rows corresponds to sliding windows of system calls.
- The columns contain the system call no.s as defined in OS.
- The last column represent the associated behaviour.

- The following equivalence classes may be created based on the information table.

$X_1 = \{w_1, w_5\}$	$X_2 = \{w_7, w_8, w_{15}, w_{16}\}$
$X_3 = \{w_4, w_9, w_{17}\}$	$X_4 = \{w_2, w_3, w_6, w_{10}, w_{11}, w_{21}, w_{27}\}$
$X_5 = \{w_{13}, w_{14}\}$	$X_6 = \{w_{18}, w_{24}, w_{25}, w_{26}, w_{29}\}$
$X_7 = \{w_{22}, w_{23}, w_{30}\}$	$X_8 = \{w_{12}, w_{19}, w_{20}, w_{31}\}$
$X_9 = \{w_{28}, w_{32}\}$	

- Considering the concept of interest as Behaviour = M,
 - The conditional probabilities of the concept with X_i is, $P(C|X_i) = P(\text{Behaviour} = M|X_i) = \frac{|\text{Behaviour} = M \cap X_i|}{|X_i|}$.
 - The probability of X_i 's are given by $P(X_i) = \frac{|X_i|}{|U|}$.

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9
$P(C X_i)$	1.0	0.75	0.67	0.57	0.5	0.4	0.33	0.25	0.0
$P(X_i)$	0.0625	0.125	0.09375	0.21875	0.0625	0.15625	0.09375	0.125	0.0625

Three-way Decisions based on GTRS

- Considering a game between Accuracy and Generality.
 - These can be defined as,

$$Accuracy(\alpha, \beta) = \frac{|(\text{POS}_{(\alpha, \beta)}(C) \cap C) \cup (\text{NEG}_{(\alpha, \beta)}(C) \cap C^c)|}{|\text{POS}_{(\alpha, \beta)}(C) \cup \text{NEG}_{(\alpha, \beta)}(C)|},$$

$$Generality(\alpha, \beta) = \frac{|\text{POS}_{(\alpha, \beta)}(C) \cup \text{NEG}_{(\alpha, \beta)}(C)|}{|U|},$$

- Three types of strategies were formulated for each player.
 - s_1 (α_{\downarrow} , decrease of α by 20%),
 - s_2 (β_{\uparrow} , increase of β by 20%),
 - s_3 ($\alpha_{\downarrow}\beta_{\uparrow}$, decrease of α and increase of β by 20%).
- The payoffs are based on the measure of accuracy and generality.
 - $u_A(s_m, s_n) = Accuracy(\alpha, \beta)$.
 - $u_G(s_m, s_n) = Generality(\alpha, \beta)$.

Payoff Table for the Game

		G		
		$s_1 = \alpha_{\downarrow}$ = 20% dec. α	$s_2 = \beta_{\uparrow}$ = 20% inc. β	$s_3 = \alpha_{\downarrow}\beta_{\uparrow}$ = 20% (dec. α & inc. β)
A	$s_1 = \alpha_{\downarrow}$ = 20% dec. α	(0.82,0.34)	(1.0,0.13)	(0.82,0.34)
	$s_2 = \beta_{\uparrow}$ = 20% inc. β	(1.0,0.13)	(0.75,0.50)	(0.75,0.50)
	$s_3 = \alpha_{\downarrow}\beta_{\uparrow} = 20\%$ (dec. α & inc. β)	(0.82,0.34)	(0.75,0.50)	(0.74,0.72)

- The cell with bold font represents the game solution.
- The corresponding thresholds are $(\alpha, \beta) = (0.6, 0.2)$.
- These thresholds can be used to induce three-way decisions.

Three-way Decisions based on ITRS

- ITRS determine thresholds for three-way decisions based on minimization of the overall uncertainty.
 - The uncertainty of a particular region, say positive region may be determined as,

$$\Delta P(\alpha, \beta) = H(\pi_C | \text{POS}_{(\alpha, \beta)}(C)) = -P(C | \text{POS}_{(\alpha, \beta)}(C)) \log P(C | \text{POS}_{(\alpha, \beta)}(C)) \\ - P(C^c | \text{POS}_{(\alpha, \beta)}(C)) \log P(C^c | \text{POS}_{(\alpha, \beta)}(C)),$$

- For the considered information table,

$$P(C | \text{POS}_{(1,0)}(C)) = \frac{\sum_{i=1}^1 P(C | X_i) * P(X_i)}{\sum_{i=1}^1 P(X_i)} = \frac{1 * 0.0625}{0.0625} = 1.0 \quad (1)$$

- The probability $P(C^c | \text{POS}_{(1,0)}(C)) = 1 - P(C | \text{POS}_{(1,0)}(C)) = 1 - 1 = 0$.
- Therefore, $H(\pi_C | \text{POS}_{(1,0)}(C)) = -1 * \log 1 - (0 * \log 0) = 0$.
- The uncertainty for other regions can be similarly obtained.

Three-way Decisions based on ITRS

- To determine a minimum value of uncertainty, we consider the domains of thresholds based on majority oriented model given by $0 \leq \beta < 0.5 \leq \alpha \leq 1.0$.
 - This leads to the domain of α , i.e., $D_\alpha = \{1.0, 0.7, 0.6, 0.5\}$ and domain of β , i.e., given by $D_\beta = \{0.0, 0.3, 0.4\}$.

	$\alpha = 1.0$	$\alpha = 0.7$	$\alpha = 0.6$	$\alpha = 0.5$
$\beta = 0.0$	0.875	0.8680	0.8607	0.8606
$\beta = 0.3$	0.8682	0.8688	0.8661	0.8768
$\beta = 0.4$	0.8544	0.8665	0.8704	0.8937

- The cell with bold font represents the minimum uncertainty which corresponds to $(\alpha, \beta) = (0.6, 0.2)$.
- These thresholds can be used to induce three-way decisions.

Conclusion and Future Work

- **Conclusions**
- We consider a three-way decision making approach to malware analysis.
- Essential change is the deferment decision option.
 - Useful for decision making under low quality information.
- An architecture for malware analysis with three-way decisions.
- A demonstrative example suggest that use of the suggested approach.
- **Future Work.**
- Deployment of the three-way approach on the production systems.
 - This will enable us to measure efficiency on large scale.
- Examination in the context of latest technology.
 - Smartphones, tablets and other digital devices.

Questions?

JingTao Yao, PhD

姚静涛

Professor

DEPARTMENT OF COMPUTER SCIENCE

University
of Regina

Regina, Saskatchewan
Canada S4S 0A2

phone: 306.585.4071

fax: 306.585.4745

email: jtyao@cs.uregina.ca

<http://www.cs.uregina.ca/~jtyao>

