

Game Theory for Cyber Deception

Jeffrey Pawlick

2018 Conference on Decision and
Game Theory for Security

October 29, 2018

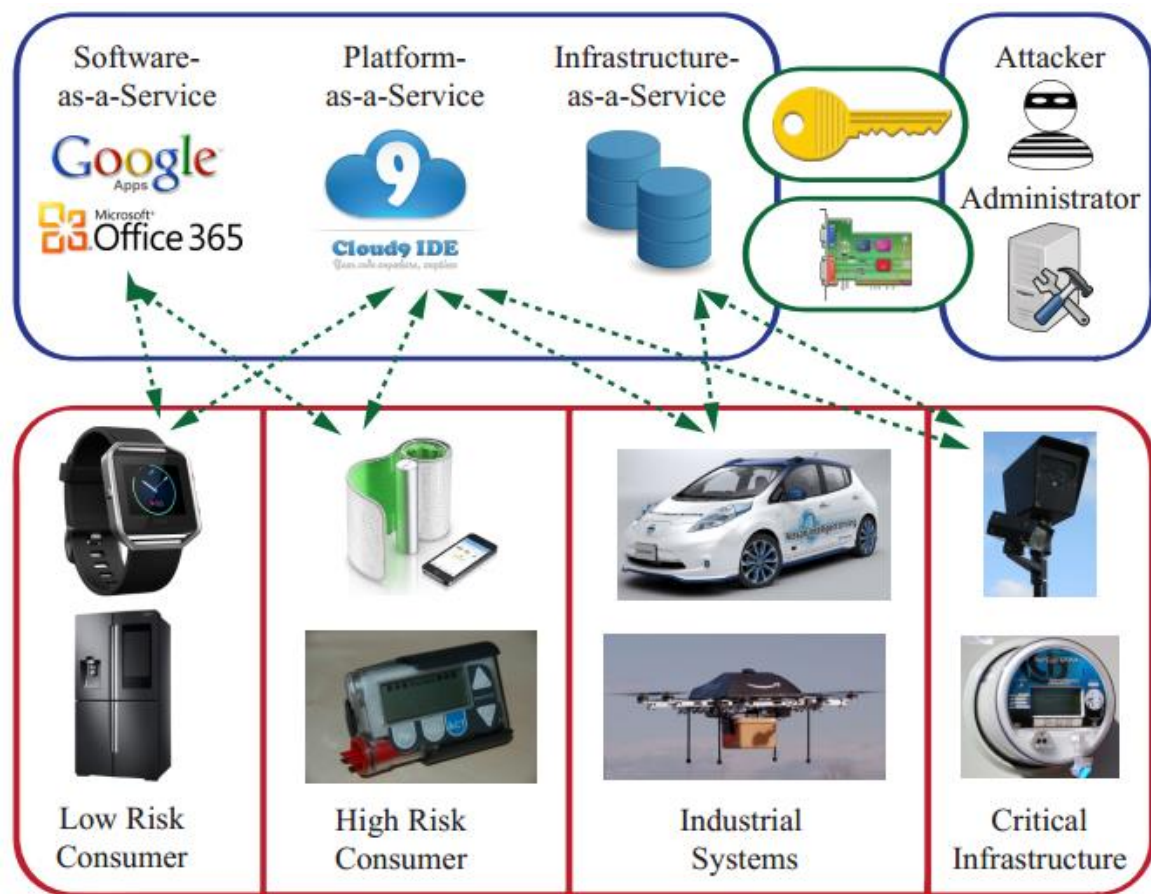
Supported in part by an NSF IGERT
grant through the Center for
Interdisciplinary Studies in Security and
Privacy (CRISSP) at New York University



Laboratory for Agile and
Resilient Complex Systems

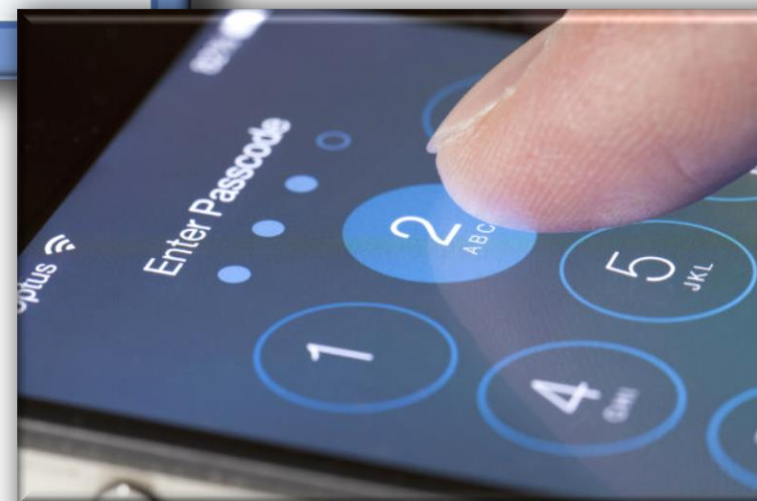
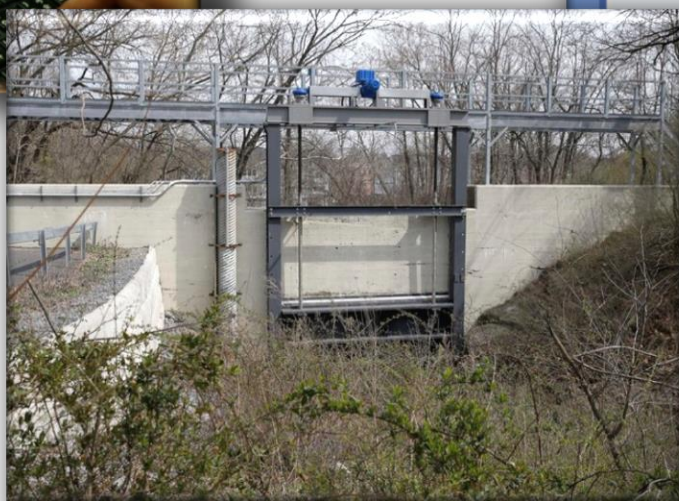
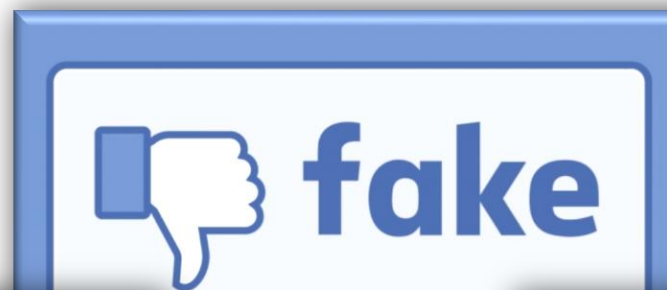


Increasing Connectivity



- Controlled systems: biological, social, physical, communication
- Cloud: offers SaaS, PaaS, IaaS
- *Internet of controlled things (IoCT):*
 - Internet of things (IoT) +
 - Wireless sensor-actuator networks (WSAN) +
 - Cyber-physical systems (CPS)

Deception Online and in the IoT



Towards a Science of Deception

- Knowledge that is wholistic, essential, transferable, quantitative
- Prediction that is relevant for law, policy, and business
- Mechanism design that is relevant for economics and technology

Game theory

Control theory

Machine
Learning

Estimation and
Detection

Signal
Processing



Outline of the Slide

- 1) Introduction
- 2) Taxonomy of defensive deception
- 3) Signaling games for mimetic deception
- 4) Strategic trust for counter-deception
- 5) Future challenges



Deception in Economics, Psychology, and Privacy

Imposing cognitive load to elicit cues to deceit: inducing the reverse order technique naturally

Aldert Vrij^{a*}, Sharon Leal^a, Samantha Mann^a and Ronald Fisher^b

^aDepartment of Psychology, University of Portsmouth, Portsmouth, UK; ^bDepartment of Psychology, Florida International University, Miami, FL, USA

(Received 10 January 2010; final version received 12 August 2010)

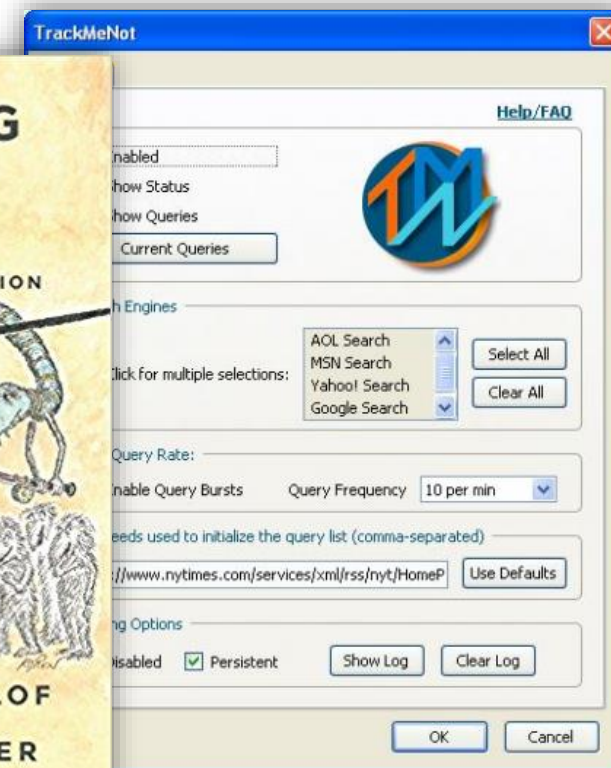
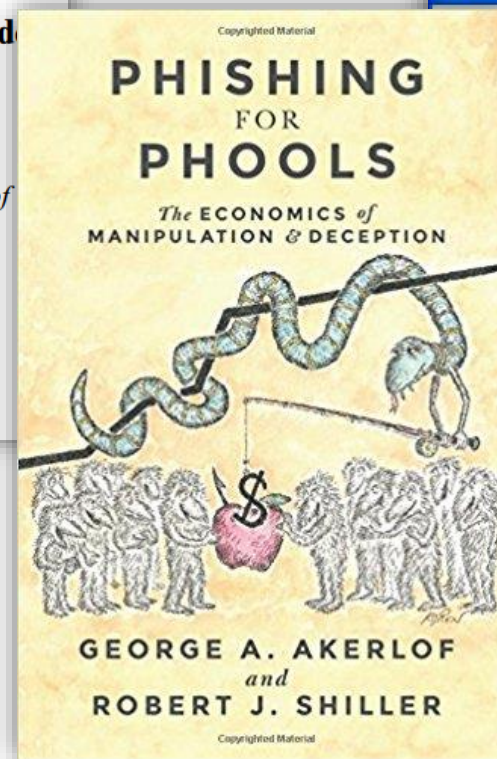
Deception: The Role of Consequences

By URI GNEEZY*

Deception is part of many economic interactions. Business people, politicians, diplomats, lawyers, and students in the experimental laboratory who make use of private information do not always do so honestly. This observation indicates that behavior often rejects the moral approach to deception. As St. Augustine wrote, "To me, however, it seems certain that every lie

sociated with lying per se. This assumption is very useful in many economic models. Consider contract theory, where it is assumed that without an explicit contract, neither side will fulfill its respective obligations. For example, George Akerlof's (1970) paper on asymmetric information and the market for lemons assumes that sellers of used cars will always lie if it is in their

between lying
in reverse
recall their
ment 1, 31



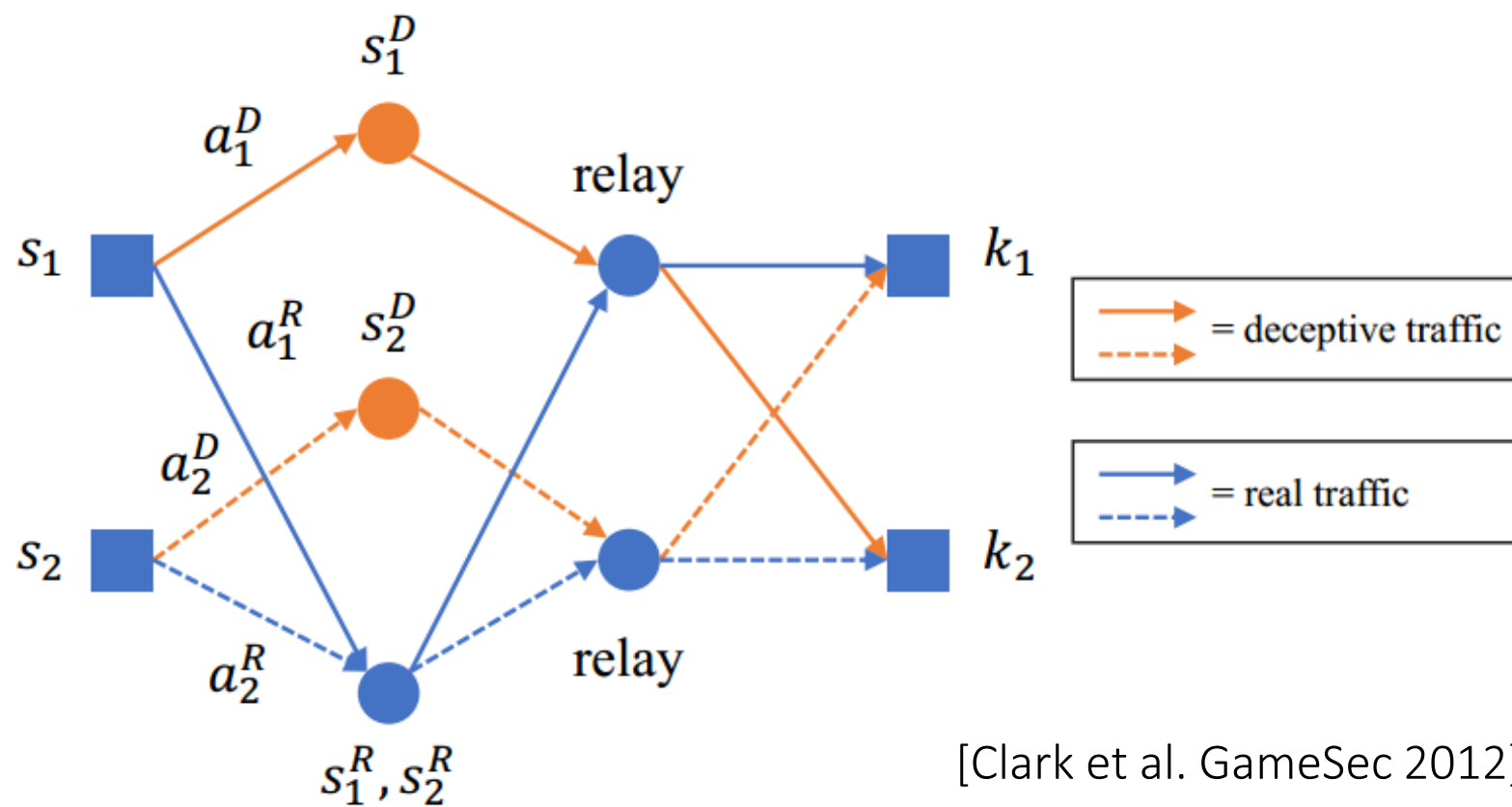
Defensive Deception in Cybersecurity & Privacy

Authors and Year	Game-Theoretic Model	Application Domain
Chessa et al. 2015	Nash	Info. Privacy
Shorki 2015	Stackelberg	Info. Privacy
Alvim et al. 2017	Nash	Info. Privacy
Theodorakopoulos et al. 2014	Bayesian Stackelberg	Location Privacy
Rass et al. 2017	Nash	General Security
Clark et al. 2015	Stackelberg & Nash	Network Security
Zhu & Basar 2013	Nash	Network Security
Feng et al. 2017	Stackelberg	General Security
Clark et al. 2012	Stackelberg	Network Security
Zhu et al. 2012	Stackelberg	Network Security
Pawlick & Zhu 2016	Stackelberg	Info. Privacy
Pawlick & Zhu 2017a	Mean-Field	Info. Privacy
Zhang et al. 2010	Best Response	Anonymity

Authors and Year	Game-Theoretic Model	Application Domain
Freudiger et al. 2009	Bayesian Nash	Location Privacy
Lu et al. 2012	Nash	Location Privacy
Carroll & Grosu 2011	Signaling	Network Security
Mohammadi et al. 2016	Signaling	Social Networks
Píbil et al. 2012	Bayesian Nash	Network Security
Kiekintveld et al. 2015	Bayesian Nash	Network Security
Pawlick & Zhu 2015	Signaling with Evidence	Network Security
Pawlick & Zhu 2017b	Signaling with Evidence	Network Security
Zhuang et al. 2010	Multi-Period Signaling	General Security
Durkota et al. 2015	Stackelberg with MDP	Network Security
Horák et al. 2017	One-Sided POMDP	Network Security

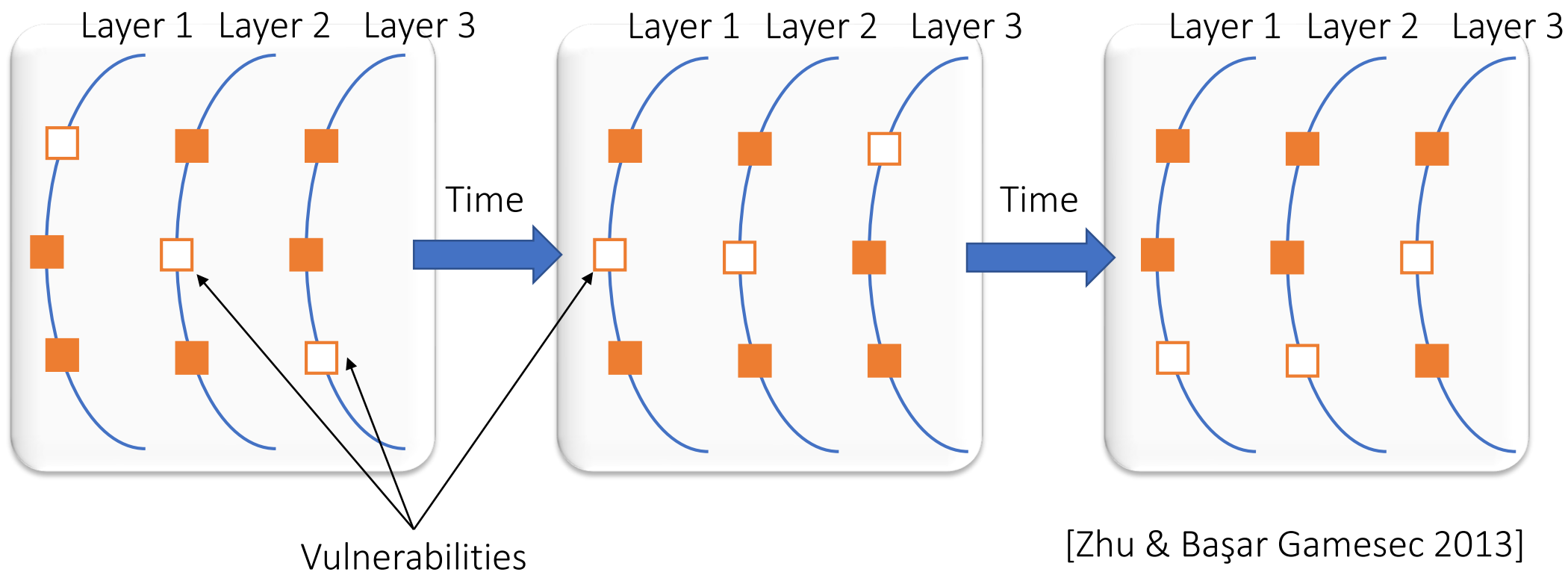


Obfuscation Example



[Clark et al. GameSec 2012]

Moving Target Defense Example



Definition of Types of Deception

- To deceive $\stackrel{\text{def}}{=}$ to intentionally cause another agent to acquire or continue to have a false belief, or to be prevented from acquiring or cease to have a true belief [Mahon 2016].

Two different types of deception: Creating a false belief vs. preventing the acquisition of a true belief?

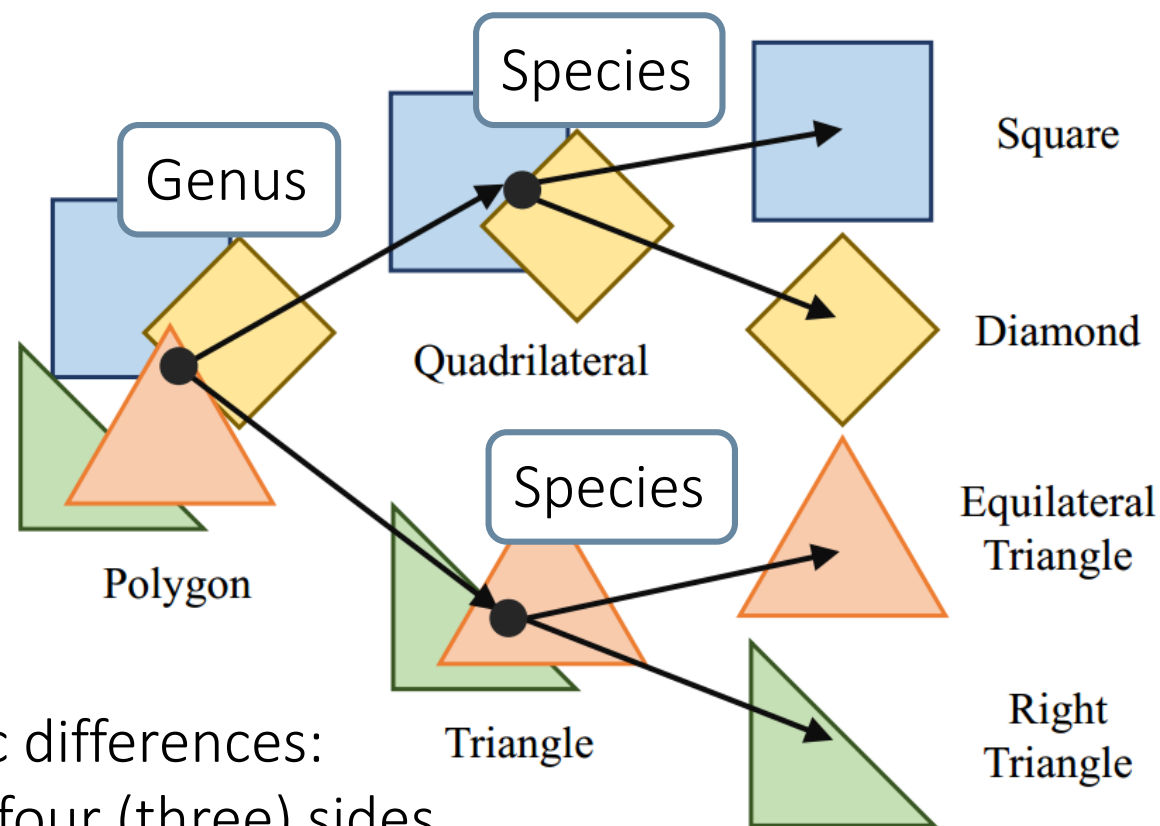
Where do “perturbation,” “obfuscation,” and “moving target defense” fit?

Goal of our taxonomy: to rigorously define types of defensive deception for cybersecurity and privacy.



Defensive Deception in Cybersecurity & Privacy

There is a need for “the construction of a common language and a set of basic concepts about which the security community can develop a shared understanding” [U.S. Dept. of Defense].



Definition of Species of Deception

Specific differences: incentives, actors, actions, and time-horizon

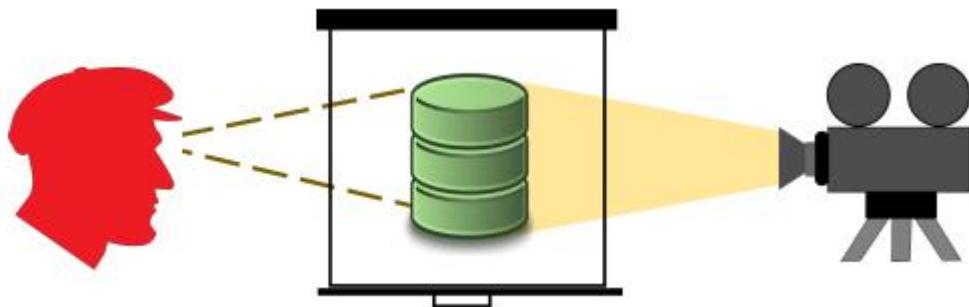
- Incentives / utility functions – what is the goal of the deception?
- Actors / players – who are the participants in the deception?
- Actions – what means are used to achieve the deception?
- Time-horizon – what is the duration of the deception?



Incentives: What is the Purpose of the Deception?

To deceive ^{def} to intentionally cause another agent to acquire or continue to have a false belief, or to be prevented from acquiring or cease to have a true belief.

Mimetic Deception

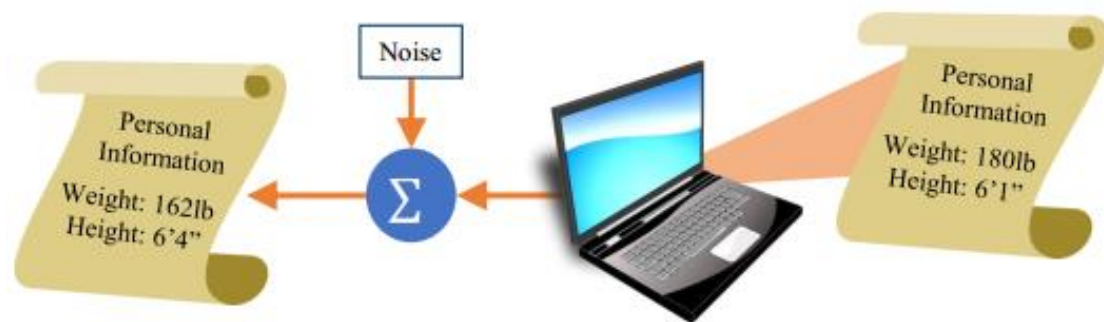


Cryptic Deception



Actors: Who are the Participants in the Deception?

Intensive Deception



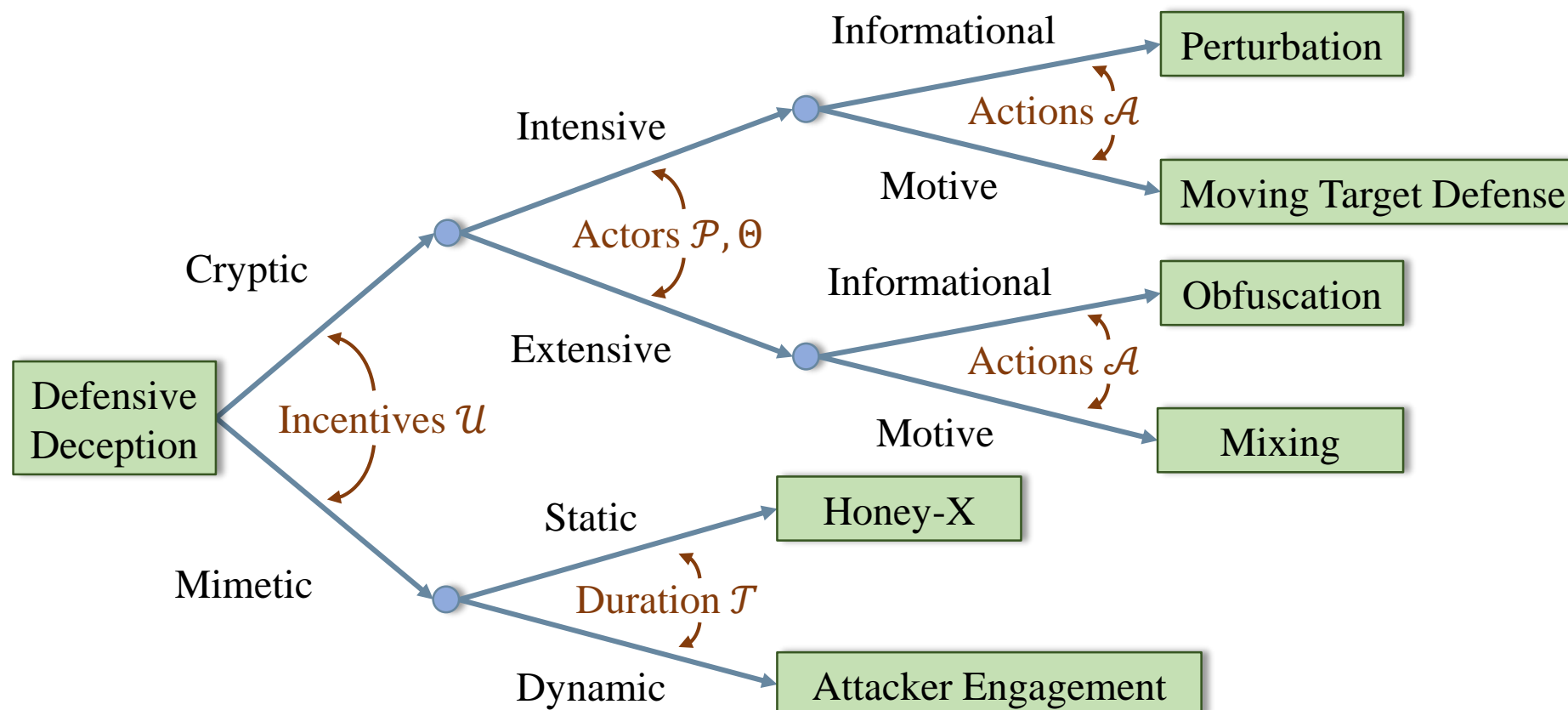
- Single target / actor

Extensive Deception



- Multiple targets / actors

Taxonomy Based on Game Theoretic Principles



Trends in Papers on Cryptic Deception

Motive	Feng. et al. 2017 – Stack. (uses MDP) Clark et al. 2015 – Stack. (with leader mixed-strategies)	Zhang et al. 2010 – Best response in multiple stages (user-adversary) Freudiger et al. 2009 – Nash (user-user) Lu et al. 2012 – Nash (user-user)
	Rass et al. 2017 – Nash (mixed-strategies) Zhu and Başar 2013 – Nash (mixed-strategies)	
Informational	Chessa et al. 2015 – Nash (user-user) Alvim et al. 2017 – Nash (utilities are <i>a priori</i>)	Clark et al. 2012 – Stack. (user-adversary) Zhu et al. 2012 – Stack (user-adversary) Pawlick and Zhu 2016 – Stack. (user-adversary) Pawlick and Zhu 2017a – Stack (user-adversary) and Mean-Field Game (user-user)
	Shorki 2015 – Stack. (user-adversary) Theodorakopoulos et al. 2017 – Stack. (user-adversary)	
	Intrinsic	Extrinsic



Trends in Papers on Mimetic Deception

<p>Carroll and Grosu 2011 – Signaling Mohammadi et al. 2016 – Signaling Pawlick and Zhu 2015 – Signaling Pawlick and Zhu 2017b – Signaling</p>	<p>Zhuang et al. 2010 – Multi-Period Signaling</p> <p>Durkota et al. 2015 – Stackelberg (with Markov decision process)</p> <p>Horák et al. 2017 – One-sided partially-observable stochastic game</p>
Static	Dynamic

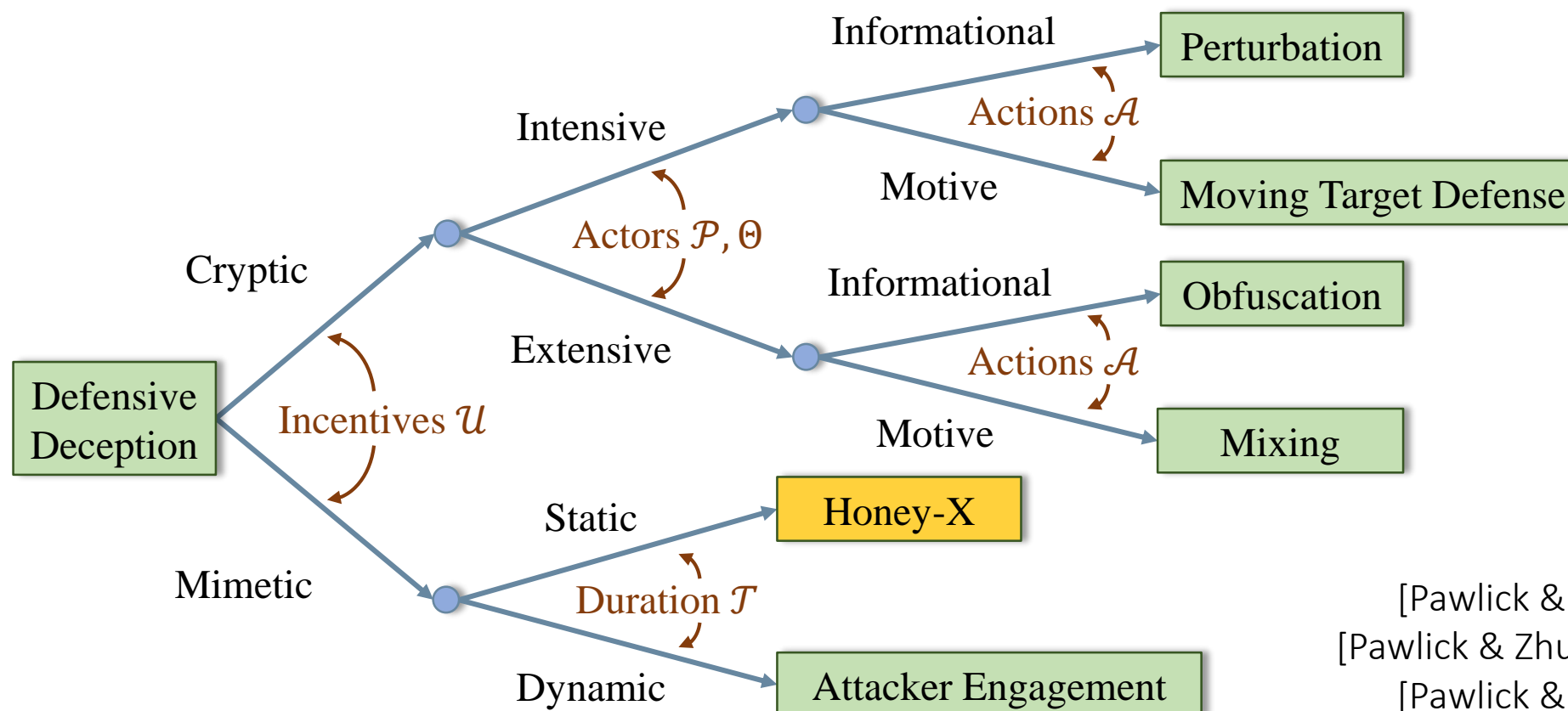


Opportunities for Future Research

- Theoretical Advances – Most papers use Nash or Stackelberg equilibrium. There are few dynamic games or studies of dynamic problems (which might arise in the IoT).
- Test implementations – These exist in physical security, but not in cybersecurity. Why?
 - Wariness of security through obscurity? But we have quantified guarantees
 - High demand for security analysts? Collaboration will be necessary.
 - Challenges of interdisciplinary security? Problems require cognitive science, psychology, sub-rationality, models of attacker preferences, criminology, etc.
- Mimetic Deception – Literature lacks it. Why? Randomization is straightforward? Law?



Taxonomy Based on Game Theoretic Principles



[Pawlick & Zhu, WEIS 2015]

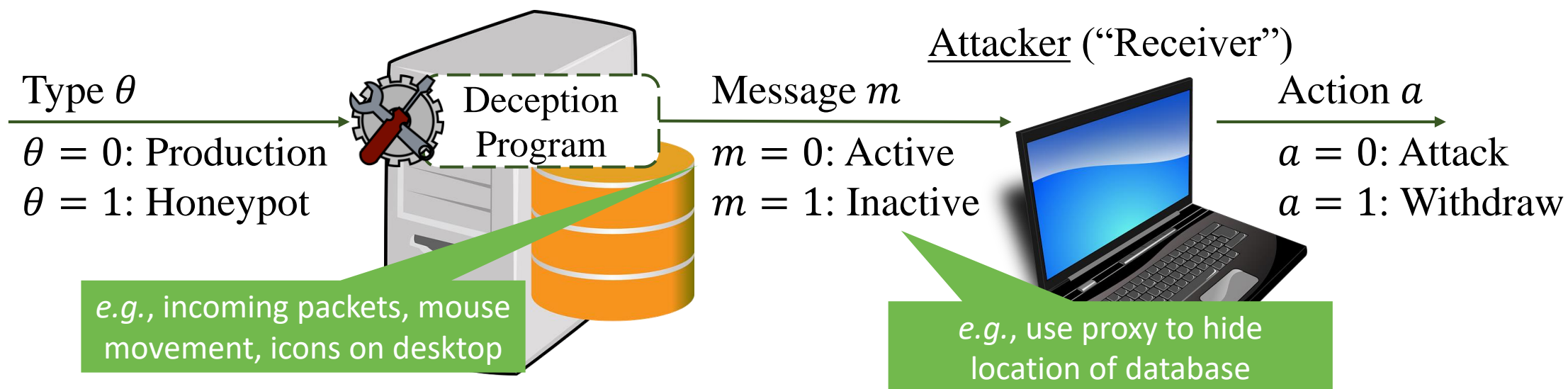
[Pawlick & Zhu, IEEE CNS 2017]

[Pawlick & Zhu, WEIS 2018]

Mimesis and Modeling Belief

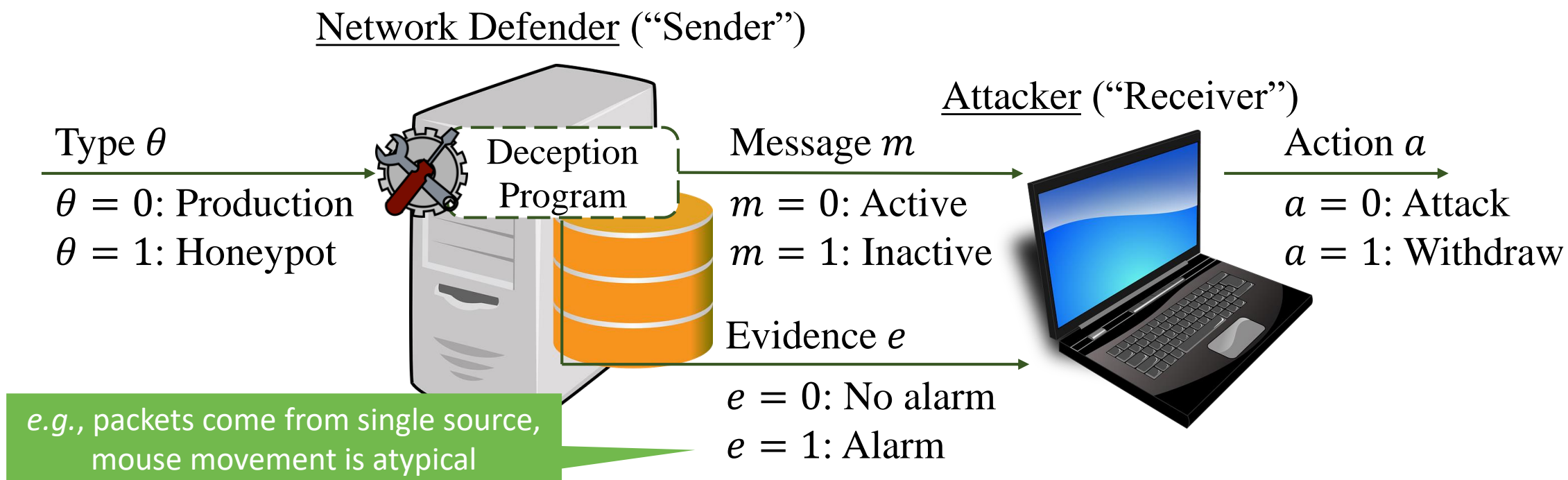
- Signaling games model belief [Lewis 1969, Crawford & Sobel 1982].

Network Defender (“Sender”)



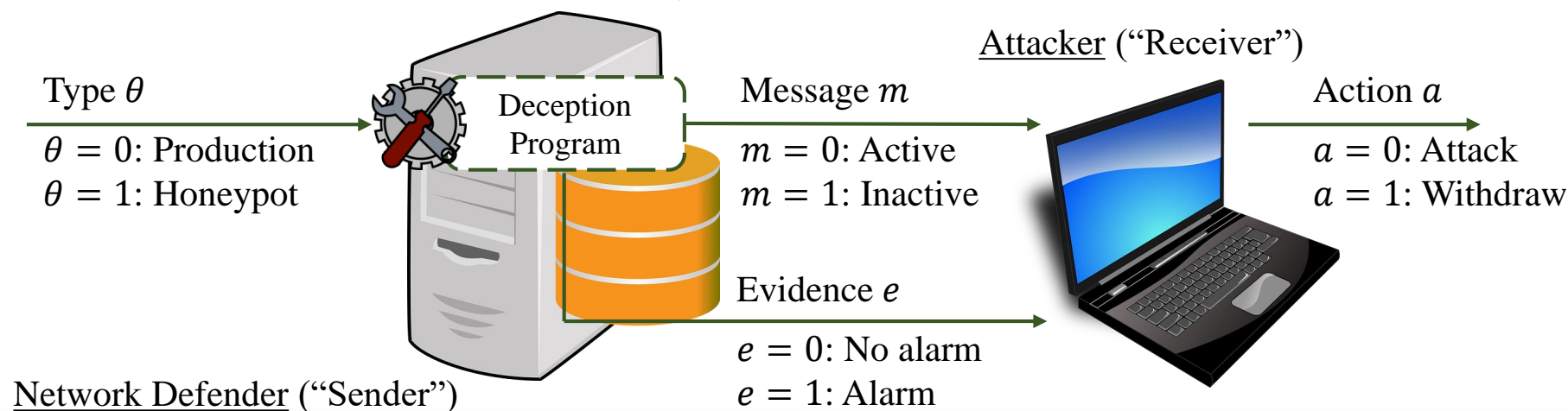
Mimesis and Modeling Belief

- But “deception program” may leak evidence.



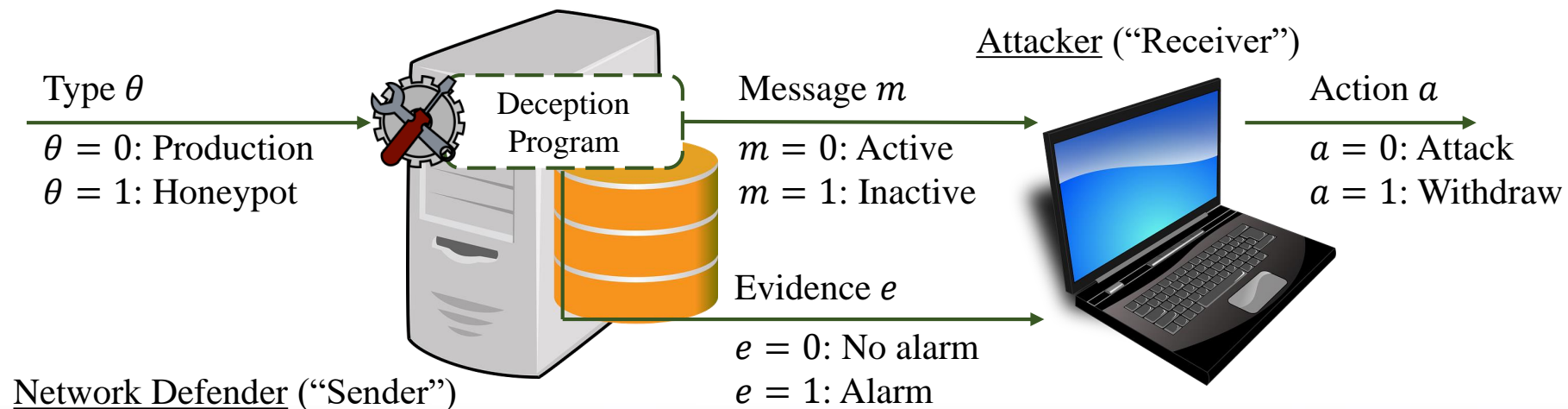
Mixed Strategies, Belief, and Expected Utility

- Attacker has prior belief of system type θ with probability (wp) $p(\theta)$.
- Defender chooses activity level m w.p. $\sigma^S(m | \theta)$.
- Defender leaks evidence e wp $\lambda(e | \theta, m)$.
- Defender forms belief $\mu^R(\theta | m, e)$ and chooses action a wp $\sigma^R(a | m, e)$.



Mixed Strategies, Belief, and Expected Utility

- System of type θ has an expected utility of $U^S(\sigma^S, \sigma^R | \theta)$.
- Attacker that observes activity level m and evidence e has an expected utility of $\sum_{\theta \in \Theta} \mu^R(\theta | m, e) U^R(\sigma^R | \theta, m, e)$.



Perfect Bayesian Nash Equilibrium

A PBNE is a strategy profile $(\sigma^{S*}, \sigma^{R*})$ and posterior beliefs $\mu^R(\theta | m, e)$ such that:

$\forall \theta \in \Theta,$

$$\sigma^{S*} \in \operatorname{argmax}_{\sigma^S \in \Gamma^S} U^S(\sigma^S, \sigma^{R*} | \theta),$$

$\forall m \in M, e \in \mathbb{E}V,$

$$\sigma^{R*} \in \operatorname{argmax}_{\sigma^R \in \Gamma^R} \sum_{\theta \in \Theta} \mu^R(\theta | m, e) U^R(\sigma^R | \theta, m, e),$$

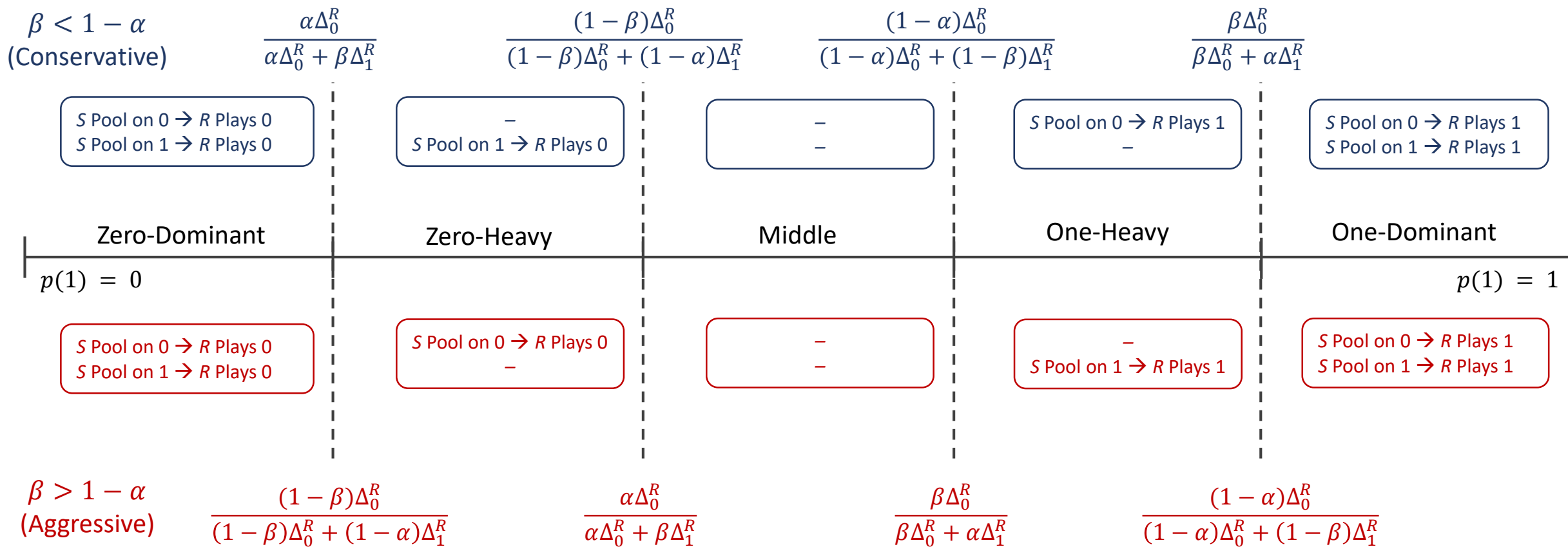
and

$$\mu^R(\theta | m, e) = \frac{\lambda(e | \theta, m) \sigma^S(m | \theta) p(\theta)}{\sum_{\tilde{\theta} \in \Theta} \lambda(e | \tilde{\theta}, m) \sigma^S(m | \tilde{\theta}) p(\tilde{\theta})},$$

when that fraction is defined.



Equilibrium Regions



Partially-Separating Equilibria in the Middle Regime

Theorem (Aggressive Detectors). For $\beta > 1 - \alpha$, within the Middle regime, there exists a PBNE in which

$$\sigma^{S^*}(m = 1 | \theta = 0) = \frac{\bar{\alpha}\bar{\beta}\Delta_1^R}{(\bar{\alpha}^2 - \bar{\beta}^2)\Delta_0^R} \left(\frac{p(1)}{1-p(1)} \right) - \frac{\bar{\beta}^2}{\bar{\alpha}^2 - \bar{\beta}^2},$$

$$\sigma^{S^*}(m = 1 | \theta = 1) = \frac{\bar{\alpha}^2}{\bar{\alpha}^2 - \bar{\beta}^2} - \frac{\bar{\alpha}\bar{\beta}\Delta_0^R}{(\bar{\alpha}^2 - \bar{\beta}^2)\Delta_1^R} \left(\frac{1-p(1)}{p(1)} \right),$$

and

$$\sigma^{R^*}(a = 1 | m = 0, e = 0) = 0, \quad \sigma^{R^*}(a = 1 | m = 0, e = 1) = \frac{1}{\alpha + \beta},$$

$$\sigma^{R^*}(a = 1 | m = 1, e = 0) = 1, \quad \sigma^{R^*}(a = 1 | m = 1, e = 1) = \frac{\alpha + \beta - 1}{\alpha + \beta},$$

and the beliefs are computed by Bayes' Law in all cases. Here $\bar{\mathbf{x}} = \mathbf{1} - \mathbf{x}$.



Partially-Separating Equilibria in the Middle Regime

Theorem (Conservative Detectors). For $\beta < 1 - \alpha$, within the Middle regime, there exists a PBNE in which

$$\sigma^{S^*}(m = 1 | \theta = 0) = \frac{\beta^2}{\beta^2 - \alpha^2} - \frac{\alpha\beta\Delta_1^R}{(\beta^2 - \alpha^2)\Delta_0^R} \left(\frac{p(1)}{1 - p(1)} \right),$$

$$\sigma^{S^*}(m = 1 | \theta = 1) = \frac{\alpha\beta\Delta_0^R}{(\beta^2 - \alpha^2)\Delta_1^R} \left(\frac{1 - p(1)}{p(1)} \right) - \frac{\alpha^2}{\beta^2 - \alpha^2},$$

and

$$\sigma^{R^*}(a = 1 | m = 0, e = 0) = \frac{1 - \alpha - \beta}{2 - \alpha - \beta}, \quad \sigma^{R^*}(a = 1 | m = 0, e = 1) = 1,$$

$$\sigma^{R^*}(a = 1 | m = 1, e = 0) = \frac{1}{2 - \alpha - \beta}, \quad \sigma^{R^*}(a = 1 | m = 1, e = 1) = 0,$$

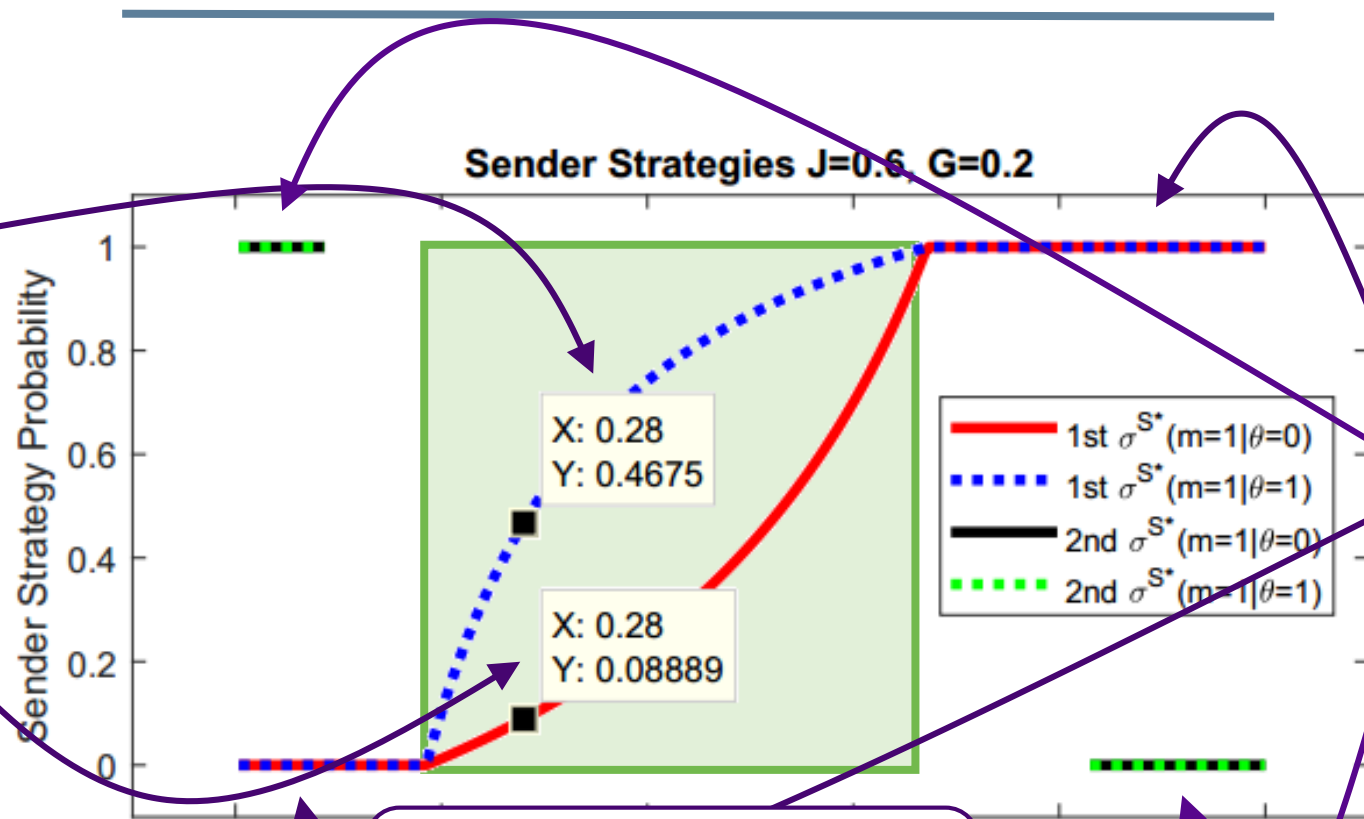
and the beliefs are computed by Bayes' Law in all cases.



Partially-Separating Strategies for S

“Reveal”
honeypot as
inactive wp 0.47.

“Reveal”
production as
active wp 0.91.



Middle Regime:
Partially Separating

Coincident lines:
Pooling
Equilibrium #2



Partially-Separating Strategies for S

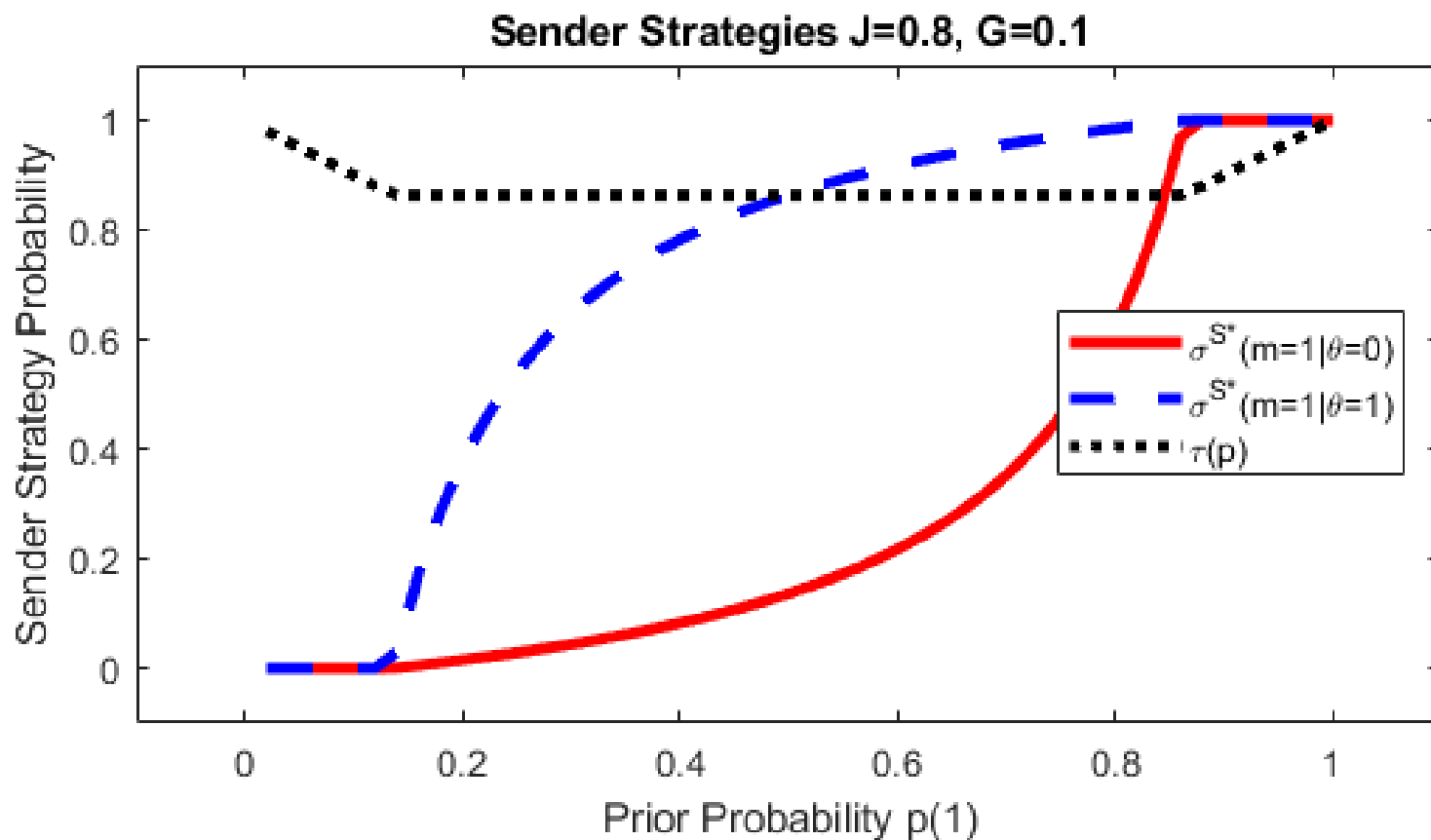
“Reveal”
honeypot as
inactive wp 0.47.

“Reveal”
production as
active wp 0.91.

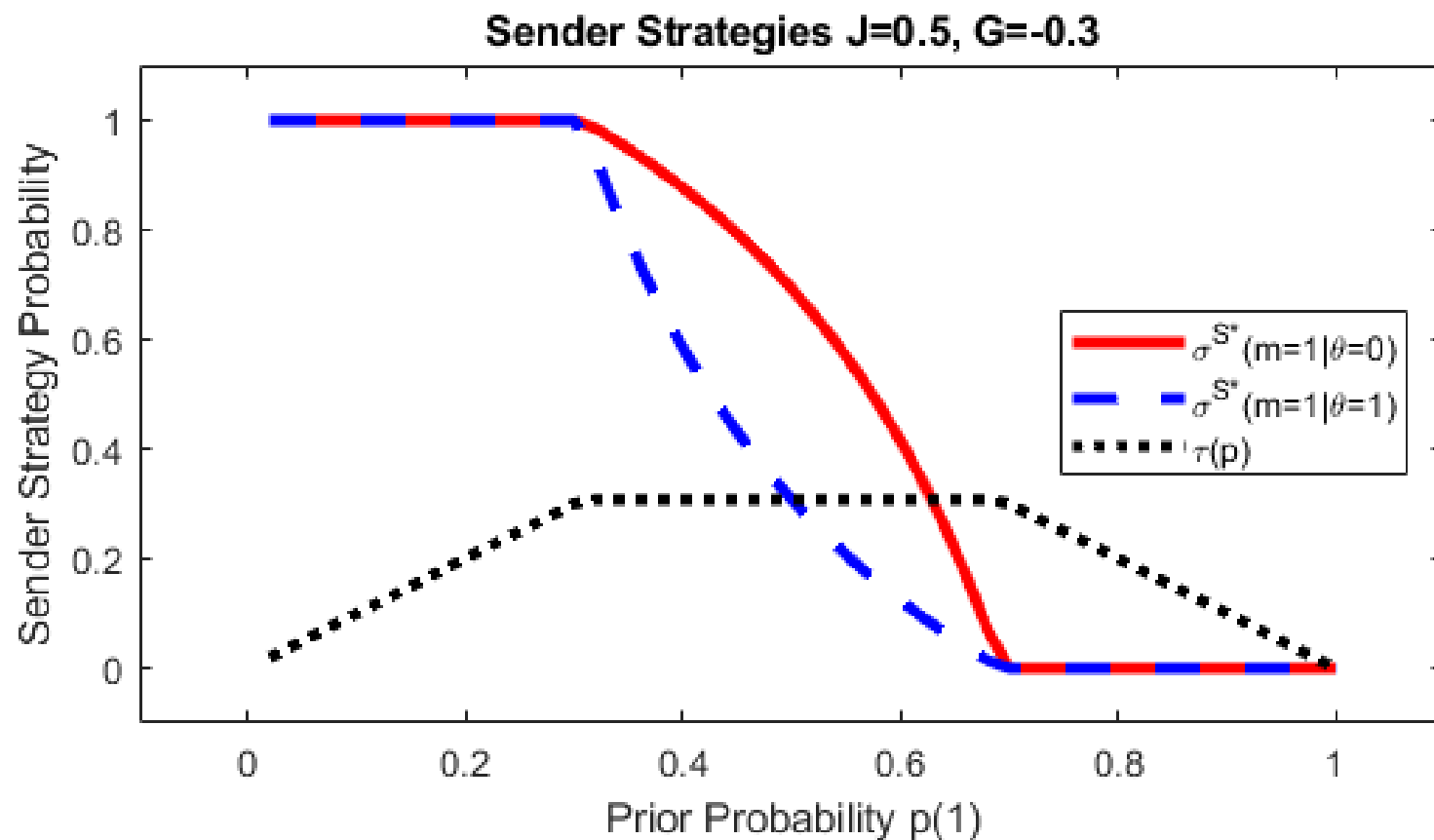
- It is incentive-compatible to reveal the true type with some probability in the Middle regime.
- [Henricks & McAfee 2006] on feints finds information communication due to lying costs.
- [Crawford 2003] on lying finds information communication due to bounded rationality.
- The present model finds information communication due to leakage / evidence.



Comparative Statics: Detector Quality $J = \beta - \alpha$



Comparative Statics: Aggressiveness $G = \beta - (1 - \alpha)$



Truth Induction

Theorem (Truth Induction). Set $\Delta_0^R = \Delta_1^R$. Within regimes that feature unique PBNE, for all $J \in [0,1]$ and for any prior probability $p(\theta)$:

$$\tau(J, G, p) \geq \frac{1}{2} \text{ for } G \in [0,1),$$

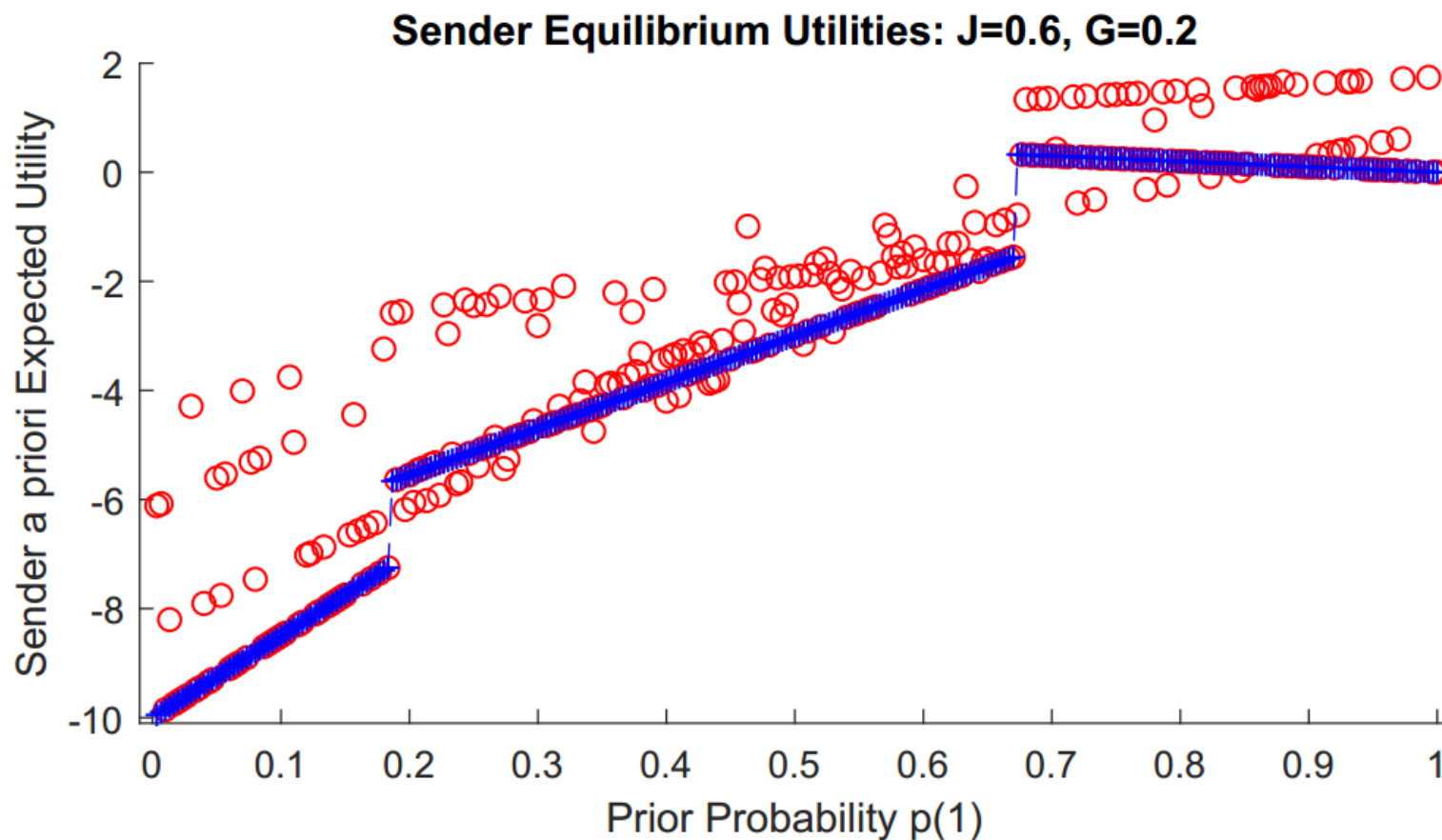
$$\tau(J, G, p) \leq \frac{1}{2} \text{ for } G \in (-1,0],$$

where

$$\tau(J, G, p) \triangleq \sum_{\theta \in \{0,1\}} p(\theta) \sigma^{S^*}(m = \theta \mid \theta; p).$$

Aggressive detectors induce a *truth-telling convention*, while conservative detectors induce a *falsification convention*.

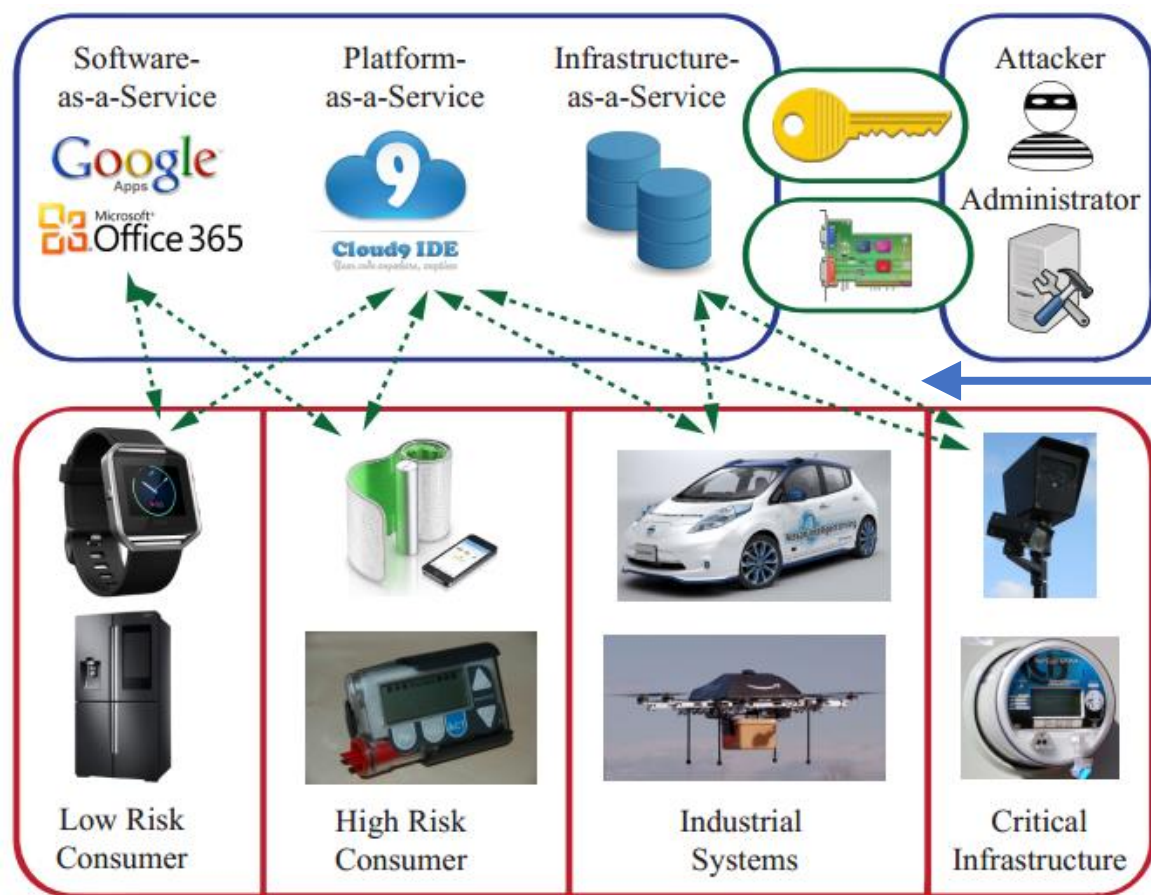
Robustness



S 's expected equilibrium utility usually improves with suboptimal actions of R .

R 's expected equilibrium utility is indifferent to suboptimal actions of S .

Strategic Trust: A Three-Player Interaction



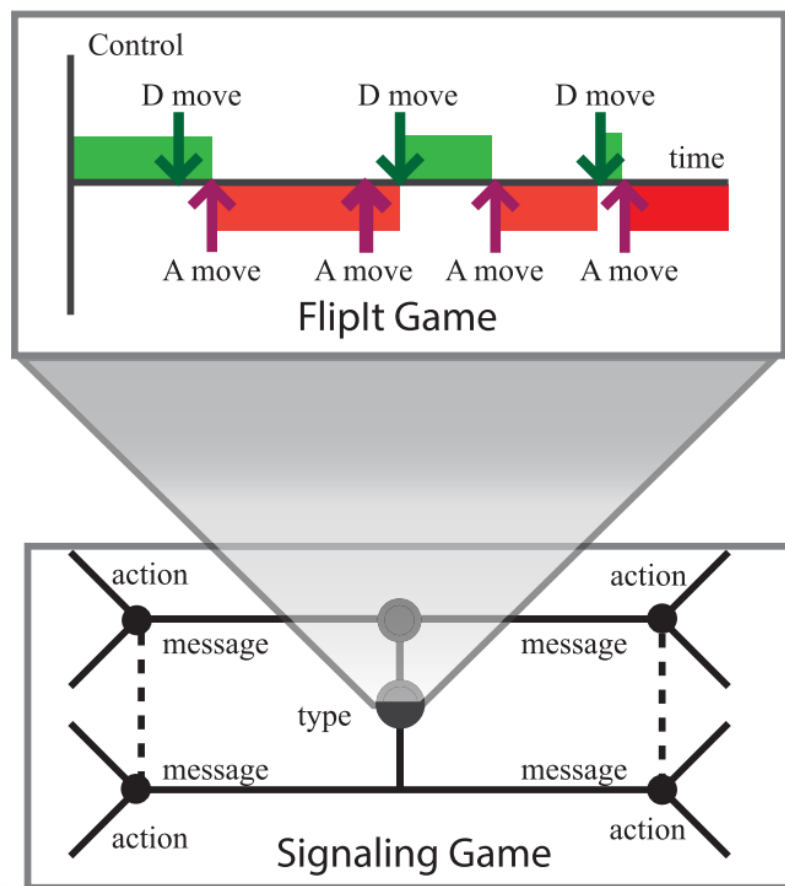
Attacker A and Defender D struggle for control of the cloud (signaling resource).

The winner sends a signal to cloud-enabled device R .

Device R decides whether to trust a possibly compromised cloud.

[Pawlick et al. GameSec 2015],
[Pawlick & Zhu IEEE T-IFS 2016]

Strategic Trust: A Three-Player Interaction



[Bowers et al.
GameSec 2012],

[van Dijk et al.
J Cryptology 2013]

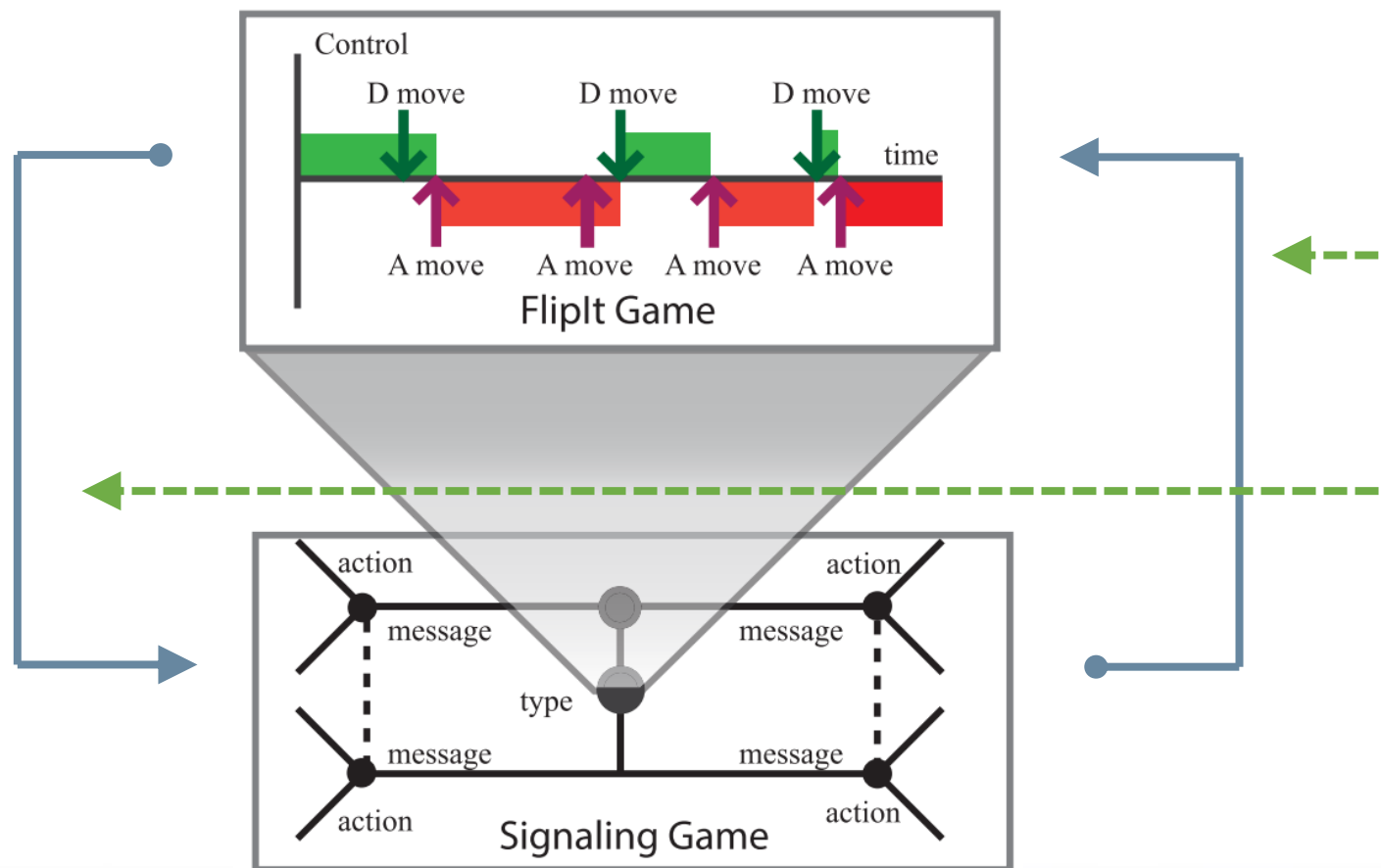
Attacker A and Defender D struggle for control of the cloud (signaling resource).

The winner sends a signal to cloud-enabled device R .

Device R decides whether to trust a possibly compromised cloud.

[Pawlick et al. GameSec 2015],
[Pawlick & Zhu IEEE T-IFS 2016]

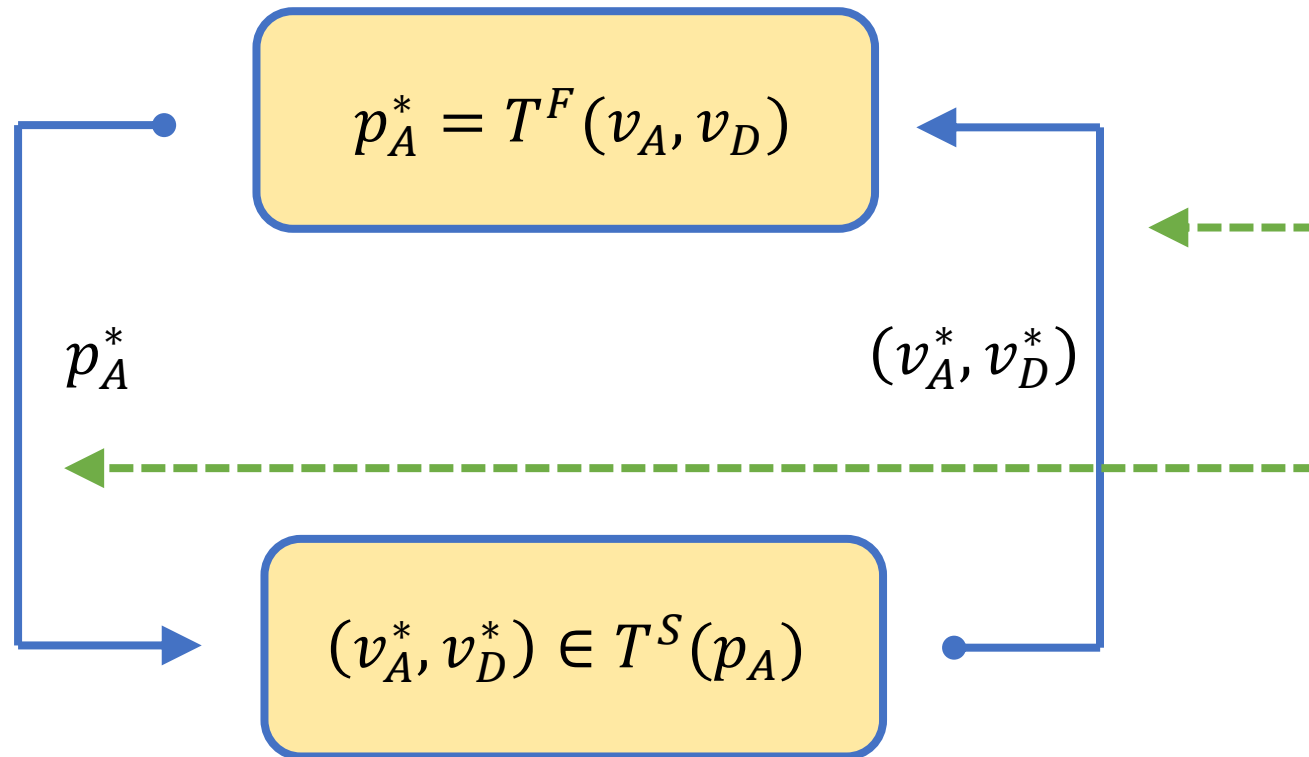
Strategic Trust: A Three-Player Interaction



The value of the cloud in the Fliplt game depends on the equilibrium of the signaling game.

The prior probability of compromise in the signaling game is based on the equilibrium of the Fliplt game.

Strategic Trust: A Three-Player Interaction

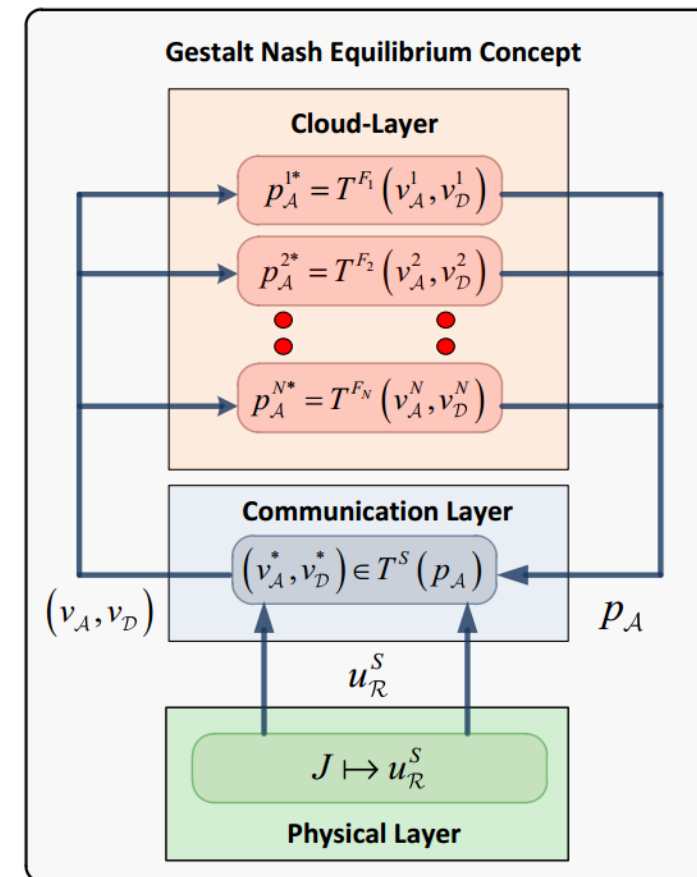


The value of the cloud in the Fliplt game depends on the equilibrium of the signaling game.

The prior probability of compromise in the signaling game is based on the equilibrium of the Fliplt game.

Strategic Trust: A $2N+1$ Player Interaction

- Consider multiple signal sources, each of which can be compromised.
- FlipIt games model attacker-defender interactions at each signal source.
- The device uses a vector signaling game to simultaneously decide which sources to trust.



Gestalt Nash Equilibrium (GNE)

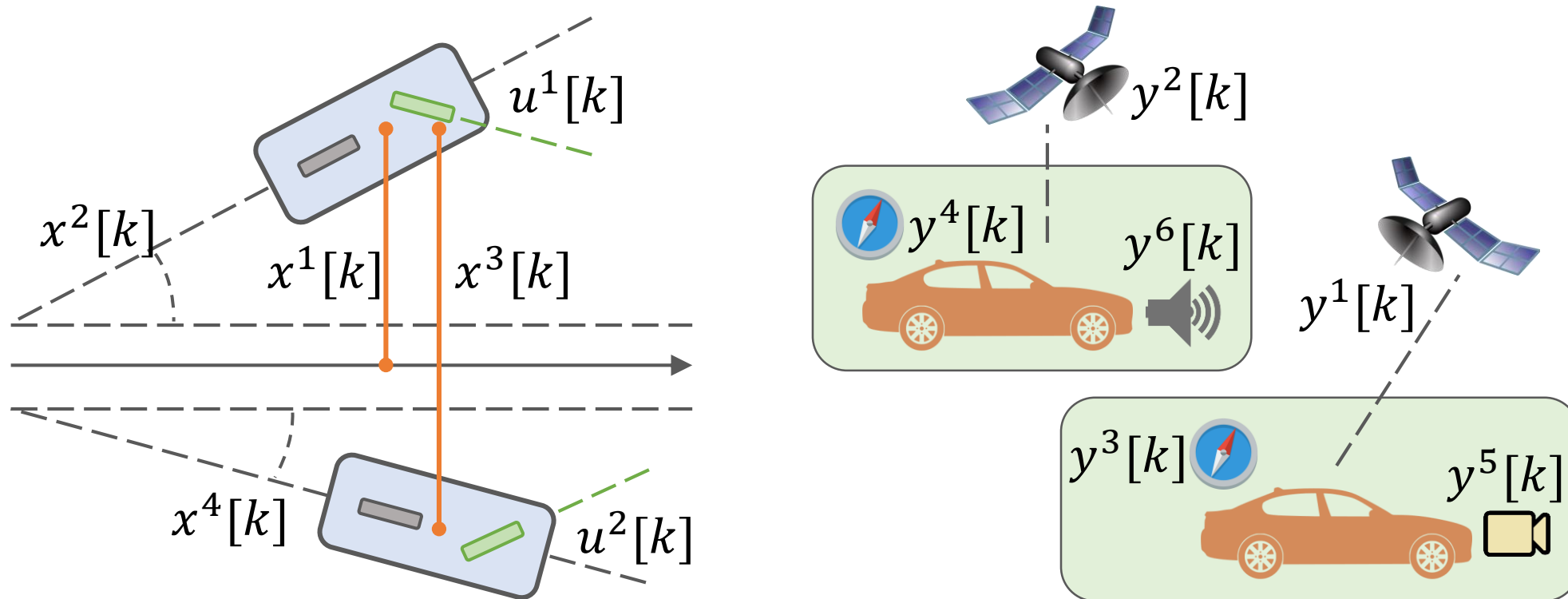
Definition (Gestalt Nash equilibrium). The triple $(p_A^\dagger, v_A^\dagger, v_D^\dagger)$ constitutes a Gestalt Nash equilibrium of the overall game if both of the following equations are satisfied:

$$\forall i \in \{1, \dots, N\}, p_A^{i\dagger} = T^{Fi}(v_A^{i\dagger}, v_D^{i\dagger})$$

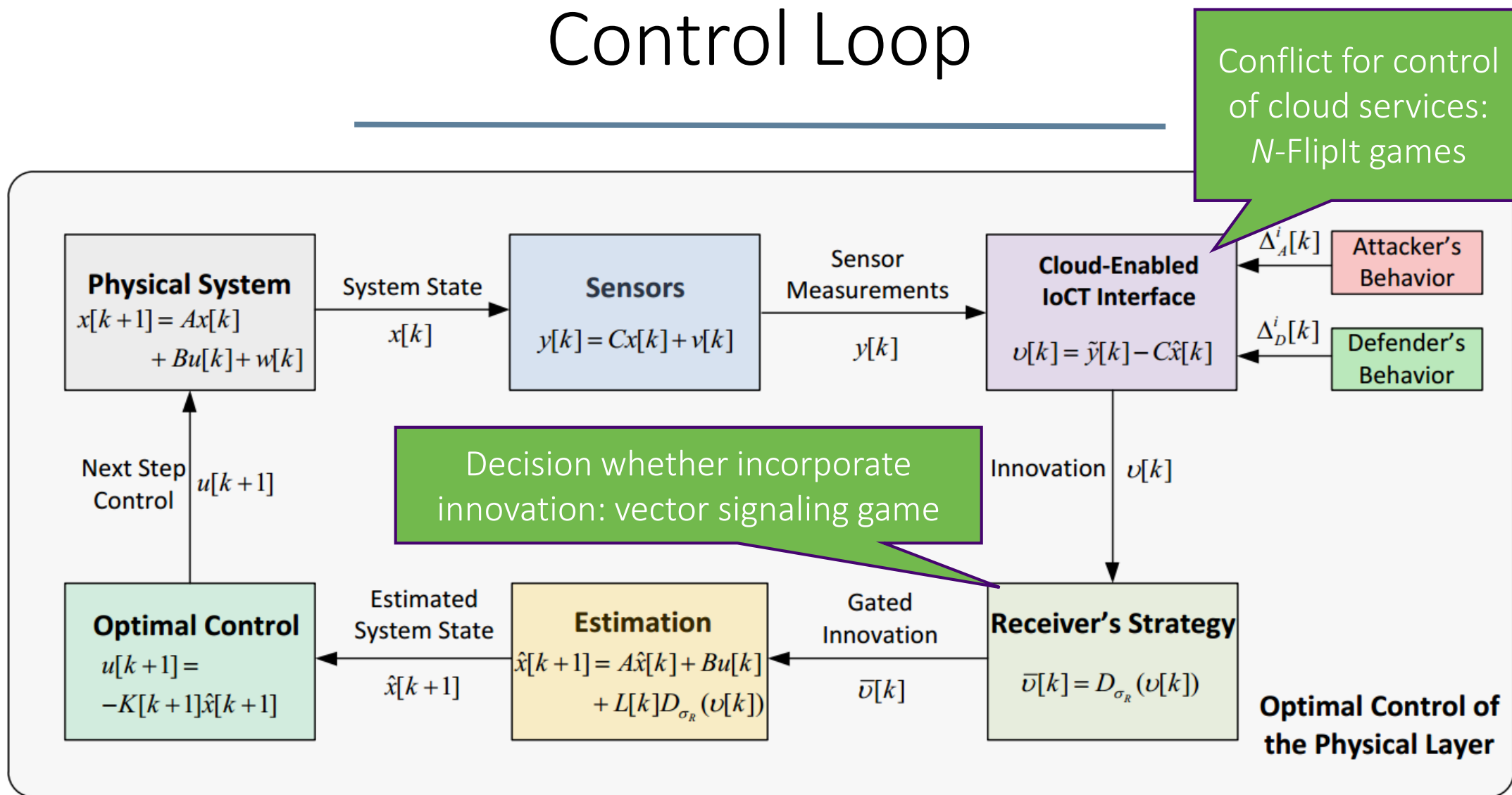
$$\left(\begin{array}{c} \left[v_A^{1\dagger} \right] \\ \vdots \\ \left[v_A^{N\dagger} \right] \end{array}, \begin{array}{c} \left[v_D^{1\dagger} \right] \\ \vdots \\ \left[v_D^{N\dagger} \right] \end{array} \right) \in T^S \left(\begin{array}{c} \left[p_A^{1\dagger} \right] \\ \vdots \\ \left[p_A^{N\dagger} \right] \end{array} \right)$$



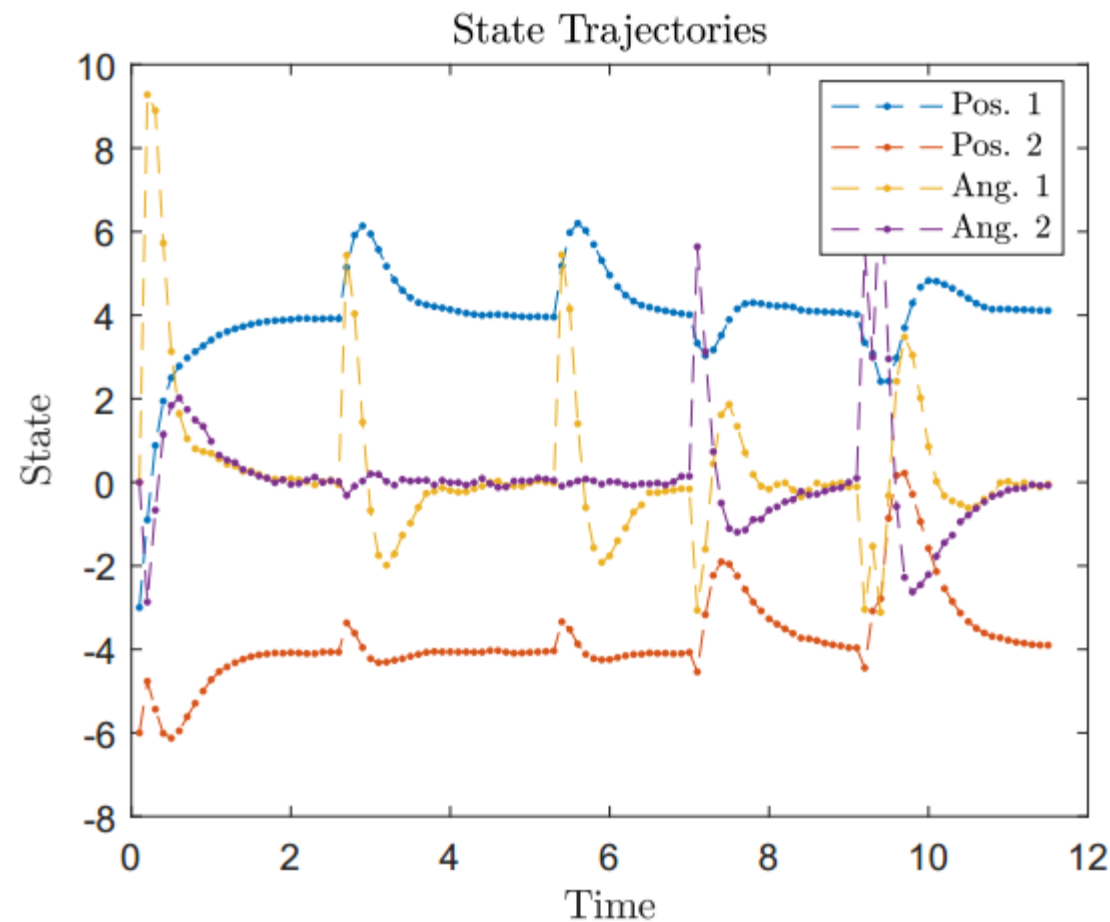
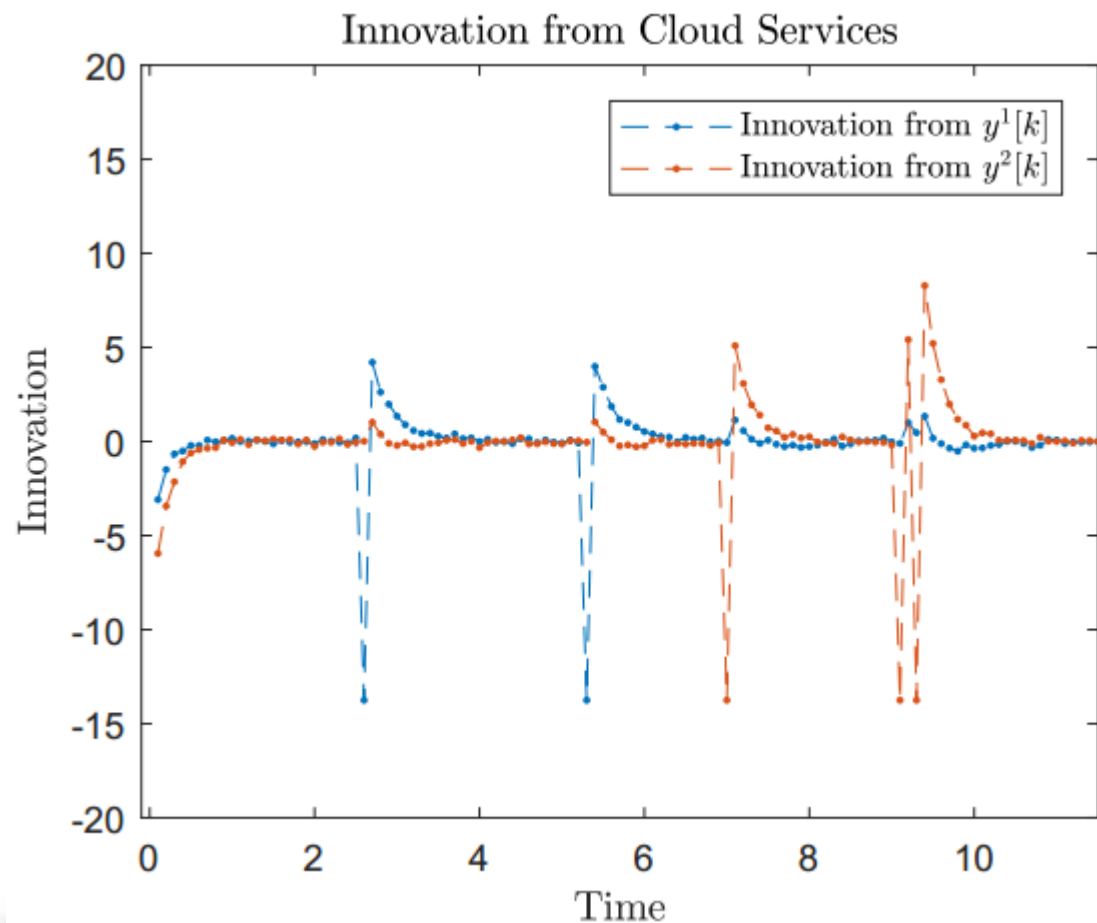
Vehicle Application: States and Measurements



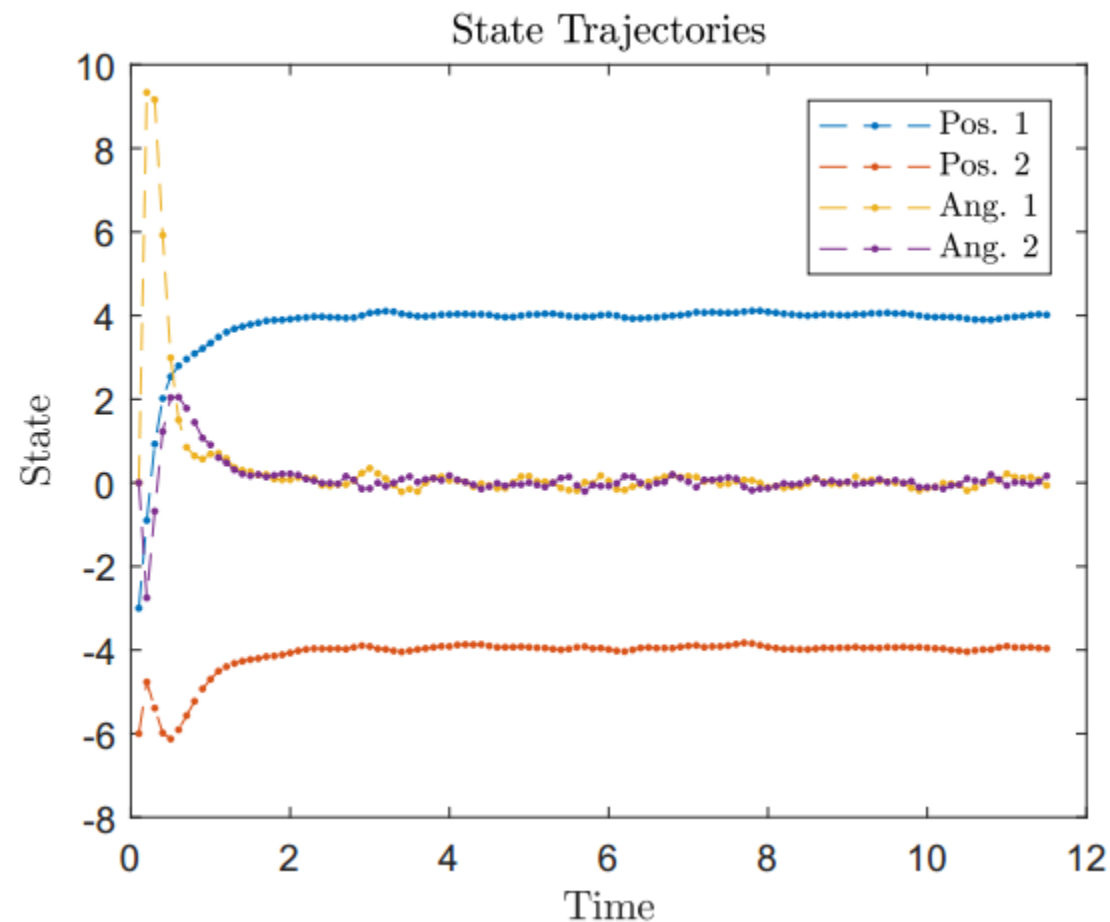
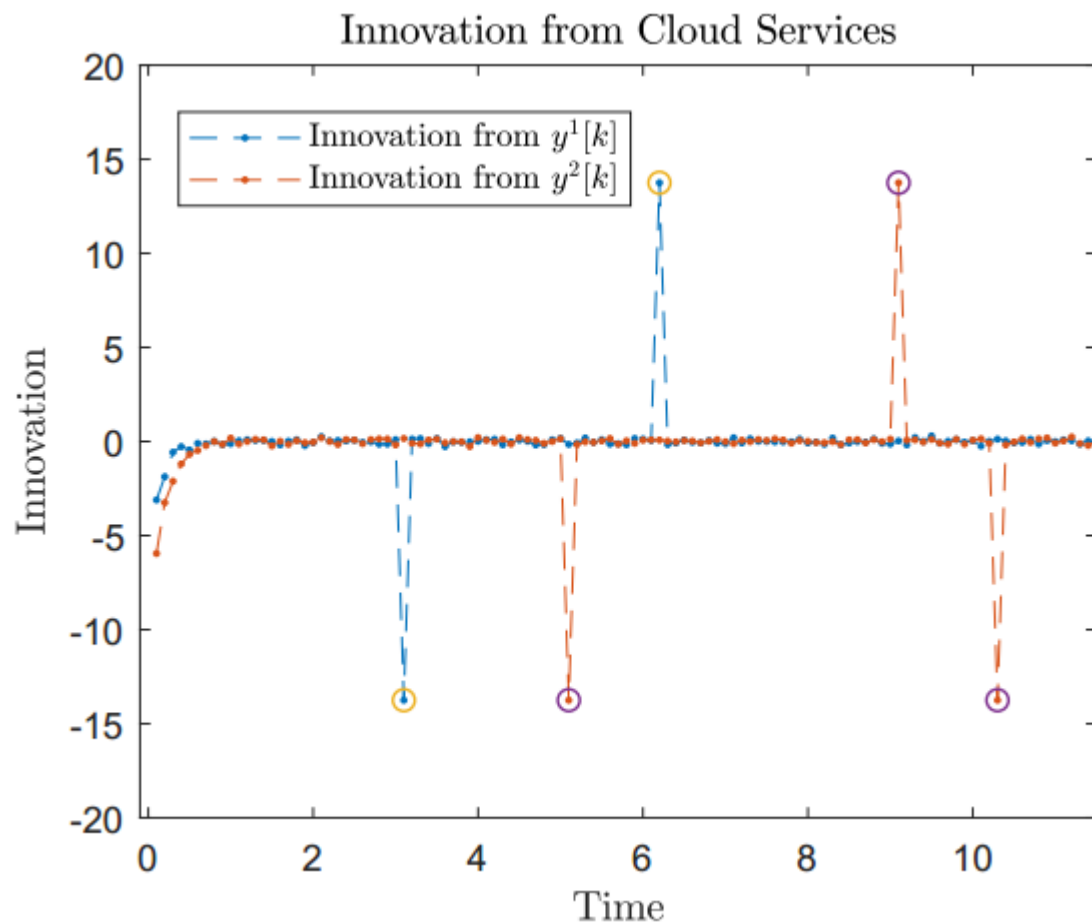
Control Loop



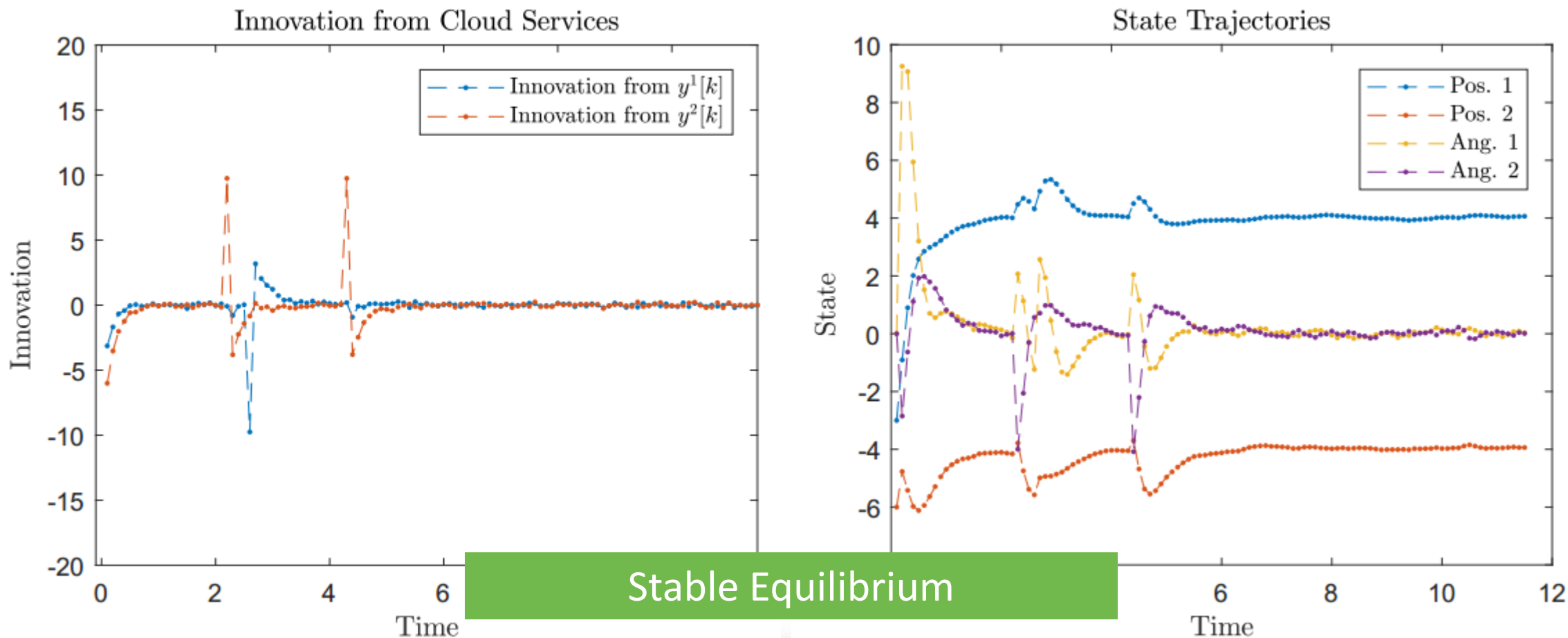
Simulation Results: High Risk, Ungated



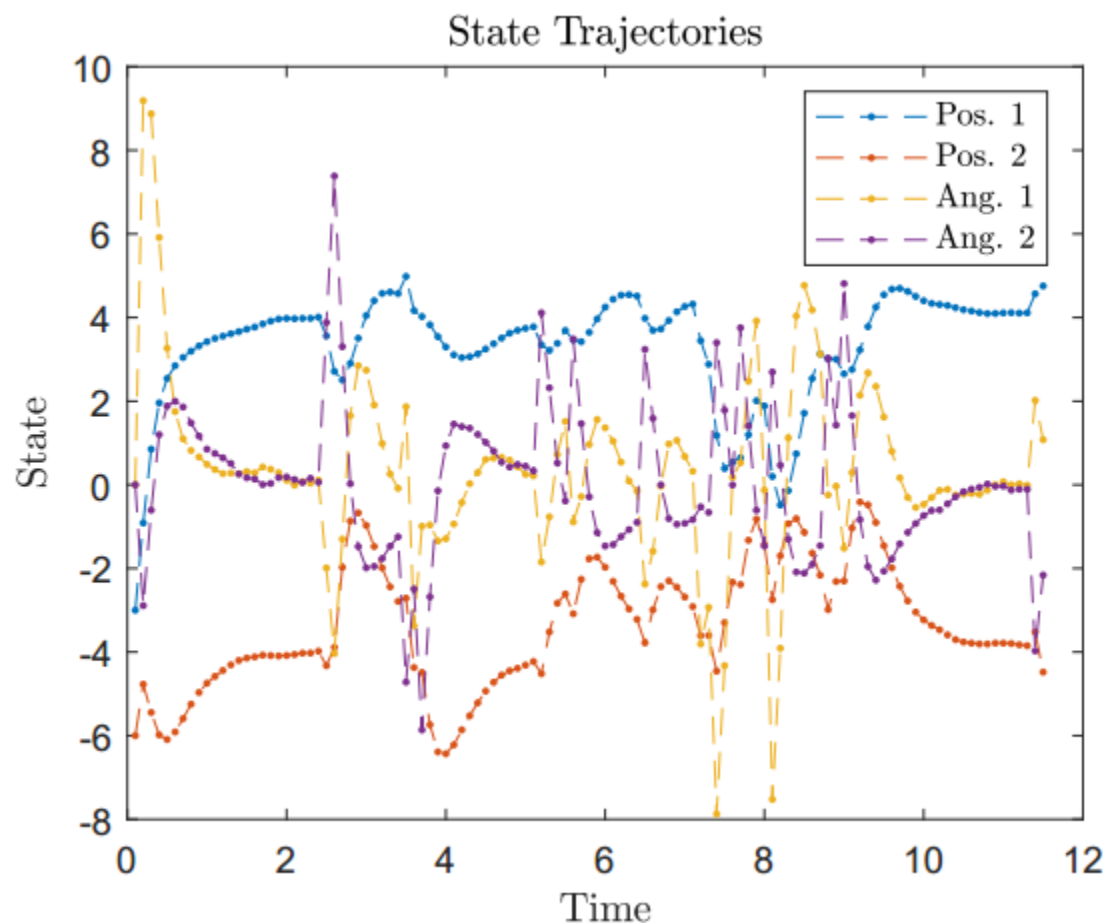
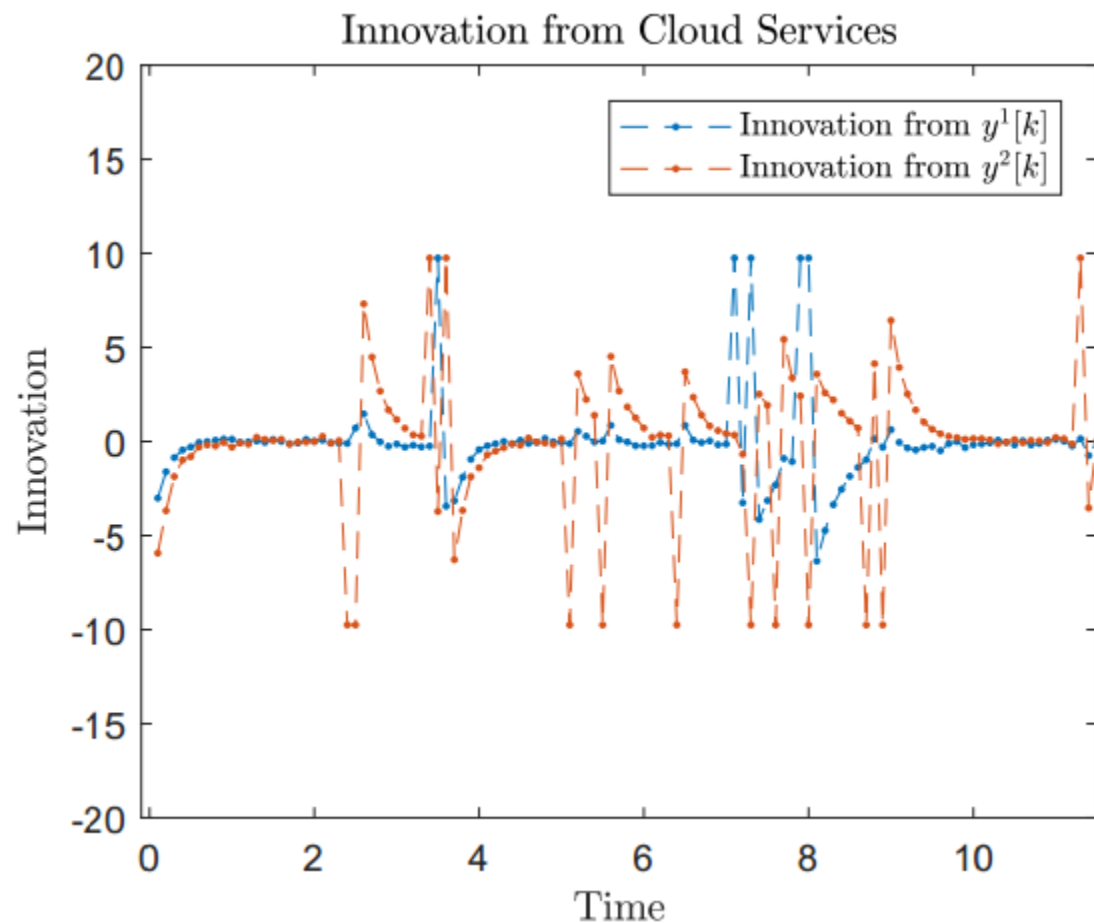
Simulation Results: High Risk, Gated



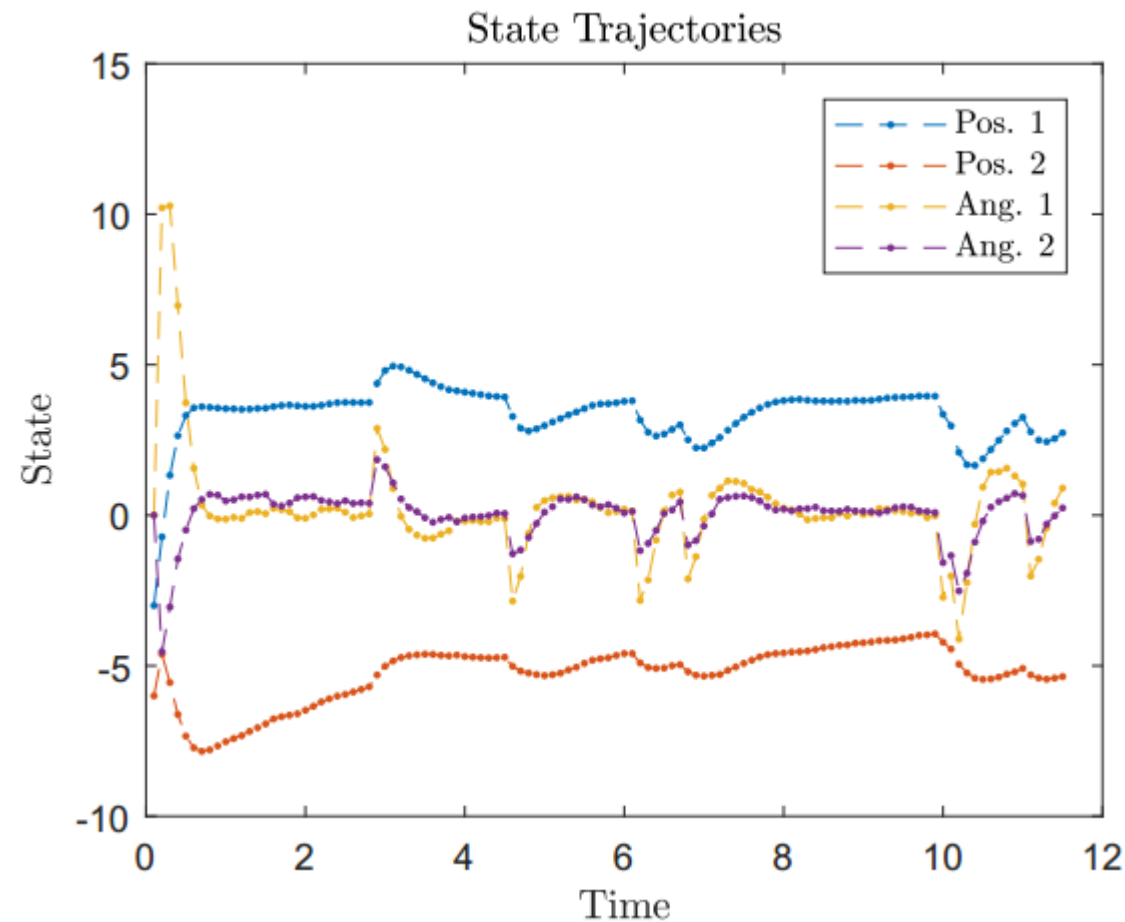
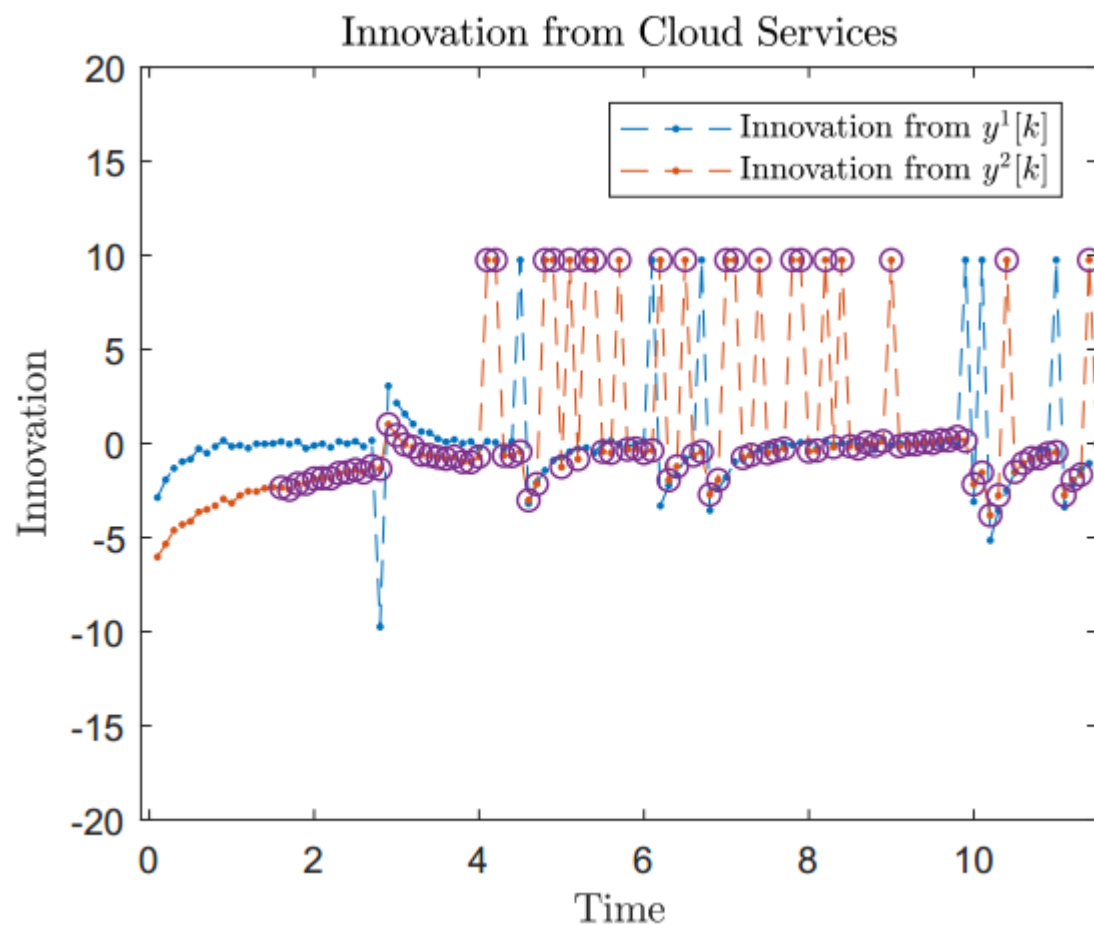
Simulation Results: Low Risk, Ungated



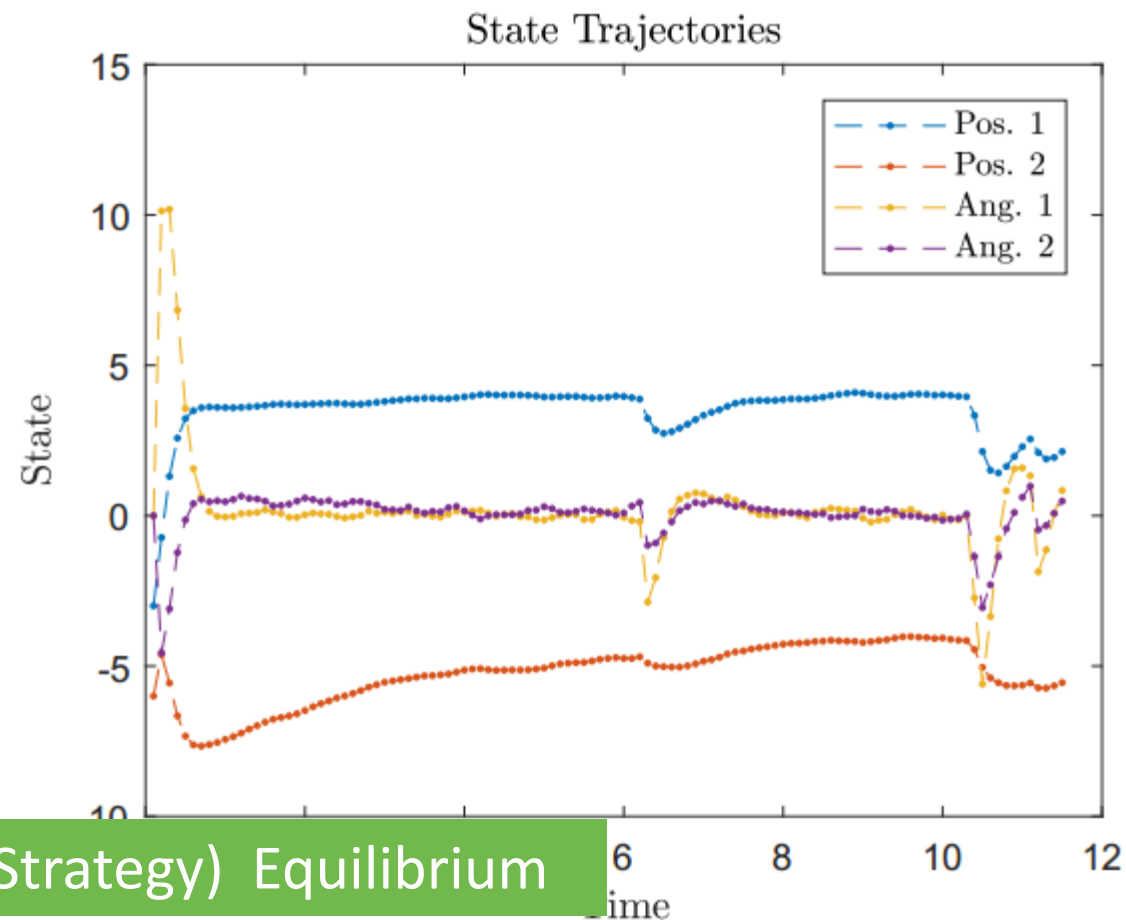
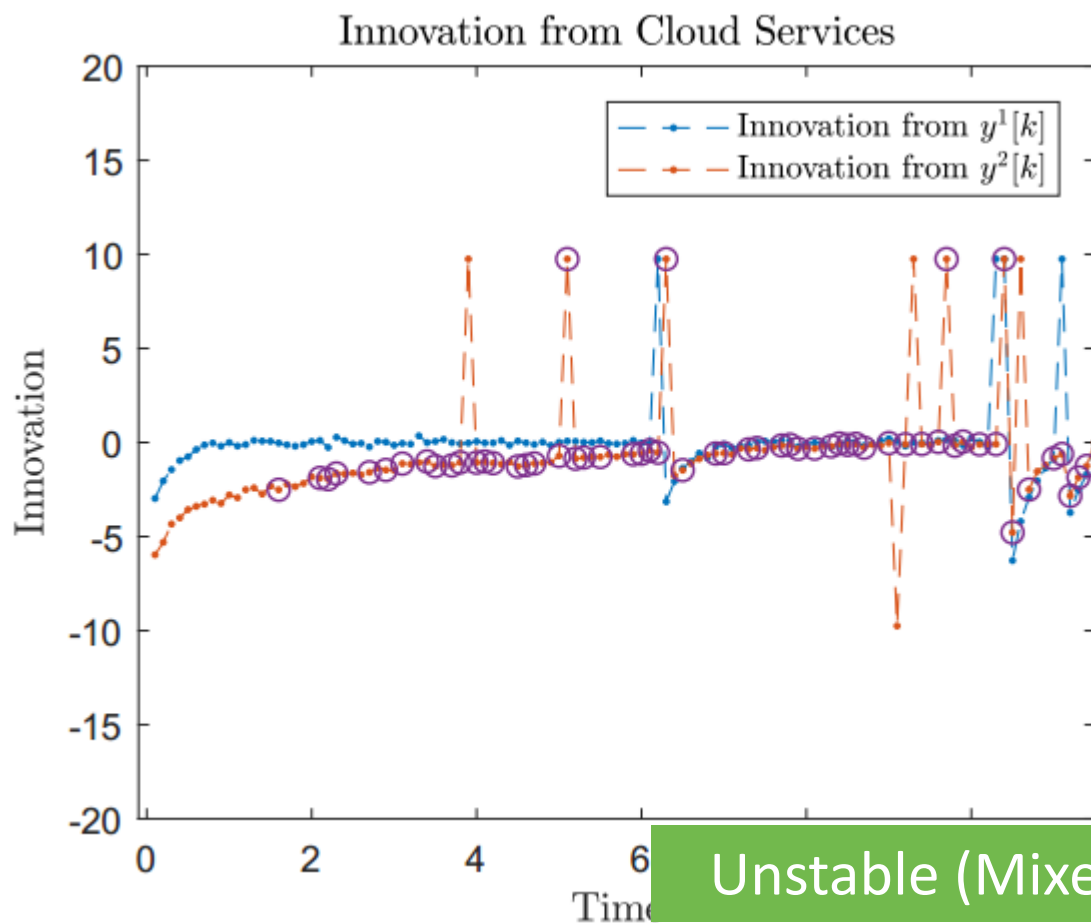
Decreased Attack Cost: Low Risk, Trusted



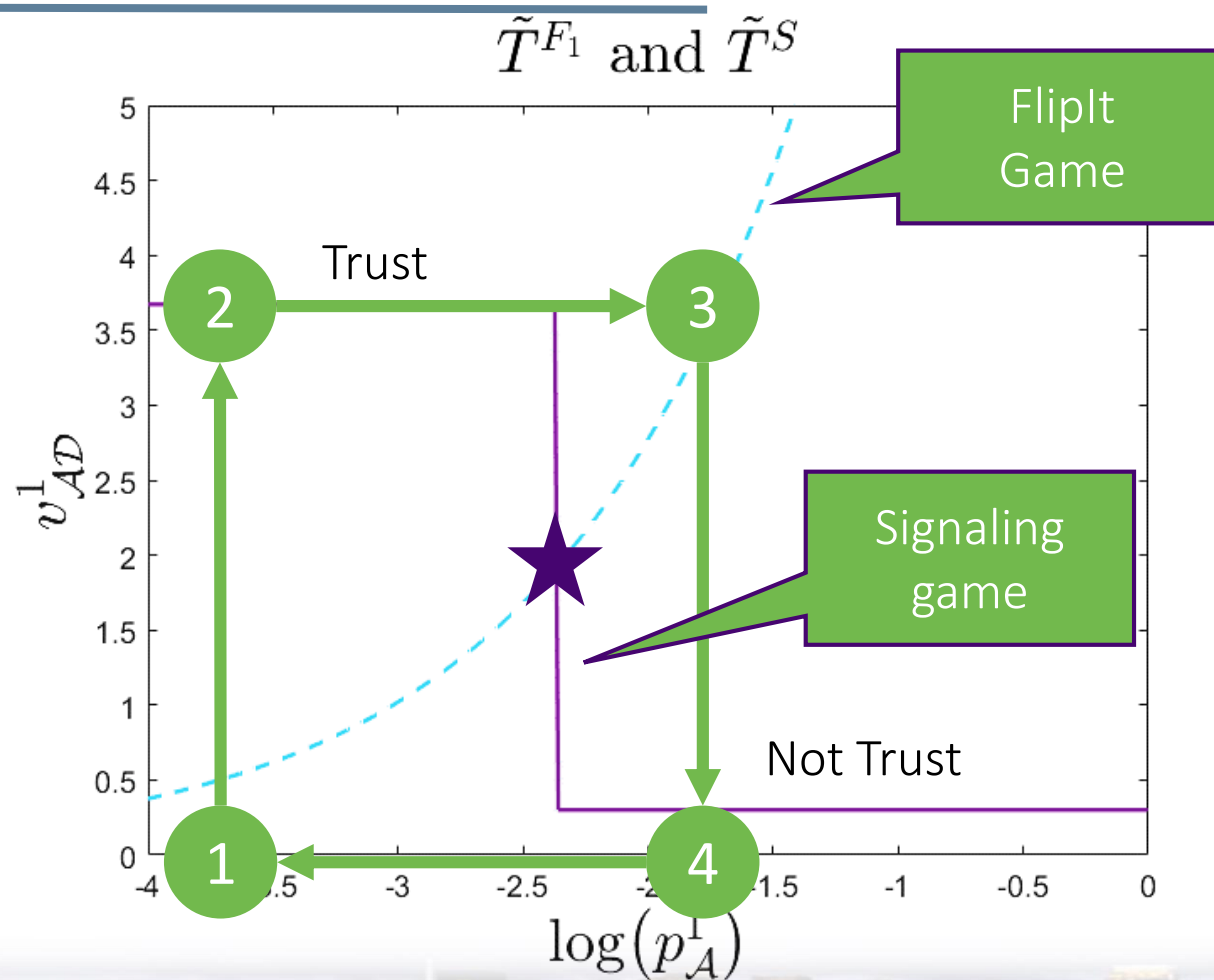
Decreased Attack Cost: Low Risk, Not Trusted



Decreased Attack Cost: Low Risk, Partially Trusted



The Telemarketer Cycle



Strategic Trust Summary

- Kalman filter handles sensor noise.
- Innovation gate rejects injected biases.
- Signaling game determines risk threshold beyond which even measurements within the innovation gate should be rejected.
- Prior probabilities of the signaling game are estimated proactively using Flipt games.
- Overall equilibrium concept: fixed point of the composition of mappings that describe all $N+1$ games.



Challenges for Future Work

- Taxonomy of counter-deception: can we include detection, trust, adversarial machine learning, and periodic renewal?
- Non-strategic trust: under what conditions can agents refrain from calculating strategies and simply *trust* other agents?
- General theory of multi-game compositions: can we formulate rules for combining games in series, in parallel, and in combinations of the two?



