



# An Investigation into Trust and Reputation Frameworks for Autonomous Underwater Vehicles

Thesis submitted in accordance with the requirements of  
the University of Liverpool for the degree of Doctor in Philosophy by

**Andrew Bolster**

May 2016



# Contents

|   |          |
|---|----------|
| Preface   | ix       |
| Abstract  | xi       |
| Acknowledgements  | xxiii    |
| <b>1 Introduction</b>   | <b>1</b> |
| 1.1 Mobile Ad-hoc Networks (MANETs) . . . . .                         | 1        |
| 1.2 Node Density in MANETs . . . . .                                  | 2        |
| 1.3 Routing in Mobile Ad-hoc Network . . . . .                        | 3        |
| 1.3.1 Proactive Routing . . . . .                                     | 3        |
| 1.3.2 Reactive Routing . . . . .                                      | 4        |
| 1.3.3 Hybrid Routing . . . . .  | 5        |
| 1.4 MANETs in Harsh Environments . . . . .                            | 5        |
| 1.5 Systems Approach to Trust and Trust Engineering . . . . .         | 5        |
| 1.6 Trust Operation Against Capable Attackers . . . . .               | 7        |
| 1.7 Contributions . . . . .   | 7        |
| 1.8 Conclusion . . . . .  | 7        |
| 1.8.1 Layout . . . . .  | 7        |
| <b>2 Background on Trust and its Applications to MANETs</b>           | <b>9</b> |
| 2.1 Trust Definitions and Perspectives . . . . .                      | 9        |
| 2.1.1 Modelling Trust Relationships . . . . .                         | 10       |
| 2.1.2 Taxonomy and Notations of Trust . . . . .                       | 12       |
| 2.1.3 Characteristics of Trust Relationships . . . . .                | 12       |
| 2.1.4 Topologies of Multi-Party Trust Networks . . . . .              | 13       |
| 2.1.5 Trust Establishment Strategies . . . . .                        | 14       |
| 2.1.6 Attacks on Trust . . . . .                                      | 14       |
| 2.2 Trusted Development and Operation of Autonomous Systems . . . . . | 15       |
| 2.2.1 Introduction . . . . .  | 15       |
| 2.2.2 Autonomy and Levels of Autonomy . . . . .                       | 15       |
| 2.2.3 Trust Perspectives in Autonomous Operation . . . . .            | 16       |
| 2.2.4 Design Trust . . . . .  | 18       |
| 2.2.5 Operational Trust . . . . .                                     | 24       |
| 2.2.6 Conclusions . . . . .   | 24       |

|          |  |           |
|----------|--|-----------|
| 2.3      | Trust in Autonomous MANETs . . . . .   | 24        |
| 2.3.1    | Trust Model Design Considerations . . . . .  | 24        |
| 2.3.2    | Attacks on MANETs . . . . .  | 25        |
| 2.3.3    | Trust Management Frameworks . . . . .  | 25        |
| 2.3.4    | Single Metric Trust Frameworks . . . . .   | 27        |
| 2.3.5    | Multi-Metric Trust Frameworks . . . . .  | 28        |
| 2.4      | Conclusion . . . . .   | 30        |
| <b>3</b> | <b>Maritime Communications and Operations</b>  | <b>31</b> |
| 3.1      | Introduction . . . . .   | 31        |
| 3.2      | Maritime Communications Environment . . . . .  | 31        |
| 3.2.1    | Mechanics of Acoustic Transmission . . . . .   | 31        |
| 3.2.2    | Velocity and density . . . . .   | 32        |
| 3.2.3    | Intensity and Power . . . . .  | 33        |
| 3.2.4    | Attenuation . . . . .  | 33        |
| 3.2.5    | Ambient Noise Model . . . . .  | 35        |
| 3.2.6    | Multipath effects . . . . .  | 35        |
| 3.2.7    | Modelling and Simulation of the Acoustic Medium / Channel . . . . .                      | 35        |
| 3.2.8    | Routing and Network Design for Underwater Acoustic Networks (UANs) . . . . .             | 36        |
| 3.3      | Marine Operations . . . . .  | 36        |
| 3.3.1    | Typical Autonomous Underwater Vehicle (AUV) mission profiles . . . . .                   | 36        |
| 3.3.2    | Potential Future Applications . . . . .  | 36        |
| 3.3.3    | Need for Trust in Maritime Networks . . . . .  | 37        |
| <b>4</b> | <b>Assessment of Trust Management Framework (TMF) Performance in Marine Environments</b> | <b>39</b> |
| 4.1      | Introduction . . . . .   | 39        |
| 4.2      | Modelling of UAN network . . . . .   | 40        |
| 4.2.1    | Mobility, Topology, and Communications . . . . .   | 40        |
| 4.2.2    | Simulation Background . . . . .  | 41        |
| 4.3      | Establishing Scale Factors in Communications Rate . . . . .                              | 41        |
| 4.3.1    | Scale Factors in Physical Node Distribution . . . . .                                    | 46        |
| 4.4      | Combined Scale Factor Analysis . . . . .   | 50        |
| 4.5      | Conclusions . . . . .  | 50        |
| <b>5</b> | <b>Use of Physical Behaviours for Trust Assessment</b>                                   | <b>53</b> |
| 5.1      | Introduction . . . . .   | 53        |
| 5.2      | AUV Mobility and Localisation . . . . .  | 54        |
| 5.2.1    | AUV operations and deployments . . . . .   | 54        |
| 5.2.2    | Localisation Technologies . . . . .  | 56        |
| 5.3      | Trust Management Frameworks . . . . .  | 57        |
| 5.4      | Physical Behaviours for Trust . . . . .  | 57        |
| 5.4.1    | Physical Metrics . . . . .   | 57        |
| 5.4.2    | Physical Misbehaviours . . . . .   | 58        |
| 5.5      | Simulation and Validation . . . . .  | 59        |
| 5.5.1    | Simulation Background . . . . .  | 59        |
| 5.5.2    | Node Control Modelling . . . . .   | 59        |

|          |  |            |
|----------|--|------------|
| 5.5.3    | Standards of Accuracy . . . . .  | 60         |
| 5.5.4    | Analysis . . . . .   | 61         |
| 5.6      | Results and Discussion . . . . .   | 63         |
| 5.6.1    | Detection of Misbehaviours . . . . .   | 65         |
| 5.6.2    | Identification of Misbehaviours . . . . .  | 66         |
| 5.6.3    | Impacts of Misbehaviour on operational performance . . . . .   | 67         |
| 5.7      | Conclusion . . . . .   | 68         |
| <b>6</b> | <b>Communications Trust Assessment in Underwater MANETs</b>  | <b>69</b>  |
| 6.1      | Introduction . . . . .   | 69         |
| 6.1.1    | Selected Misbehaviours . . . . .   | 69         |
| 6.2      | Simulation Results and Discussion . . . . .  | 69         |
| 6.2.1    | Comparison between Multi-parameter Trust Framework for MANETs (MTFM), Hermes and Objective Trust Management Framework (OTMF) . . . . . | 70         |
| 6.2.2    | Metric Weighting . . . . .   | 73         |
| 6.2.3    | Weight Significance Analysis for Behaviour Classification . . . . .  | 74         |
| 6.3      | Conclusions . . . . .  | 76         |
| 6.3.1    | Future Work . . . . .  | 77         |
| <b>7</b> | <b>Multi-Domain Trust Assessment in Collaborative Marine MANETs</b>  | <b>79</b>  |
| 7.1      | Introduction . . . . .   | 79         |
| 7.2      | Initial Optimisation of Multi-Domain Trust with Predefined Domains . . . . .   | 79         |
| 7.2.1    | Communications Trust Metrics . . . . .   | 80         |
| 7.2.2    | Physical Trust Metrics . . . . .   | 80         |
| 7.2.3    | Cross Domain Trust Metrics . . . . .   | 80         |
| 7.2.4    | Metric Weight Analysis Scheme . . . . .  | 81         |
| 7.2.5    | Significance Analysis . . . . .  | 82         |
| 7.2.6    | Weight Assessment . . . . .  | 82         |
| 7.3      | Conclusion . . . . .   | 95         |
| <b>A</b> | <b>Orphan Sections</b>   | <b>97</b>  |
| A.1      | Metric Weighting . . . . .   | 97         |
| A.2      | UNEDITED PROSE: Real Time Grey Systems . . . . .   | 97         |
| A.3      | From end of Defense Trust Conclusions . . . . .  | 99         |
| <b>B</b> | <b>Human Factors related to Trusted Operation of Autonomous Systems</b>  | <b>101</b> |
| B.1      | Information Overload . . . . .   | 101        |
| B.2      | Adaptive Automation . . . . .  | 101        |
| B.3      | Distributed Decision Making . . . . .  | 102        |
| B.4      | Complexity . . . . .   | 103        |
| B.5      | Cognitive Biases and Failing Heuristics . . . . .  | 103        |
| <b>C</b> | <b>Grey System Theory and Grey Trust Assessmen</b>   | <b>105</b> |
| C.1      | Grey numbers, operators and terminology . . . . .  | 105        |
| C.2      | Whitenisation and the Grey Core . . . . .  | 106        |
| C.3      | Grey Sequence Buffers and Generators . . . . .   | 106        |

|                     |   |            |
|---------------------|---|------------|
| C.4                 | Grey Trust  | 106        |
| <b>D</b>            | <b>Additional Graphs</b>  | <b>109</b> |
| D.1                 | From Subsection 7.2.6: Mean-Weighted Multi-Domain Trust Results                 | 109        |
| D.2                 | From Subsection 7.2.6: Multi-Domain Trust Results Targeting Non-malicious nodes | 116        |
| D.2.1               | Per Node Breakdowns   | 116        |
| D.2.2               | Averaged across remaining cohort  | 122        |
| <b>Todo</b>         | <b>list</b>   | <b>122</b> |
| <b>Bibliography</b> |   | <b>133</b> |

# Illustrations

## List of Figures

|     |  |    |
|-----|--|----|
| 2.1 | Model of Trust . . . . .   | 11 |
| 2.2 | Trust Topologies; Direct, Indirect, Recommender, etc. from the perspective of Node A . . . . .   | 14 |
| a   | Sample topology showing logical connections between nodes (Range of $A$ shown in red dashed line) . . . . .  | 14 |
| b   | Direct Relationships, the two possible trust assessments from $A$ to its connected neighbours, $B, C$ . . . . .  | 14 |
| c   | Indirect Relationships, showing the four possible trust assessments from $A$ or the three disconnected leaf nodes, $D, E, F$ . . . . .                                     | 14 |
| d   | Recommender Relationship, showing the two discrete paths trust assessments travel to $A$ ; $T_{A,B}^R = T_{A,C} \otimes T_{C,B}$ and $T_{A,C}^R = T_{A,B} \otimes T_{B,C}$ | 14 |
| 2.3 | American Society of Testing and Materials (ASTM) F41 Unmanned Maritime Vehicle System (UMVS) Architecture (with relevant substandards in parenthesis) . . . . .            | 23 |
| 2.4 | The inclusion of additional metrics and domains in trust assessment reduces the systems exposed threat surface . . . . .   | 28 |
| 3.1 | Thorp's formula . . . . .  | 33 |
| 3.2 | Ainslie & McColm Absorption Model . . . . .  | 34 |
| 3.3 | Fisher-Simmons Absorption Model . . . . .  | 34 |
| 3.4 | Non-Linear Marine Propagation in an isothermal profile . . . . .   | 36 |
| 4.1 | Initial layout with nodes spaced an average of 100m apart . . . . .  | 41 |
| 4.2 | Throughput performance overview for all mobilities under varying emission rates <i>IS THIS ENOUGH?</i> . . . . .   | 43 |
| a   | Static . . . . .   | 43 |
| b   | All Mobile . . . . .   | 43 |
| c   | All-but-one Mobile . . . . .   | 43 |
| d   | Single Mobile . . . . .  | 43 |
| 4.3 | Network performance varying packet emission rates for the static case . . . . .  | 44 |
| a   | Packet delivery . . . . .  | 44 |
| b   | Probability of arrival . . . . .   | 44 |
| c   | End-to-end delay . . . . .   | 44 |
| d   | RTS Ratios . . . . .   | 44 |
| 4.4 | Network performance varying packet emission rates for the all mobile case . . . . .  | 44 |
| a   | Packet delivery . . . . .  | 44 |
| b   | Probability of arrival . . . . .   | 44 |
| c   | End-to-end delay . . . . .   | 44 |

|      |   |    |
|------|---|----|
| d    | RTS Ratios . . . . .  | 44 |
| 4.5  | Network performance varying packet emission rates for the all-but-one mobile case . . . . .                     | 45 |
| a    | Packet delivery . . . . .   | 45 |
| b    | Probability of arrival . . . . .  | 45 |
| c    | End-to-end delay . . . . .  | 45 |
| d    | RTS Ratios . . . . .  | 45 |
| 4.6  | Network performance varying packet emission rates for the single mobile case                                    | 45 |
| a    | Packet delivery . . . . .   | 45 |
| b    | Probability of arrival . . . . .  | 45 |
| c    | End-to-end delay . . . . .  | 45 |
| d    | RTS Ratios . . . . .  | 45 |
| 4.7  | Throughput performance overview for all mobilities under varying separation<br><i>IS THIS ENOUGH?</i> . . . . . | 47 |
| a    | Static . . . . .  | 47 |
| b    | Single Mobile . . . . .   | 47 |
| c    | All-but-one Mobile . . . . .  | 47 |
| d    | All Mobile . . . . .  | 47 |
| 4.8  | Network performance varying node separation for the static case . . . . .                                       | 47 |
| a    | Packet delivery . . . . .   | 47 |
| b    | Probability of arrival . . . . .  | 47 |
| c    | End-to-end delay . . . . .  | 47 |
| d    | RTS Ratios . . . . .  | 47 |
| 4.9  | Network performance varying node separation for the all mobile case . . . . .                                   | 48 |
| a    | Packet delivery . . . . .   | 48 |
| b    | Probability of arrival . . . . .  | 48 |
| c    | End-to-end delay . . . . .  | 48 |
| d    | RTS Ratios . . . . .  | 48 |
| 4.10 | Network performance varying node separation for the all-but-one mobile case                                     | 48 |
| a    | Packet delivery . . . . .   | 48 |
| b    | Probability of arrival . . . . .  | 48 |
| c    | End-to-end delay . . . . .  | 48 |
| d    | RTS Ratios . . . . .  | 48 |
| 4.11 | Network performance varying node separation for the single mobile case . . . . .                                | 49 |
| a    | Packet delivery . . . . .   | 49 |
| b    | Probability of arrival . . . . .  | 49 |
| c    | End-to-end delay . . . . .  | 49 |
| d    | RTS Ratios . . . . .  | 49 |
| 4.12 | Normalised Throughput-Delay Product for all mobilities under varying separation and emission rate . . . . .     | 51 |
| a    | Static . . . . .  | 51 |

|     |  |    |
|-----|--|----|
| b   | Single Mobile . . . . .  | 51 |
| c   | All-but-one Mobile . . . . .   | 51 |
| d   | All Mobile . . . . .   | 51 |
| 5.1 | Visual representation of the basic Boidean collision avoidance rules used . . . . .  | 61 |
| a   | Cohesion . . . . .   | 61 |
| b   | Repulsion . . . . .  | 61 |
| c   | Alignment . . . . .  | 61 |
| 5.2 | Observed Metric Values for one simulation of each behaviour ( $x_{i,j}^{m,t}$ from (5.10))   | 63 |
| 5.3 | <i>Unnecessary but included for draft discussion</i> Observed Metric Values for one simulation of each behaviour ( $d_{i,j}^{m,t}$ from Fig. ??) . . . . .                       | 64 |
| 5.4 | Normalised Deviance values from one simulation of each behaviour ( $\alpha_{i,j}^{m,t}$ from (5.11)) . . . . .   | 64 |
| 5.5 | Per-Node-Per-Run deviance for each metric, normalised in time ( $\sum \alpha/T$ ) . . . . .  | 65 |
| 6.1 | MTFM Trust assessments of $n_1$ ( $T_{1,X}$ ), showing Direct, Recommender and Indirect relationships, as well as the Aggregate trust assessments from combining these . . . . . | 71 |
| a   | Fair Static . . . . .  | 71 |
| b   | Fair Mobile . . . . .  | 71 |
| c   | Malicious (MPC) Static . . . . .   | 71 |
| d   | Malicious (MPC) Mobile . . . . .   | 71 |
| e   | Selfish (STS) Static . . . . .   | 71 |
| f   | Selfish (STS) Mobile . . . . .   | 71 |
| 6.2 | $T_{1,0}$ for Hermes, OTMF and MTFM assessment values for fair and malicious behaviours in the fully mobile scenario (mean of MTFM also shown) . . . . .                         | 72 |
| a   | Fair Scenario . . . . .  | 72 |
| b   | Malicious Power Control (MPC) Scenario . . . . .   | 72 |
| c   | Selfish Target Selection (STS) Scenario . . . . .  | 72 |
| 6.3 | $T_{1,MTFM}$ in the All Mobile case for the Malicious Power Control behaviour, including dashed $\pm\sigma$ envelope about the fair scenario . . . . .                           | 73 |
| a   | Delay Emphasised . . . . .   | 73 |
| b   | PLR Emphasised . . . . .   | 73 |
| c   | RX Power Emphasised . . . . .  | 73 |
| d   | TX Power Emphasised . . . . .  | 73 |
| e   | Throughput Emphasised . . . . .  | 73 |
| f   | Offered Load Emphasised . . . . .  | 73 |
| 6.4 | $T_{1,MTFM}$ in the All Mobile case for the Selfish Target Selection behaviour, including dashed $\pm\sigma$ envelope about the fair scenario . . . . .                          | 74 |
| a   | Delay Emphasised . . . . .   | 74 |
| b   | PLR Emphasised . . . . .   | 74 |
| c   | RX Power Emphasised . . . . .  | 74 |

|      |   |     |
|------|---|-----|
| d    | TX Power Emphasised . . . . .   | 74  |
| e    | Throughput Emphasised . . . . .   | 74  |
| f    | Offered Load Emphasised . . . . .   | 74  |
| 6.5  | Random Forest Factor Analysis of Malicious (Malicious Power Control (MPC),<br>Selfish (Selfish Target Selection (STS)) and Fair behaviours compared against<br>each-other . . . . . | 76  |
| 7.1  | Plot of Communications Metric Feature Extraction ( $X_{comms}$ ) . . . . .  | 83  |
| 7.2  | Plot of Physical Metric Feature Extraction ( $X_{phys}$ ) . . . . .   | 83  |
| 7.3  | Multi Domain Metric Features Extraction ( $X_{merge}$ ) . . . . .   | 84  |
| 7.4  | MPC Comms Metric Trust (showing cohort trust assessments) . . . . .   | 87  |
| 7.5  | MPC Physical Metric Trust (showing cohort trust assessments) . . . . .  | 87  |
| 7.6  | MPC Full Metric Trust (showing cohort trust assessments) . . . . .  | 88  |
| 7.7  | STS Comms Metric Trust (showing cohort trust assessments) . . . . .   | 88  |
| 7.8  | STS Physical Metric Trust (showing cohort trust assessments) . . . . .  | 89  |
| 7.9  | STS Full Metric Trust (showing cohort trust assessments) . . . . .  | 89  |
| 7.10 | Shadow Comms Metric Trust (showing cohort trust assessments) . . . . .  | 90  |
| 7.11 | Shadow Physical Metric Trust (showing cohort trust assessments) . . . . .   | 90  |
| 7.12 | Shadow Full Metric Trust (showing cohort trust assessments) . . . . .   | 91  |
| 7.13 | SlowCoach Comms Metric Trust (showing cohort trust assessments) . . . . .   | 91  |
| 7.14 | SlowCoach Physical Metric Trust (showing cohort trust assessments) . . . . .  | 92  |
| 7.15 | SlowCoach Full Metric Trust (showing cohort trust assessments) . . . . .  | 92  |
| 7.16 | Assumptions made about the relevant domains of impact / detectability of<br>misbehaviours, and domain relevance of metrics, may not be optimal . . . . .                            | 95  |
| A.1  | MTFM Trust assessments for varying mobility options in the selfish case . .   | 97  |
| a    | All Nodes Static . . . . .  | 97  |
| b    | $n_1$ Randomly Walking . . . . .  | 97  |
| c    | All Nodes but $n_1$ Randomly Walking . . . . .  | 97  |
| d    | All Nodes Randomly Walking . . . . .  | 97  |
| A.2  | Beta Trust time varying assessments for of $n_1$ varying mobility options . .   | 98  |
| a    | All Nodes Static . . . . .  | 98  |
| b    | $n_1$ Randomly Walking . . . . .  | 98  |
| c    | All Nodes but $n_1$ Randomly Walking . . . . .  | 98  |
| d    | All Nodes Randomly Walking . . . . .  | 98  |
| D.1  | MPC Comms Metric Trust (showing mean of non-misbehaving nodes) . . . . .  | 109 |
| D.2  | MPC Physical Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 110 |
| D.3  | MPC Full Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 110 |
| D.4  | STS Comms Metric Trust (showing mean of non-misbehaving nodes) . . . . .  | 111 |
| D.5  | STS Physical Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 111 |
| D.6  | STS Full Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 112 |
| D.7  | Shadow Comms Metric Trust (showing mean of non-misbehaving nodes) . .   | 112 |

|      |   |     |
|------|---|-----|
| D.8  | Shadow Physical Metric Trust (showing mean of non-misbehaving nodes) . . . . .  | 113 |
| D.9  | Shadow Full Metric Trust (showing mean of non-misbehaving nodes) . . . . .  | 113 |
| D.10 | SlowCoach Comms Metric Trust (showing mean of non-misbehaving nodes) . . . . .  | 114 |
| D.11 | SlowCoach Physical Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 114 |
| D.12 | SlowCoach Full Metric Trust (showing mean of non-misbehaving nodes) . . . . .   | 115 |
| D.13 | MPC Comms Metric Trust (targeting non-malicious node) . . . . .   | 116 |
| D.14 | MPC Physical Metric Trust (targeting non-malicious node) . . . . .  | 116 |
| D.15 | MPC Full Metric Trust (targeting non-malicious node) . . . . .  | 117 |
| D.16 | STS Comms Metric Trust (targeting non-malicious node) . . . . .   | 117 |
| D.17 | STS Physical Metric Trust (targeting non-malicious node) . . . . .  | 118 |
| D.18 | STS Full Metric Trust (targeting non-malicious node) . . . . .  | 118 |
| D.19 | Shadow Comms Metric Trust (targeting non-malicious node) . . . . .  | 119 |
| D.20 | Shadow Physical Metric Trust (targeting non-malicious node) . . . . .   | 119 |
| D.21 | Shadow Full Metric Trust (targeting non-malicious node) . . . . .   | 120 |
| D.22 | SlowCoach Comms Metric Trust (targeting non-malicious node) . . . . .   | 120 |
| D.23 | SlowCoach Physical Metric Trust (targeting non-malicious node) . . . . .  | 121 |
| D.24 | SlowCoach Full Metric Trust (targeting non-malicious node) . . . . .  | 121 |
| D.25 | MPC Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .          | 122 |
| D.26 | MPC Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .       | 123 |
| D.27 | MPC Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .           | 123 |
| D.28 | STS Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .          | 124 |
| D.29 | STS Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .       | 124 |
| D.30 | STS Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .           | 125 |
| D.31 | Shadow Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .       | 125 |
| D.32 | Shadow Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .    | 126 |
| D.33 | Shadow Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .        | 126 |
| D.34 | SlowCoach Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .    | 127 |
| D.35 | SlowCoach Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . . | 127 |
| D.36 | SlowCoach Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node) . . . . .     | 128 |

## List of Tables

|     |   |    |
|-----|---|----|
| 1.1 | Summary of Characteristics of MANETs . . . . .  | 2  |
| 1.2 | Selection of Proactive Routing Protocols . . . . .  | 3  |
| 1.3 | Selection of Reactive Routing Protocols . . . . .   | 4  |
| 1.4 | Selection of Hybrid Routing Protocols . . . . .   | 5  |
| 1.5 | Comparison of Routing Strategy Classes . . . . .  | 6  |
| 2.1 | Definitions of Trust . . . . .  | 10 |
| 2.2 | Factors of Trust . . . . .  | 10 |
| 2.3 | Factors of Trust for Autonomous Systems . . . . .   | 11 |
| 2.4 | Definitions of Autonomy . . . . .   | 17 |
| 2.5 | Levels of Decision Making Automation . . . . .  | 18 |
| 2.6 | Levels of Automation . . . . .  | 19 |
| 2.7 | Level of Interoperability (LOI) for STANAG 4586 Compliant UCS . . . . .   | 21 |
| 3.1 | Contributing factors to Ocean Ambient Acoustic Noise . . . . .  | 35 |
| 4.1 | Comparison of system model constraints as applied between Terrestrial and Marine communications . . . . .               | 42 |
| 4.2 | Tabular view of data from Fig. 4.11, including ideal propagation time . . . . .   | 49 |
| 5.1 | REMUS 100 Mobility Constraints as applied in simulation . . . . .   | 60 |
| 5.2 | Overall Q-Test Outlier Correct Detection Accuracy . . . . .   | 65 |
| 5.3 | Per-Metric Q-Test Outlier Detection Accuracy . . . . .  | 66 |
| 5.4 | Metric Confidence Responses for known behaviours (5.12) . . . . .   | 66 |
| 5.5 | Successful Identification rates on untrained results using (5.13) . . . . .   | 67 |
| 5.6 | Successful Identification rates on untrained results using (5.13), with outlier consensus checks . . . . .              | 68 |
| 6.1 | Correlation Coefficients between metric weights and behaviour detection targets . . . . .                               | 75 |
| 7.1 | Multi Domain Metric Feature Correlation ( $X_{merge}$ ) . . . . .   | 84 |
| 7.2 | $\Delta T$ across domains and “proposed” behaviours targeting known misbehaving node . . . . .                          | 86 |
| 7.3 | $\Delta T^-$ False Positive assessments across domains and “proposed” behaviours across non-misbehaving nodes . . . . . | 93 |
| 7.4 | Optimised metric vector weights per domain trained upon and behaviour targeted . . . . .                                | 94 |





# Glossary

**ABR** Associativity Based Routing. [xiii, 4, 6](#)

**ACS** Autonomous Collaborative System. [xiii, 101](#)

**AODV** Ad hoc On-demand Distance Vector. [xiii, 4, 6](#)

**ASTM** American Society of Testing and Materials. [vii, xiii, 22, 23](#)

**AUV** Autonomous Underwater Vehicle. [iv, xiii, 36, 37, 40](#)

**BLOS** Beyond Line of Sight. [xiii, 21](#)

**CBRP** Cluster Based Routing Protocol. [xiii, 4, 6](#)

**CGSR** Cluster-head Gateway Switch Routing. [xiii](#)

**CMRE** Centre for Maritime Research and Experimentation. [xiii, 41](#)

**DDR** Distributed Dynamic Routing. [xiii, 5](#)

**DLI** Data Link Interface. [xiii, 21](#)

**DREAM** Distance Routing Effect Algorithm for Mobility. [xiii, 3, 6](#)

**DSDV** Destination-Sequences Distance Vector. [xiii, 3](#)

**DSR** Dynamic Source Routing. [xiii, 4, 26](#)

**DST** Distributed Spanning Trees. [xiii, 5](#)

**DSTL** Defence Science and Technology Laboratory. [xiii, 41](#)

**EOD** Explosive Ordnance Disposal. [xiii, 22](#)

**FSR** Fisheye State Routing. [xiii](#)

**GPS** Global Positioning System. [xiii, 4, 6](#)

**GRG** Grey Relational Coefficient. [xiii, 29, 70, 77](#)

- GSR** Global State Routing. [xiii](#)
- HRI** Human Robot Interaction. [xiii](#), [17](#)
- HSC** Human Supervisory Control. [xiii](#), [21](#), [101](#)
- HSR** Hierarchical State Routing. [xiii](#)
- INDD** Inter-Node Distance Deviation. [xiii](#), [80](#)
- INHD** Inter-Node Heading Deviation. [xiii](#), [80](#)
- ISR** Intelligence, Surveillance and Reconnaissance. [xiii](#), [21](#)
- JAUS** Joint Architecture for Unmanned Systems. [xiii](#), [22](#)
- LAR** Location Aided Routing. [xiii](#), [4](#), [6](#)
- LOA** Level of Automation. [xiii](#), [16](#), [18](#), [19](#), [101](#), [102](#)
- LOI** Level of Interoperability. [xii](#), [xiii](#), [21](#)
- LOS** Line of Sight. [xiii](#), [22](#)
- MAC** Medium Access Control. [xiii](#), [42](#)
- MANET** Mobile Ad-hoc Network. [iii–v](#), [xii](#), [xiii](#), [1–3](#), [5](#), [7–29](#), [37](#), [39](#), [40](#), [53](#), [68–77](#), [79](#), [82](#), [108](#)
- MCM** Mine-Counter Measure. [xiii](#)
- MHC** Maritime Hydrography Capability. [xiii](#)
- MMWN** Multimedia support in Mobile Wireless Networks. [xiii](#)
- MPC** Malicious Power Control. [x](#), [xi](#), [xiii](#), [69](#), [75](#), [76](#), [85](#), [87](#), [88](#), [109](#), [110](#), [116](#), [117](#), [122](#), [123](#)
- MTFM** Multi-parameter Trust Framework for MANETs. [v](#), [ix](#), [xiii](#), [26](#), [29](#), [30](#), [40](#), [70–72](#), [76](#), [107](#)
- OLSR** Optimised Link State Routing. [xiii](#), [3](#), [6](#)
- OSPF** Open Shortest Path First. [xiii](#), [6](#)
- OTMF** Objective Trust Management Framework. [v](#), [ix](#), [xiii](#), [7](#), [26](#), [27](#), [30](#), [40](#), [70–72](#), [75](#), [76](#)
- P2P** Peer to Peer. [xiii](#), [26](#)

- PKI** Public Key Infrastructure. [xiii](#), [5](#), [25](#)
- PLR** Packet Loss Rate. [xiii](#), [26](#), [40](#), [70](#), [72–75](#), [80](#)
- ROAM** Routing On-demand Acyclic Multipath. [xiii](#), [4](#), [6](#)
- RTS** Request To Send. [xiii](#), [42](#), [46](#)
- SAE** Society of Automotive Engineers. [xiii](#), [21](#), [22](#)
- SLURP** Scalable Location Update Routing Protocol. [xiii](#), [5](#)
- SoS** System of Systems. [xiii](#), [15](#)
- STANAG** NATO Standardization Agreement. [xiii](#)
- STAR** Source Tree Adaptive Routing. [xiii](#), [6](#)
- STS** Selfish Target Selection. [x](#), [xi](#), [xiii](#), [69](#), [75](#), [76](#), [86](#), [88](#), [89](#), [111](#), [112](#), [117](#), [118](#), [124](#), [125](#)
- TBRPF** Topology Dissemination Based On Reverse-Path Forwarding. [xiii](#), [3](#), [6](#)
- TMF** Trust Management Framework. [iv](#), [xiii](#), [5](#), [7](#), [28](#), [37](#), [39–50](#), [71](#), [72](#), [76](#), [77](#)
- TPP** Trusted Third Party. [xiii](#), [5](#), [25](#)
- UAN** Underwater Acoustic Network. [iv](#), [xiii](#), [7](#), [36](#), [40](#)
- UAV** Unmanned Aerial Vehicle. [xiii](#), [16](#), [21](#)
- UMVS** Unmanned Maritime Vehicle System. [vii](#), [xiii](#), [22](#), [23](#)
- V2V** Vehicle to Vehicle. [xiii](#), [15](#)
- VSM** Vehicle Specific Module. [xiii](#), [21](#)
- WRP** Wireless Routing Protocol. [xiii](#), [3](#)
- ZHLS** Zone-based Hierarchical Link State. [xiii](#), [5](#)
- ZRP** Zone Routing Protocol. [xiii](#), [5](#)



# Preface

This thesis is primarily my own work. The sources of other materials are identified.



# **Abstract**

As Autonomous underwater vehicles (AUVs) become technically more competent, and fiscally more attainable, their use has been applied to a great many areas within defence, commercial and environmental areas of concern. Increasingly, these applications are tending towards utilising independent collective behaviour of teams or fleets of these platforms.



# Acknowledgements

There are many people who deserve the highest thanks for their support, patience, kindness and understanding. The greatest thanks have to be distributed among my family and friends, for putting up with my madness; both the madness of starting it and the madness of seeing it through. Maybe I'll get a job that you can actually explain! Next, I must thank Professor Marshall, without whom this work wouldn't have been attempted let alone completed. Finally, even though I swore I'd never do it, this work is dedicated to R, who knows why.

Alan-hu Akbar



# Chapter 1

## Introduction

### 1.1 MANETs

With the explosive growth in the use of mobile telephony and the increasing miniaturisation and efficiency gains of portable communications devices, the classical paradigm of a broadcast/receiver or server/client has given way to an increasing use of decentralised, ad-hoc networks that not only accommodate but take advantage of network mobility.

Whether these networks are decentralised cellular / RF / 802.11 WiFi networks for use in disaster relief areas [1] or biologically inspired wireless sensor networks for low-energy, low-maintenance environmental monitoring [2][3], **MANET** theory developed over the past 30 years has gone from it's first formal definition, emerging from DARPA's Packet Radio Network research [4], to being an integral part of modern practical communications.

Minimally, a **MANET** consists of a collection of mobile physical entities (nodes) that communicate cooperatively to collect, distribute, disseminate, and collate data and/or influence across an area. In most cases **MANET** nodes incorporate bi-directional transceivers to send and receive data, however this bi-directionality is not a requirement (for example in the area of Wireless Sensor Networks [5]). **MANETs** may utilise omnidirectional, static, or steerable communications antennae, and a selection of protocols such as WiFi, Bluetooth, GSM, UMTS, Optical or Acoustics, and may incorporate a range of mobilities across nodes, from static devices, terrestrial and marine surface platforms, and aerial and underwater platforms. A core characteristic of **MANETs** is the inclusion and integration of heterogeneous node collections, i.e where different nodes or groups of nodes in a network may have different capabilities in terms of propulsion, sensor apparatus, communications capability, etc.

**MANETs** may be totally independent with no external connections; include independent per-node communications backhauls (e.g. Cellular Modems in mobile phones as part of a Bluetooth Personal Area Network), or include static nodes that provide infrastructure based backhaul. However, this multiplicity of variations and options presents several challenges to users and operators; the physical topology of **MANETs** can vary wildly over short periods of time. A particular challenge to **MANET** operation is that

given any node may operate as a routing / gateway node, if/when that node moves to a different region, network segments that had previously used that node as a path must renegotiate / re-establish their routes. These situations, if not appropriately managed, lead to opportunities for subversion and selfishness.

The characteristics of **MANETs** as defined by Corson et al. are paraphrased in Table 1.1.

**Table 1.1:** Summary of Characteristics of **MANETs**[6]

|  |   |
|--|---|
| Dynamic Topologies                     | Nodes are free to move arbitrarily; thus, the typically multi-hop network topology may change randomly and rapidly at unpredictable times, and may consist of both bidirectional and unidirectional links.  |
| Bandwidth Constrained, Varied Capacity | Wireless links will continue to have significantly lower capacity than their hardwired counterparts. In addition, the realized throughput of wireless communications, after accounting for the effects of multiple access, fading, noise, and interference conditions, etc., is often much less than a radio's maximum transmission rate.<br>One effect of the relatively low to moderate link capacities is that congestion is typically the norm rather than the exception, i.e. aggregate application demand will likely approach or exceed network capacity frequently. |
| Energy Constrained Operation           | Some or all of the nodes in a <b>MANET</b> may rely on batteries or other exhaustible means for their energy. For these nodes, the most important system design criteria for optimization may be energy conservation.   |
| Limited physical security              | Mobile wireless networks are generally more prone to physical security threats than are fixed cable nets. The increased possibility of eavesdropping, spoofing, and denial-of-service attacks should be carefully considered.<br>Existing link security techniques are often applied within wireless networks to reduce security threats. As a benefit, the decentralized nature of network control in <b>MANETs</b> provides additional robustness against the single points of failure of more centralized approaches.  |

## 1.2 Node Density in **MANETs**

One fundamental compromise in the operation of wireless **MANETs** is the tradeoff between the number of hops required between source and destination nodes and the effective bandwidth available to the network overall[7]. This compromise is encapsulated in the relative density of a given network; that is, the number of nodes in a given node's one-hop locality, drawing direct links between wireless transmission strength / reception sensitivity, the environmental noise floor, environmental channel characteristics, the mobility of the nodes and the number of nodes deployed in a region.

Expand Node Density discussion to include examples of sparse, dense, long/deep, fully connected networks

**Table 1.2:** Selection of Proactive Routing Protocols

| Name  | Description   |
|-------|---|
| DSDV  | <b>Destination-Sequences Distance Vector</b> is a loop free derivative of the Distributed Bellman-Ford algorithm where each node maintains two tables; one that attempts to maintain a globally accurate next-hop routing table for all destination nodes (the routing table) and a route advertisement table, monitoring routes that the node itself can provide. These tables are updated both periodically and opportunistically. Loop-free status is maintained by monitoring a monotonic “sequence number”, which guarantees that if a long-loop returned packet is observed, it is discarded in favour of a route with a higher sequence number (i.e. newer route) [8]. |
| OLSR  | <b>Optimised Link State Routing</b> reduces the traffic-overhead of truly distributed link-state exchange and monitoring by establishing a multipoint replaying strategy (MRP) where nodes select a subset of their one-hop network relay to retransmit their packets, based on the two-hop connectivity of the network, thereby reducing contention and overheads by reducing local retransmitters. However, <b>OLSR</b> does not monitor link <i>quality</i> beyond binary “active/failed” state which can lead to non-optimal MRP and route selection in wireless networks for instance.   |
| WRP   |   |
| TBRPF |   |
| DREAM |   |

## 1.3 Routing in Mobile Ad-hoc Network

Given the decentralised nature of **MANET** operations, routing protocols are an active area of research. This research is classified according to the strategies used for discovering, monitoring and updating routes within the network, and are usually grouped into three classes; proactive (or Table Driven), reactive (or On Demand) and hybrid. A summary of the generalised characteristics of these classes is shown in [Table 1.5](#).

### 1.3.1 Proactive Routing

In Proactive routing, protocols attempt to maintain a up-to-date, global topology awareness of the network, where every node knows how the best next-hop to contact any other node in the network. This is extremely efficient for relatively small, static networks, with minimal storage and time requirements []. When the network topology is significantly modified by a shift in topology, either due to a node “dropping out” or moving, route renegotiation and optimisation is extremely resource consuming, as this global state is converged upon in a distributed manner by nodes exchanging their local knowledge of the “new” topology. The decomposition and updating of the node-knowledge of the network state, and the method of updating these state-tables, is the primary differentiator between proactive protocols, a selection of which are summarised in [Table 1.2](#).

**Table 1.3:** Selection of Reactive Routing Protocols

| Name | Description  |
|------|--|
| DSR  | [9]  |
| AODV |  |
| ROAM |  |
| ABR  |  |
| LAR  | <b>Location Aided Routing</b> incorporates location information (usually from <b>GPS</b> ), and generates a heuristic based on either the distance from the current node <i>towards</i> the destination location, or the distance from the current node <i>away from</i> the original source, minimising and maximising this distance respectively. These methods limit control overheads and usually accurately determine the shortest path. However, in highly mobile networks this behaviour appears increasingly flood-like (similar to <b>DSR</b> and <b>AODV</b> ), and the general requirement for highly accurate and timely positional information restricts the application of this protocol |
| CBRP | <b>Cluster Based Routing Protocol</b> uses a hierarchical topology where each cluster has a cluster-head which coordinates routing within that cluster. As only cluster-heads coordinate routing across clusters, transmission overheads are minimised compared to other route distribution methods. However, the negotiation and maintenance overheads and propagation delays associated with hierarchical clustering make the network susceptible to temporary routing loops as nodes may have inconsistent residual routing information during cluster re-negotiation   |

Finish Proactive Routing Protocols Table

### 1.3.2 Reactive Routing

In contrast to Proactive Routing, Reactive (or “on-demand”) routing establishes routing information when it is required, rather than in advance. This route establishment is usually based on a request-response exchange where the node requesting routing information “floods” its local network with next-hop requests. The format of this flooding (and the context of any responses) are the main differentiators between protocols, as shown in [Table 1.3](#). The on-demand nature of route discovery can lead to significantly lower traffic than proactive routing protocols, but this is often a trade-off between lower average traffic and larger pre-transmission discovery delays. As such, reactive routing lends itself to low-traffic, delay tolerant, dynamic mobile applications as it does not require rediscovery after every *topology* change, but only on transmission along a new or stale route.

Finish Reactive Routing Protocols Table

**Table 1.4:** Selection of Hybrid Routing Protocols

| Name  | Description |
|-------|-------------|
| DST   |             |
| DDR   |             |
| ZRP   |             |
| ZHLS  |             |
| SLURP |             |

### 1.3.3 Hybrid Routing

Write Hybrid Routing Protocols

Finish Hybrid Routing Protocols Table

## 1.4 MANETs in Harsh Environments

As [Mobile Ad-hoc Networks \(MANETs\)](#) grow beyond the terrestrial arena, their operation and the protocols designed around them must be reviewed to assess their suitability to different communications environments, ensuring their continued security, reliability, and performance.

The distributed and dynamic nature of [MANETs](#) mean that it is difficult to maintain a [Trusted Third Party \(TTP\)](#) or evidence based trust system such as Certificate Authorities or using [Public Key Infrastructure \(PKI\)](#).

more background on the operation of TTP/CA/PKI?

Therefore, a distributed, collaborative system must be applied to these networks. Such distributed trust management frameworks aim to detect, identify, and mitigate the impacts of malicious actors by distributing per-node assessments and opinions to collectively self-police behaviour. As such, [TMFs](#) can be used to predict and reason on the future interactions between entities in a system.

[TMFs](#) provide information to assist the estimation of future states and actions of nodes within [MANETs](#). This information is used to optimize the performance of a network against malicious, selfish, or defective misbehaviour by one or more nodes. Previous research has established the advantages of implementing [TMFs](#) in 802.11 based [MANETs](#), particularly in terms of preventing selfish operation in collaborative systems [11], and maintaining throughput in the presence of malicious actors [12]

## 1.5 Systems Approach to Trust and Trust Engineering

Trust as Assurance

**Table 1.5:** Comparison of Routing Strategy Classes (from Abolhasan et al. [10])

| \ \ Class<br>Area      | Proactive  | Reactive   | Hybrid  |
|------------------------|--|--|---|
| Routing Structure      | Both flat and hierarchical structures are available  | Mostly flat except <b>CBRP</b>   | Mostly hierarchical   |
| Route Availability     | Always available if nodes are reachable  | Determined when needed   | Depends on the location of the destination  |
| Control Traffic Volume | Usually high, attempt at reduction is made.<br>e.g. <b>OLSR</b> , <b>TBRPF</b>                   | Lower than Global routing and further improved using <b>GPS</b> .<br>e.g. <b>LAR</b>   | Mostly lower than proactive and reactive  |
| Periodic Updating      | Yes, some may be conditional e.g. <b>STAR</b>  | Not required, however some nodes may require periodic beacons. e.g. <b>ABRs</b>  | Usually used within each zone or between gateway nodes  |
| Mobility Handling      | Usually updates occur at fixed intervals. <b>DREAM</b> alters periodic updates based on mobility | <b>ABR</b> uses localised broadcast queries, <b>ROAM</b> uses threshold updates, <b>AODV</b> routing uses local route discovery              | Usually more than one path may be available. Single point of failures are reduced by working as a group     |
| Storage Requirements   | High   | Dependent on number of nodes kept or required; usually lower than proactive protocols  | Usually depends on cluster or zone size; may become as large as proactive if clusters are big               |
| Delay Level            | Short routes are predetermined   | Higher than proactive  | Short for destinations in the same zone/-cluster as source. Interzone may be as large as Reactive protocols |
| Scalability            | Up to 100 nodes; <b>OSPF</b> and <b>TBRPF</b> may scale higher                                   | Source routing protocols; up to a few hundred nodes. Point-to-point may scale higher. Depends on level of traffic and levels of multihopping | Designed for up to or more than 1000 nodes  |

## 1.6 Trust Operation Against Capable Attackers

Trust operation against capable attackers

## 1.7 Contributions

Contributions

## 1.8 Conclusion

Conclusions including Layout

### 1.8.1 Layout

#### Chapter 2

In this chapter the current literature and research on the concepts, theory, and applications concerning Trust and Trust Management is explored, specifically leaning towards the applications of Trust within Autonomous [MANETs](#).

In [Section 2.1](#), the abstract quantity of “trust” is explored, In [Section 2.2](#), Autonomy and “Trusted Operation” of autonomous systems is investigated from a system architects and a system operators perspective. In [Section 2.3](#), current use and applications of Trusted operation of [MANETs](#) is explored, including current [TMFs](#).

#### Chapter 4

In this chapter, the need for multi-metric trust assessment in [UAN](#) is demonstrated as an example of a harsh network environment.

The operation of a selection of traditional [MANET TMFs](#) in this environment is investigated. These challenges are characterised and results are presented that demonstrate a multi-metric approach to Trust greatly enhances the effectiveness of [TMFs](#) in these environments.

In [Section 4.2](#) an experimental configuration for the marine space is established, and the scenarios and results presented in [13] are reviewed for comparison. In [Section 6.2](#) findings in trust establishment and malicious behaviour detection are presented and comparing with other current [TMFs](#) (Hermes and [OTMF](#)) and the use of this multi-parameter approach to detecting malicious and selfish behaviour in autonomous marine networks is analysed.

The contributions of this chapter are a study on the comparative operation and performance of [TMFs](#) in marine acoustic networks, and a review of metric suitability for [TMFs](#) in marine environments, informing future metric selection for experimenters and theorists. Finally, a methodology to assess the usefulness of metrics in discriminating against misbehaviours in such constrained, delay-tolerant networks is demonstrated.

Key parts of this chapter were presented at TrustCom-BigDataSE-ISPA 2015 as “Single and Multi-Metric Trust Management Frameworks for use in Underwater Autonomous Networks.” [14]

## Chapter 2

# Background on Trust and its Applications to MANETs

### 2.1 Trust Definitions and Perspectives

For a term that is so common in every-day speech, “Trust”<sup>1</sup> is a challenging discussion area, particularly given the wealth of proposed definitions (Table 2.1).

Beyond these dry, vague, and often “fuzzy” definitions, there is a significant ontological conflict between the subjective and objective perspectives of trust; is “trust” an attribute of the actor performing a given action, or of the observer of such an action? Or indeed is trust itself an action upon a relationship between actors? Is it qualitative or quantitative? These questions have challenged philosophers, psychologists and social scientists for decades.

In human trust relationships it is recognized that there can be several domains of trust for example organizational, sociological, interpersonal, psychological and neurological [15].

These domains of trust are, from a human perspective, quite natural and are formed during the earliest stages of linguistic integration. This leads to recognisable deviations in the experiential concept of “trust” across cultures with differing linguistic histories. This has led to a wealth of work in the social sciences (as well as management schools across the world) in to how to develop, understand, and repair trust across cultural boundaries [16].

As such it is important to explore the following areas of trust definitions, the characteristics of trust relationships and the impact of topology on the information available to assess trust within an abstract network before approaching the application of Trust towards Autonomous Systems and finally to MANETs:

---

<sup>1</sup>As a point of notation, in this work ”Trust” and ”trust” are used interchangeably to refer to the concept, action, or belief of a specified trusting relationship. Where Trust is capitalised outside of grammatical convention, it is to emphasise “trust as a concept” rather than a particular value or relationship

**Table 2.1:** Definitions of Trust

| Definition  | Source            |
|---|-------------------|
| Assured reliance on the character, ability, strength, or truth of someone or something.   | Merriam-Webster   |
| Firm belief in the reliability, truth, or ability of someone or something   | OED               |
| The willingness of a party to be vulnerable to the actions of another party based on the expectation that the other will perform a particular action important to the trustor, irrespective of the ability to monitor or control that other party | Mayer et al. [17] |
| An expectancy held by an individual or a group that the word, promise, verbal or written statement of another individual or group can be relied upon  | Rotter [18]       |

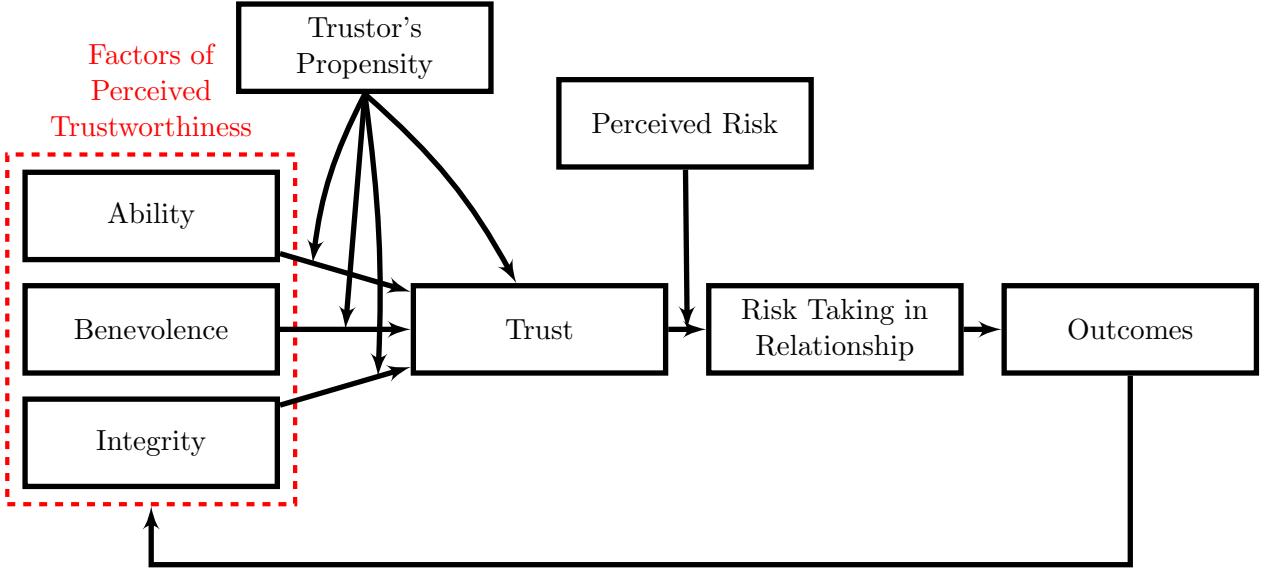
### 2.1.1 Modelling Trust Relationships

Mayer et al. [17] proposed a model of trust that encapsulates generalised factors of perceived trustworthiness of a *trustee* in interpersonal relationships (Table 2.2), accommodating a subjective trustworthiness and risk-taking potentiality on the part of the *trustor*. This formulation of trust allowed a wider discussion of the characteristics of trust relationships, both between individuals and within networks or communities.

**Table 2.2:** Factors of Trust (from Mayer et al. [17])

| Factor      | Definition   |
|-------------|--|
| Ability     | Collection of skills, competencies, capabilities and characteristics that enable a party to have influence or action within some specific domain |
| Benevolence | The extent to which a trustee is believed to want to do good to or by the trustor beyond a selfish profit motive                                 |
| Integrity   | Acceptance or adherence to a common set of principals of operation that the trustor finds acceptable   |

As shown in Fig. 2.1, Mayer primarily focuses on the Trustor's perspective and processes with respect to a given trust-based relationship. Three primary factors of perceived trustworthiness; based on previous outcomes, are assessed and synthesised along with the Trustor's own interanalysed propensity to Trust with respect to the different factors observed, to generate a given trust value. This trust value is incorporated with the risk / reward as assessed by the trustor to conclude what level of risk taking (Trust) can be assumed in the relationship between this trustor and a given trustee.

**Figure 2.1:** Model of Trust (from [17])**Table 2.3:** Factors of Trust for Autonomous Systems (from Lee and See [15])

| Factor      | Definition   | Mayer Term  |
|-------------|--|-------------|
| Performance | 'The current and historical operation of the automation, including characteristics such as reliability, predictability, and ability' | Ability     |
| Process     | The degree to which the automation's algorithms are appropriate for the situation and able to achieve the operators goals.           | Integrity   |
| Purpose     | The degree to which the automation is being used within the realm of the designers intent  | Benevolence |

Fig. 2.1: Only reasonable delimitation of Trustee Operation that doesn't corrupt  
Mayers thesis is curring half way though Risk Taking, Outcomes and *maybe* per-  
ceived risk (as this is also affected by outcomes; may make more sense to have a sep-  
arate augmented diagram)

Lee and See [15] extended and synthesised Mayer et al's approach to personal and interpersonal trust towards a generalised concept of trust for human and autonomic/autonomous systems with alternative contextual definitions shown in Table 2.3 (including their approximate mappings to Mayer et al's approach).

Sun et al. [19] suggests that there are two overarching forms of trust:

- Behavioural: That one entity voluntarily depends on another entity in a specific situation
- Intentional: That one entity would be willing to depend on another entity

It is suggested that these overarching forms are supported by and indeed are drawn from four major constructs within social and networked environments, as identified by McKnight and Chervany [20]:

- Trusting Belief: the subjective belief within a system that the other trusted components are willing and able to act in each others best interests
- Dispositional Trust: a general expectation of trustworthiness over time
- Situational Decision Trust: in-situ risk assessment where the benefits of trust outweigh the negative outcomes of trust
- System Trust: the assurance that formal impersonal or procedural structures are in place to ensure successful operation.

Sun argues that only System Trust and Behavioural Trust are relevant to trusted networking applications. However, it is arguable that in any communications network where the operation of that network is not the only concern, or where that network has to interact with any operator, then all of these factors come into play; as we will see ([Subsection 2.2.3](#)). Both System and Behavioural trust rely on what Sun calls a Belief Formation Process, or a trust assessment, while the other trust constructs deal with the interactions between trust and decision making against an internal assessment of network trustworthiness.

### 2.1.2 Taxonomy and Notations of Trust

Liu and Wang do lots on this [21] as well as discussion regarding the entropic/probabilistic models of trust. This may be too much to throw in, might inject it later

Talk about trust vs untrust vs nontrust

Explore notations of transitivity and abstract trust synthesis

### 2.1.3 Characteristics of Trust Relationships

There are five commonly considered characteristics or attributes of Trust relationships in general, but not all relationships exhibit them and they are not assumed to be a complete specification of Trust (synthesised from [17, 20, 22, 23]):

- *Multi-Party* - One-to-one; one-to-many; many-to-one; many-to-many. Trust is not an absolute characteristic of a lone individual. Trust may include multi-agent abstractions (one-to-many), such as a preferential trust/distrust towards a group exhibiting a particular attribute, e.g. members of the armed forces / police services. Likewise, there can be trustor/trustee attributes that can generalise relationships between collectives (many-to-many), e.g. Jets and Sharks[24].

- *Transitive* - Trust assessments can be shared (i.e. recommendations), where this second order trust assessment incorporates both the observed trustworthiness of the trustee, as well as the trustworthiness of the intermediate trustor. In some models this is further extended to include out-of-network intermediate trustors that have some other defined authority, e.g. PKI Certificate Authority
- *Evidential* - Trust must be based on some form of evidence-based observation or assessment, such as historical success rates of performing a certain action, or second-hand observations of trust from a third party.
- *Directional Asymmetry* - The majority of relationships are bi-directional but are asymmetric, i.e. between two entities who “trust” each other, there are two independent trust relationships that may have very different “values” or extents.
- *Contextual* - Trust can be variable and loosely coupled between contexts with respect to the action being assessed or the environment within which the trustee is operating, e.g. Doctors are trusted to perform medical procedures but that trust may not improve their success at correctly wiring an electrical plug. However there are plenty of counter-examples to this, as from [17], two of the three listed factors of trust are “Benevolence” and “Integrity” and these are unrelated to the ability of a trustee to perform a particular action, so it is reasonable to make an initial assumption that if a trustee is being benevolent in one activity or context, that that benevolence *should* extend to other contexts.

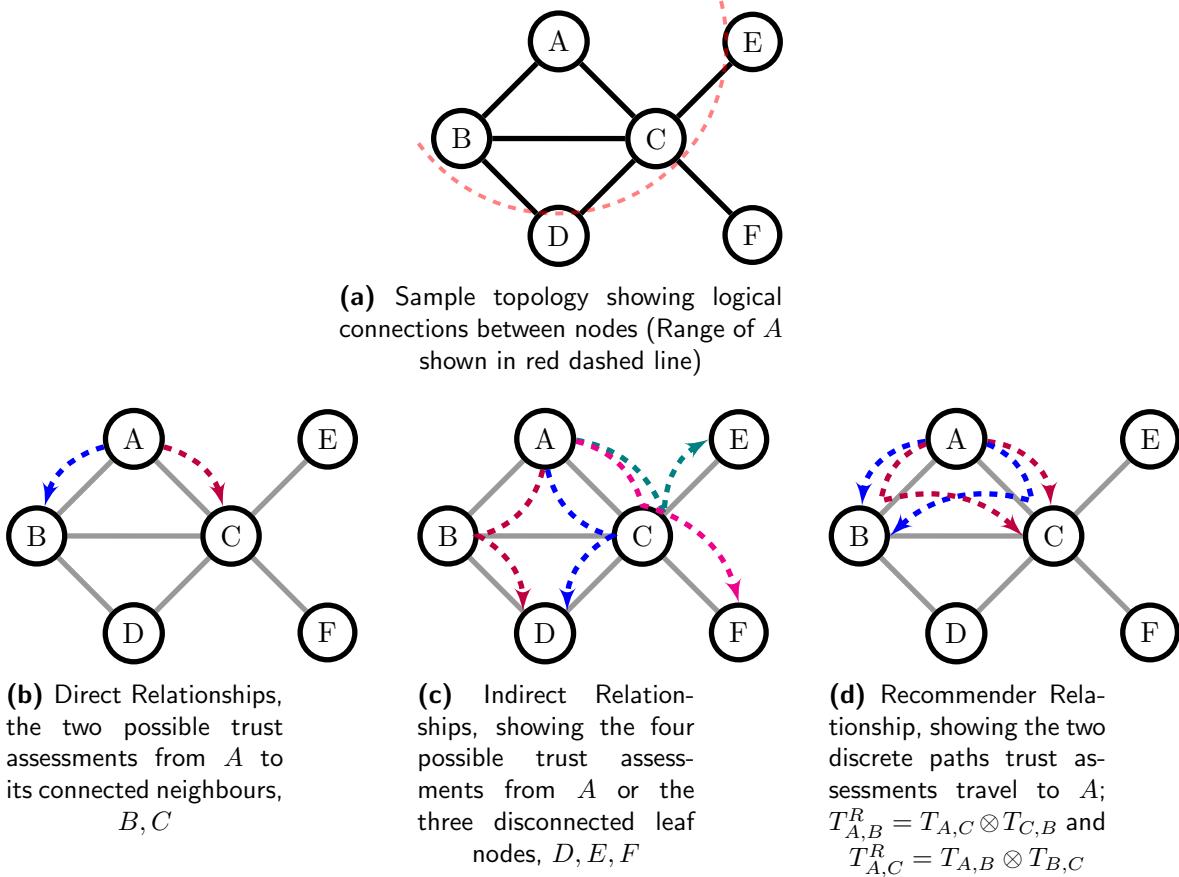
#### 2.1.4 Topologies of Multi-Party Trust Networks

Beyond the attributes or characteristics of an individual trust relationship, within any multi party sparsely connected network or community, topological context is useful in both establishing trust and in disseminating observations for collaborative assessment.

Within sparsely connected networks, there are three primary types of relationship, minimally demonstrated in Fig. 2.2;

- *Direct* - Whereby two nodes have a 1-hop communications link between them ( $A, B, C$  in the given figure)
- *Indirect* - Where two nodes have a  $n > 1$  hop communications link ( $E, D$  from  $A$  or  $C$ 's perspective in the given figure)
- *Recommendation* - Where three nodes are fully connected so as to enable the exchange of direct opinions and form composite opinions based on the target and reporter (i.e.  $A$  has both its own Direct assessment of  $C$ , as well as it's knowledge of  $B$ 's Direct assessment of  $C$ )

Redo relationship examples after Notation finished



**Figure 2.2:** Trust Topologies; Direct, Indirect, Recommender, etc. from the perspective of Node A

### 2.1.5 Trust Establishment Strategies

Need to discuss how trust is established a) initially among a co-launched group, b) with a newcomer and c) with a returner (Li and Singhal [11], Liu [22], Theodorakopoulos and Baras [25])

### 2.1.6 Attacks on Trust

In [21], Liu and Wang identify four types of attacks on Trust within networks that generate collaborative trust assessments through the exchange of recommendations; On-Off, Conflicting-behaviour, Badmouthing, and Sybil/Newcomer attacks.

The all of these attacks can be abstracted as “non-isotropic attacks” i.e. attacks that attempt to hide malicious / selfish behaviour behind the expected statistical variation in observations within a cohort. In each case, a different dimension of this assumed statistical normality is exploited; in On-Off, the attacker attempts to “hide” in the time dimension by only occasionally misbehaving, in the Badmouth attack the attacker is relying on it’s false recommendation being equitably received as its targets true actions. In the Conflicting behaviour attack, the attacker effectively “badmouths” a subset of nodes, hiding itself amid the “false” reports coming from the conflicting subsets of nodes.

Finally, in Sybil/Newcomer attacks the attacker takes advantage of an assumed naivety of the collective by presenting itself as a “new”, and therefore, zero-history entity that can initially neither be trusted or untrusted.

## 2.2 Trusted Development and Operation of Autonomous Systems

### 2.2.1 Introduction

The aim of the section is to explore where trust is likely to impact [System of Systems \(SoS\)](#) that contain autonomous elements incorporating Human Factors, Command and Control considerations, and [Vehicle to Vehicle \(V2V\)](#) distributed communication, from the perspective of trusted and semi-trusted operation.

Expand introduction and plan the rest of the section

### 2.2.2 Autonomy and Levels of Autonomy

Autonomy, like trust, is a nebulous term applied across research, defense and commercial circles that has its origins in human experience and interactions.

Autonomy, coming from the Greek roots *auto-* (self) and *nomos* (law) is the concept of a self-driven agency, and can be considered the concept of a “rational” individuals capacity to make un-coerced decisions in an informed manner. This autonomy is distinct from *freedom*, where freedom is the *ability* to perform an action, not the *capability to choose* which action to perform. That is not to say that autonomy or autonomous action exists in an ideal vacuum with perfect and complete information with no coercive factors or outside influences. The ability to recognise, process, weight and filter inputs, knowledge, “responsibilities”, influences and outside factors and come to an effective decision is a key skill for any self-governing agent, however this is above and beyond the concept of “basic autonomy”. From the implicit variability and complexity of environment and context that classically autonomous entitie<sup>2</sup> inhabit, there is little assumption that “autonomy” always produces a categorically “correct” or “good” decision, but is instead a case of an agent choosing the action that is *in its own best interests based on available information*[26].<sup>3</sup>

This understanding of individual autonomy has been scaled up through social systems and has been studied at length to understand the emergence of post-Marxist proto-anarchistic movements [27] and from a higher perspective, international politics, especially in the cases of quasi-federalised collections of states such as the United States of America [28] and the European Union/Eurozone/Schengen Area [29]

<sup>2</sup>That's *Homo Sapiens*

<sup>3</sup>Arply discusses a counter example of this “goodness” assessment as Huckleberry Finns’ release of Jim against his “best judgement”, and that rather than this action being an instance of morally justified or self-congratulatory autonomy, it was “the right thing to do” from an abstract moralistic perspective rather than a justifiably beneficial action, and it is a case of *akrasia*; the lacking of self-governance and the antonym of autonomy.

In the most general case in the world of artificial systems, Autonomy is understood as a graduated spectrum of allocation of functionality between a system (or system of systems) and a human operator assigned with performing a given task. Where a system is more “autonomous”, more of the sensing, planning, decision and action operations are performed by the system. (See Table 2.4 for a review of current definitions of autonomy and autonomous systems) This graduated spectrum of allocated functionality is generally termed the **Level of Automation (LOA)**, where an increasing LOA correlated to increasing control and decision making freedom to the autonomous system from the human operator.(Table 2.5)

While Autonomy is largely taken to be a robotics term based in the case of one human operator and one robotic entity, the development of more generalised cyber-physical systems has expanded this definition; from over-the-horizon human operation of **Unmanned Aerial Vehicles (UAVs)** to global networks of collaborating machines such as Google and beyond.

As such, the interactions *between* autonomous agents are becoming increasingly relevant to the operating efficiencies of overall collaborative systems, whether or not a human operator is “in-the-loop”.

Possibly expand this discussion

See Appendix B for a more thorough discussion on the Human Psychological Factors related to the planning, use, and integration of trusted autonomous systems in classical command and control contexts.

Ref Table 2.5 there may be a case to discuss the breakdown of “Plan, Decide, Execute, Inform”, possibly a nice onion-style graphic

### 2.2.3 Trust Perspectives in Autonomous Operation

For the purposes of this work, two perspectives on trust for autonomous systems are defined: Design Trust and Operational Trust.

- *Design Trust* - When an autonomous system is under development a level of Trust is established in it through the manner in which it has been designed and tested. This is the same as conventional systems. Given that systems that have high-levels of autonomy are designed to behave adaptively to dynamic environments, it is challenging to fully predict such non-deterministic behaviours prior to operational deployment. For example, in a navigation system it is difficult to predict the dynamic environment it will need to adapt to. Trust needs to be developed so that the design and testing of such systems are sufficient to predict that operation will be, if not optimal, at least satisfactory.
- *Operational Trust* - Trust at runtime or in-situ that both the individual nodes within a system are operating as expected and that the interfaces between the operator and the system are as expected. This latter aspect covers issues such as

**Table 2.4:** Definitions of Autonomy

| Definition  | Source                       |
|---|------------------------------|
| ...should be able to carry out its actions and to refine or modify the task and its own behaviour according to the current goal and execution context of its task   | Alami et al. [30]            |
| Autonomy refers to systems capable of operating in the real-world environment without any form of external control for extended periods of time   | Bekey [31]                   |
| ...a system situated within and a part of an environment that senses that environment and acts on it, over time, in pursuit of its own agenda and so as to effect with it senses in the future.<br>...Exercises control over its own actions  | Franklin and Graesser [32]   |
| An unmanned systems own ability of sensing, perceiving, analyzing, communicating, planning, decision-making, and acting, to achieve goals as assigned by its human operator(s) through designed HRI. ...The condition or quality of being self-governing  | Huang [33]                   |
| ...that the robot can operate self-contained, under all reasonable conditions without requiring recourse to the human operator. Autonomy means that a robot can adapt to change in its environment ...or itself ...and continue to reach a goal.  | Murphy [34]                  |
| ...it should learn what it can to compensate for partial or incorrect prior knowledge   | Russell and Norvig [35]      |
| Autonomy refers to a robot's ability to accommodate variations in its environment. Different robots exhibit different degrees of autonomy; the degree of autonomy is often measured by relating the degree at which the environment can be varied to the mean time between failures and other factors indicative of the robots performance. | Thrun [36]                   |
| ...agents operate without the direct intervention of humans or others, and have some kind of control over their actions and internal states.  | Wooldridge and Jennings [37] |

physical/wireless links and interpretation of data at each end of such a communication link. Can be subdivided into two types of perspective;

- *Hard Trust* or technical trust - The quantitative measurement and communication of the expectation of an actor performing a certain task, based on historic performance and through consensus building within a networked system. Can be thought of as a de-risking strategy to measure and monitor the ability of a system, or another actor within a system, to perform a task unsupervised.
- *Soft Trust* or common trust - The qualitative assessment of the ability of an actor to perform a task or operation consistently and reliably based on social or experiential factors. This is the human form of trust and is the main

**Table 2.5:** Levels of Decision Making Automation (Sheridan and Verplank [38])

| LOA | Description  |
|-----|--|
| 1   | The computer offers no assistance; the human must make all decisions and actions |
| 2   | The computer offers a complete set of decision/action alternatives, or           |
| 3   | Narrows the selection down to a few, or  |
| 4   | Suggests one alternative and   |
| 5   | Executes that suggestion if the human operator approves, or                      |
| 6   | Allows the human a restricted time to veto before automatic execution, or        |
| 7   | Executes automatically, then necessarily informs the human, and                  |
| 8   | Informs the human only if asked, or  |
| 9   | Informs the human only if it, the computer, decides to.                          |
| 10  | The computer decides everything and acts autonomously, ignoring the human.       |

motivational driver for the human-factors trust discussion in [Appendix B](#). Can be viewed as the abstract level of confidence an operator has in an actor to perform a task unsupervised.

Possibly worth looking at the Definition environment from amsthm to look after definitions like this

Operational Trust is functionally derived from, but distinct from Design Trust.

It is already clear that these two definitions are extremely close in their construction, but represent fundamentally different approaches to trust, one coming from a sociological perspective of person-to-person and person-to-group relationships from day to day life, and the other coming from a statistical or formal appraisal of an operation by a system.

Need to provide a linking section to the next blocks about Design/Operational Trust

#### 2.2.4 Design Trust

Five aspects of Design Trust have been identified:

No idea how to phrase this citation correctly; it's "my" work that was generated for DSTL and don't want to waste any more space backing it up; can I get away with just citing myself?

Rethink using these questions at all; opens up to awkward questioning that isn't answered in the thesis

1. **Formal Specification of Dynamic Operation:** Autonomous Systems (AS) may be required to operate in complex, uncertain environments and as such their specification may need to reflect an ability to deal with unspecified circumstances.

**Table 2.6:** Levels of Automation (paraphrased from Endsley and Kaber [39])

| LOA                       | Description   |
|---------------------------|---|
| Manual Control            | The human monitors, generates options, selects options (makes decisions), and physically carries out options.   |
| Action Support            | The automation assists the human with execution of selected action. The human does perform some control actions.  |
| Batch Processing          | The human generates and selects options; then they are turned over to automation to be carried out (e.g., cruise control in automobiles)  |
| Shared Control            | Both the human and the automation generate possible decision options. The human has control of selecting which options to implement; however, carrying out the options is a shared task.              |
| Decision Support          | The automation generates decision options that the human can select. Once an option is selected, the automation implements it.  |
| Blended Decision Making   | The automation generates an option, selects it, and executes it if the human consents. The human may approve of the option selected by the automation, select another, or generate another option.    |
| Rigid System              | The automation provides a set of options and the human has to select one of them. Once selected, the automation carries out the function.   |
| Supervisory Control       | The automation selects and carries out an option. The human can have input in the alternatives generated by the automation.   |
| Automated Decision Making | The automation generates options, selects, and carries out a desired option. The human monitors the system and intervenes if needed (in which case the level of automation becomes Decision Support). |
| Full Automation           | The system carries out all actions.   |

This includes engaging with dynamic systems of systems environments where an autonomous system may cooperate with a system not envisaged at design time.  
*How can systems that are required to demonstrate that they meet their requirement be specified flexibly enough to permit adaptive behaviours?*

2. **Security:** Any unmanned system has the potential to be used for illegitimate purposes by unscrupulous third parties who could exploit security vulnerabilities to gain control of the system or sub-systems. Any system that has the potential to cause harm from such actions must have security designed in from the start to ensure that the system can be trusted to be resilient from cyber attack. Current accreditation schemes rely on a security assessment of a known architecture and there are mutual accreditation recognition schemes that could be encoded in dynamic discovery handshake protocols. This would produce a secure network assured through the accreditation of its component systems. For example, the

Multinational Security Accreditation Board (MSAB) deals with Combined Communications Electronics Board (CCEB) and NATO Accreditations to provide security assurance of internationally connected networks. Encoding such agreements into secure handshakes could enable dynamic accreditation of autonomous systems cooperating in a coalition environment. It is not known whether these have been demonstrated, so the question is: *Can autonomous systems be designed to understand the security situation when interfacing with known or unknown systems?*

Need to check in with JP/JGF on status of JANUS. IIRC Janus dropped the whole idea of negotiating capabilities

3. **Verification and Validation of a Flexible Specification:** Following on from the description of a flexible specification, establish that the AS conforms and performs in accordance to the specification. This has direct implication for the trust in the resultant system. How can systems demonstrate that they will behave acceptably when the environment is unknown?
4. **Trust Modelling and Metrics:** This could be argued as part of the Verification and Validation of the system. However, models are increasingly being embedded into system design as a reference. Thus it is useful to consider this element separately. *How can trust be modelled sufficiently to span the space of most potential behaviours to help ensure that systems will be trusted when moved into operational environments? Can this be measured to allow comparison and minimum requirements set?*
5. **Certification:** The certification requirements placed on specific systems will vary depending on domain and national approaches to certification. However, the common element in the requirement for certification is that a certified system is deemed as sufficiently trustworthy for use within its context of certification. Additionally Certification also relies on the predictability of a system. Because the aim of autonomous systems is to deal effectively with uncertain environments, *can they (autonomous systems) be certified without being demonstrated in the environment within which they will adapt new behaviour?*

Design against and Compliance with existing standards can contribute significantly to the demonstrable trustworthiness of any systems design. If a system has been designed to a Standard then it has known properties that have been accepted as good practice. However, current standards do not address the issue of the five areas listed above.

Need to squeeze in something about Block 4 above is the focus of this work. Possibly could live in the conclusions

There are three main organisations that are developing or have developed assurance standards for Unmanned Systems in commercial, civil and military applications:

- NATO Standardization Office (NSO)

---

| LOI | Description  |
|-----|--|
| 1   | Indirect receipt/transmission of <b>UAV</b> related payload data   |
| 2   | Direct receipt of <b>ISR</b> data where direct covers reception of <b>UAV</b> payload data by the UCS when it has direct communication with the <b>UAV</b> |
| 3   | Control and monitoring of the <b>UAV</b> payload in addition to direct receipt of <b>ISR</b> /other data   |
| 4   | Control and monitoring of the <b>UAV</b> , less launch and recovery  |
| 5   | Launch and Recovery in addition to LOI 4   |

---

**Table 2.7:** Levels of Interoperability for STANAG 4586 Compliant UCS [40]

- Society of Automotive Engineers (SAE)
- American Society of Testing and Materials (ASTM)

**NATO Standardization Office** Faced with the growing adoption of similar but disparate **UAV** systems within NATO territories and coalition nations, STANAG 4586[40] was promulgated in 2005 and defined a logistic and interoperability framework to provide commonality in the command and control architecture and implementations of **UAV**/Ground station communications.

This included a particularly interesting development in the form of **Society of Automotive Engineers (SAE) Vehicle Specific Module (VSM)** interoperability, whereby existing systems could be grandfathered into STANAG 4586 compliance by the addition of a **VSM** to operate as a protocol translator. This **VSM** could be mounted on the remote system directly, utilising a compliant **Data Link Interface (DLI)**, or mounted on the ground-based controller, retaining the proprietary **DLI** to the remote system. The standard describes five **Level of Interoperability (LOI)** for compliant **UAV** systems, shown in Table 2.7. This structure has been criticised as being short sighted and at odds with the reality of modern and proposed autonomous vehicle operations [41], specifically that in modern autonomous systems, there is no such thing as direct control or Operator-in-the-loop, especially in the case of **Beyond Line of Sight (BLOS)** systems, and that in increasingly autonomous systems, operation is done as **Human Supervisory Control (HSC)**, or more commonly described as Operator-on-the-loop, whereby the operator interacts with the intermediate autonomous system and that autonomous system eventually performs that task on the hardware.

Further, the standard predominantly deals with a one-to-one mapping between operators and nodes, when this is quite against the current state of the art; greater focus is being made in collective and collaborative assignment and having a single operating agent managing a group of autonomous nodes in-field, and handing off vehicle management responsibilities to the individual nodes.

**SAE** The AS-4 steering group is responsible for the development and maintenance of the **Joint Architecture for Unmanned Systems (JAUS)** standards, which provide several service sets for Inter-System cooperation and interoperability, either in the form of a specified design language (JSIDL<sup>4</sup>) or as a direct framework implementation, such as the **JAUS** Mobility, Mission Spooling, Environment Sensing, or Manipulator Service Sets<sup>5</sup>. This provides a stack-like interoperability model akin to the OSI inter-networking standard, providing logical connections between common levels across devices regardless of how subordinate layers are implemented. Importantly, **JAUS** service models are open-sourced under the BSD-license, and a development toolkit is available for anyone to develop **JAUS**-compatible communications and control protocols[42].

It is also important to note that **JAUS** is part funded, and heavily utilised by, US Army and Marine Robotic Systems Joint Project Office (RS-JPO), which manage the development, testing, and fielding of unmanned (ground) systems for those respective forces. This includes now legacy M160 mine clearance platform and the highly popular (both with forces and their in-field operators) iRobot Packbot inspection and **Explosive Ordnance Disposal (EOD)** family of robotic platforms.

Needs references

**ASTM** The **ASTM** F38 committee has developed a **Line of Sight (LOS)**, single-asset-single-operator stove-piped framework for Unmanned Air Systems that is too constrained in scope for applicability to a more heterogeneous operating environment[43]. However, the F41 Committee, focused on **UMVSSs** has collectively developed a range of interoperable standards, covering Communications, Autonomy and Control, Sensor Data Formats, and Mission Payload Interfacing. Of particular interest is the Autonomy and Control standard which highlighted a requirement on the vehicle system to be able to recognise an authorised client, be that a human operator or an additional collaborating vehicle [44]. Further, the standard states that the responsibility of the safety and integrity of any payload remains with the vehicle. This standard was withdrawn in 2015 due to **ASTM** regulations requiring standards to be updated within 8 years of approval, and has no direct replacement within **ASTM**, but stands as a useful guiding perspective on autonomy standards within industry.

### Summary of Design Trust

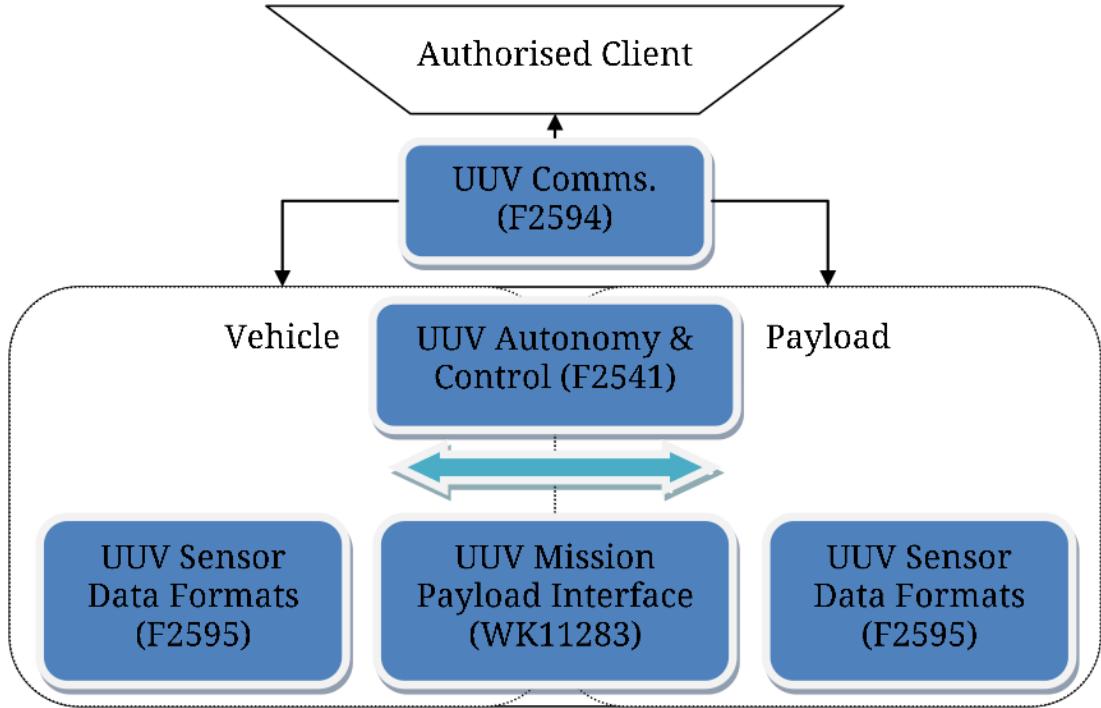
The implications of trust in autonomy beyond securing communications and data are an area in need of further research (BAE Systems, 2013. Maritime Autonomy Final Report - Combined Response,)

Need to check security status of this source

Of particular concern is the verification of autonomous behaviours. Technology Readiness Level deficiencies were identified in the Maritime Capability Contribution of

<sup>4</sup> **JAUS** Service Interface Definition Language

<sup>5</sup> SAE AS6009, AS 6062, AS 6060, and AS 6057 respectively



**Figure 2.3:** ASTM F41 UMVS Architecture (with relevant substandards in parenthesis)

Unmanned Systems (MCCUS) Osprey Phase 1 report(Clark, H. et al., 2012. Maritime Capability Contribution of Unmanned Systems,)

Need to check security status of this source

, with a particular focus on failsafe behaviour. The addition of increased on-board autonomy in MUxS, properly understood and verified, would greatly improve this future capability, similar to recent developments in the UAS arena[41].

There are opportunities for increased decentralisation and in-field collaboration(Walton, R., 2012. Maritime Autonomy PDR Pack.)

Need to check security status of this source

, however, difficulties in Trust between human operators and autonomous systems have already been clearly identified[45],and this has been demonstrated by the recent decision by the German government to renege on its €500M investment in the Euro Hawk programme, due to concerns about civil certification of the onboard autonomy[46] In order for these new distributed structures to be relied upon to provide operational performance, reliability and to maintain in-field situational awareness, vulnerabilities to disruption, interruption, and subversion need to be understood and minimised.

### 2.2.5 Operational Trust

#### Summary of Human Factors impacting Operational Trust in Defence Contexts

When dealing with human supervision of autonomous or semi-autonomous systems, there is an inherent conflict between the expectations of the operator, and the hopes of system architects. System architects aim to provide more and more information to the operator to justify a systems operation, and Operators in reality need less and less information to be efficient when things are going well, and responsive in a dynamic environment. This places huge demands on Human Interface design and indeed on communications design to provide this timely, relevant, interactive connection between any autonomous system and the end operator(s). Recent work has presented the idea of taking user interface (UI) inspiration from the entertainment sector, in terms of UI best practises developed over two decades of Real-Time Strategy game development [47], and follow up work into automated mission debrief demonstrated that such operational support could improve causal situational awareness of an operator when compared to a human-baseline [48]. In terms of the human factors challenges (See Appendix B for a discussion of these challenges), they are often contradictory in their direction, particularly when contrasting between Adaptive Automation and Cognitive Biases challenges. This is a key part of the “soft trust” perspective, where the operators and commanders need to be able to implicitly and explicitly trust the operation of a remote system with limited feed-back bandwidth, high latency, or long-term operation such that direct remote operation is infeasible or undesirable. To be able to trust that systems ability to continue on a course, survey an area, notify on detection of an anomaly, etc. is going to be the corner stone of any autonomous systems justification in the future.

### 2.2.6 Conclusions

ReDo this later

## 2.3 Trust in Autonomous MANETs

### 2.3.1 Trust Model Design Considerations

From the previous sections, Trust can be redefined as “the level of confidence one agent has in another to perform a given action on request or in a certain context”. Trust in the autonomous or semi-autonomous realm is the ability of a system to establish and maintain this level of confidence in itself or another systems’ operations.

There are five topics that are important to address in any MANETs trust model [49]:

Could really do with a better / additional cite than this...

- The trust model should be without infrastructure. Because the network routing infrastructure is formed in an ad-hoc fashion, the trust management can not depend on, e.g., a Trusted Third Party (TTP). There is no PKI, where some center nodes monitor the network, and publish illegal nodes periodically. In a MANET, there are no certification authorities (CA) or registration authorities (RA) with elevated privileges etc.
- The trust model should be anonymous because of the anonymity of mobile nodes in MANETs.
 

This isn't actually explained or justified in Kamvar so it may have been pulled out of his ass
- The trust model should be robust. That is, it can be robust to all kinds of unfriendly attacks and the network itself should not be susceptible to attacks by unfriendly nodes. Moreover, in the presence of malicious nodes, they may attempt to subvert the model in order to get an unfairly good trust value.
- The trust model should have minimal control overhead in accordance with computation, storage, and complexity.
- The trust model should be self-organized. MANETs are characterized to have dynamic, random, rapidly changing and multi-hop topologies composed of variably bandwidth-constrained links

### 2.3.2 Attacks on MANETs

Standard table

Emphasise Threat Surface discussion

### 2.3.3 Trust Management Frameworks

Distributed trust management frameworks for MANETs aim to detect, identify, and mitigate the impacts of malicious or selfish actors by generating, distributing and integrating per-node assessments and opinions to collectively self-police behaviour. From the settled upon definition of trust (From Subsection 2.3.1), these opinions are attempting to model the confidence of success in a particular actor for a particular future action.

This predictive behaviour attempts to solve four important problems (paraphrased from [19]):

- *Decision support* - For example; making informed routing table decisions based on past successes/failures.
- *Adaptability* - Ongoing prediction of the networks future trust states directly determines the risk faced by the network. Internalised knowledge of the expected risk can aid in selecting appropriate measures/ countermeasures such as automatically varying the level of authentication required for network activities.

- *Misbehaviour Detection* - Trust evaluation leads to a the natural policy that highly variable or low-trust nodes within a network should be subject to higher scrutiny; triggering this response indicates that a node is damaged or misbehaving.
- *Abstraction of Collective security characteristics* - Through per-node trust evaluation, the generalised trustworthiness of a set or subset of nodes can be derived to encapsulate the “health” of the network as a whole.

Various models and algorithms for describing trust and developing trust management in distributed systems, [Peer to Peer \(P2P\)](#) communities or wireless networks have been considered.

Taking some examples;

- *Hermes Trust Establishment Framework* uses a Bayesian Beta function to model per-link [Packet Loss Rate \(PLR\)](#) over time, combining “Trust” and “Confidence of Assessment” into a single value [50].
- *Objective Trust Management Framework (OTMF)* takes a Bayesian approach and introduces the idea of applying a Beta function to changes in the per-link [PLR](#) over time, combining “Trust” and “Confidence of Assessment” into a single value [51]. OTMF however does not appropriately combat multi-node-collusion in the network [52].
- *Trust-based Secure Routing* demonstrated an extension to [DSR](#), incorporating a Hidden Markov Model of the wider ad-hoc network, reducing the efficacy of Byzantine attacks, particularly black-hole attacks but is limited by focusing on single metric observation ([PLR](#)) [52, 53].
- *CONFIDANT*; presented an approach using a probabilistic estimation of normal observations, similar to [OTMF](#). Also introduced a greedy topology weighting scheme that internally weighted incoming trust assessments based on historical experience of the reporter [12].
- *Fuzzy Trust-Based Filtering*; presented a method using Fuzzy Inference to cope with imperfect or malicious recommendation based on a probabilistic estimation of performance using conditional similarity to classify performance using overlapping Fuzzy Set Membership functions to collaboratively filter reputations across a network [54].
- *Multi-parameter Trust Framework for MANETs (MTFM)* uses a number of communications metrics together for form a vector of trust, apply grey information theory to allow a system to detect and identify the tactics being used to undermine or subvert trust [13].

### 2.3.4 Single Metric Trust Frameworks

Where  $\alpha$  and  $\beta$  represent the number of successful and unsuccessful interactions respectively.

Expand background detail on more frameworks

The Hermes trust establishment framework [50] uses Bayesian reasoning to generate a posterior distribution function of “belief”, or trust, given a sequence of observations of that behaviour,  $p(B|O)$ (2.1).

$$p(B|O) = \frac{p(O|B) \times p(B)}{\rho} \quad (2.1)$$

Where  $p(B)$  is the prior probability density function for the expected normal behaviour, and  $\rho$  is a normalising factor.

This  $\rho$  bugs me; it should really be  $p(O)$  based on Bayes Theorem

Due to its flexibility and simplicity, Hermes assumes that  $p(B)$  is a Beta function ((2.2)), and therefore the evaluation of this trust assessment is based around the expectation value of the distribution (2.4) where  $\alpha$  and  $\beta$  represent the number of successful and unsuccessful interactions respectively for a particular node  $i$ .

$$\text{beta}(p|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} \quad (2.2)$$

$$E(p) = \frac{\alpha}{\alpha + \beta} \quad (2.3)$$

where  $0 \leq p \leq 1; \alpha, \beta > 0$

A secondary measurement of the confidence factor of the trust assessment  $t$  is generated as (2.5) and these measurements are combined to form a “trustworthiness” value  $T$  (2.6).

$$t_i \rightarrow E[\text{beta}(p|\alpha, \beta)] = \frac{\alpha_i}{\alpha_i + \beta_i} \quad (2.4)$$

$$c_i = 1 - \sqrt{\frac{12\alpha_i\beta_i}{(\alpha_i + \beta_i)^2(\alpha_i + \beta_i + 1)}} \quad (2.5)$$

$$T_i = 1 - \frac{\sqrt{\frac{(t_i-1)^2}{x^2} + \frac{(c_i-1)^2}{y^2}}}{\sqrt{\frac{1}{x^2} + \frac{1}{y^2}}} \quad (2.6)$$

In (2.6),  $x$  and  $y$  are constants to weight the two-dimensional polar mapping of trust and confidence assessments  $(t_i, c_i)$ , and from [50], are taken as  $x = \sqrt{2}, y = \sqrt{9}$ .

This makes absolutely no sense without a few diagrams

Upon this per-node assessment methodology, OTMF overlays an observation distribution protocol so as to make the measurements  $\alpha_i$  and  $\beta_i$  representative of the direct

and 1-hop networks observations of the target node  $i$ , as well as expiring old observations from assessment and eliminating observations from “untrustworthy” nodes.

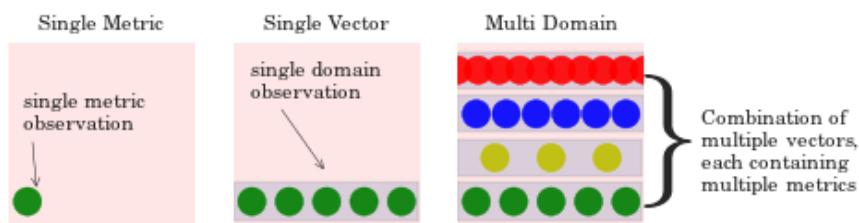
Want at least CONFIDANT and Fuzzy in here for contrast

To date this work has been mostly limited to terrestrial, RF based networks. There are many situations where the observed metrics will include significant noise and occur at irregular, sparse, intervals. Conventional approaches such as probabilistic estimation do not produce trust values that reflect the underlying reality and context of the metrics available, as they require a-priori assumption that the trust value under exploration has an expected distribution, that that distribution is mono-modal, and the input metrics are binary. In scenarios with variable, sparse, noisy metrics, estimating the distribution is difficult to accomplish a-priori. These single metric TMFs provide malicious actors with a significant advantage if their activity is undetectable by that one assessed metric, especially if the attacker is aware of the observed metric in advance.

The objective of operating a TMF is to increase the confidence in, and efficiency of, a system by reducing the amount of undetectable negative operations an attacker can perform. In the case where the attacker can subvert the TMF, the metric under assessment by that TMF does not cover the threat mounted by the attacker. In turn, this causes a super-linearly negative effect in the efficiency of the network as the TMF is assumed to have reduced the possible set of attacks when in fact it has only made it more advantageous to attack a different aspect of the networks operation. An example of such a behaviour would be the case in a TMF focused on PLR where an attacker selectively delays packets going through it, reducing the overall throughput of one or more network routes. Such behaviour would not be detected by the TMF.

### 2.3.5 Multi-Metric Trust Frameworks

Given the potential incentives to a selfish attacker and potential threats to trust and fairness in sparse, noisy, and constrained environments, single metric trusts discussed above do not suitably cover the exposed threat surface.



**Figure 2.4:** The inclusion of additional metrics and domains in trust assessment reduces the systems exposed threat surface

Replace Fig. 2.4 with vector one

A multi-metric approach may be more appropriate to capture and monitor the realities of harsh and sparse communications environments.

**MTFM**[13] uses Grey Theory (see [Appendix C](#)) to perform cohort based normalization of metrics at runtime, providing a “grey relational grade” of trust compared to other observed nodes in that interval for individual metrics, while maintaining the ability to reduce trust values down to a stable assessment range for decision support without requiring every environment entered into to be characterised. This presents a stark difference between the Grey and Probabilistic approaches. Grey assessments are relative in both fairly and unfairly operating networks. All nodes will receive mid-range trust assessments if there are no malicious actors as there is nothing “bad” to compare against, and variations in assessment will be primarily driven by topological and environmental factors. Guo et al. [13] demonstrated the ability of **Grey Relational Analysis (GRC)** to normalise and combine disparate traits of a communications link such as instantaneous throughput/load, received signal strength, etc. into a **Grey Relational Coefficient (GRC)**, or a “trust vector” in this instance.

The grey relational vector is given as

$$\begin{aligned}\theta_{k,j}^t &= \frac{\min_k |a_{k,j}^t - g_j^t| + \rho \max_k |a_{k,j}^t - g_j^t|}{|a_{k,j}^t - g_j^t| + \rho \max_k |a_{k,j}^t - g_j^t|} \\ \phi_{k,j}^t &= \frac{\min_k |a_{k,j}^t - b_j^t| + \rho \max_k |a_{k,j}^t - b_j^t|}{|a_{k,j}^t - b_j^t| + \rho \max_k |a_{k,j}^t - b_j^t|}\end{aligned}\quad (2.7)$$

where  $a_{k,j}^t$  is the value of an observed metric  $x_j$  for a given node  $k$  at time  $t$ ,  $\rho$  is a distinguishing coefficient set to 0.5,  $g$  and  $b$  are respectively the “good” and “bad” reference metric sequences from  $\{a_{k,j}^t, k = 1, 2 \dots K\}$ , i.e.  $g_j = \max_k(a_{k,j}^t)$ ,  $b_j = \min_k(a_{k,j}^t)$  (where each metric is selected to be monotonically positive for trust assessment, e.g. higher throughput is presumed to be always better).

Weighting can be applied before generating a scalar value ([C.3](#)) allowing the detection and classification of misbehaviours.

$$[\theta_k^t, \phi_k^t] = \left[ \sum_{j=0}^M h_j \theta_{k,j}^t, \sum_{j=0}^M h_j \phi_{k,j}^t \right] \quad (2.8)$$

Where  $H = [h_0 \dots h_M]$  is a metric weighting vector such that  $\sum h_j = 1$ , and in unweighted case,  $H = [\frac{1}{M}, \frac{1}{M} \dots \frac{1}{M}]$ .  $\theta$  and  $\phi$  are then scaled to  $[0, 1]$  using the mapping  $y = 1.5x - 0.5$ . To minimise the uncertainties of belonging to either best ( $g$ ) or worst ( $b$ ) sequences in ([C.2](#)) the  $[\theta, \phi]$  values are reduced into a scalar trust value by  $T_k^t = (1 + (\phi_k^t)^2 / (\theta_k^t)^2)^{-1}$  [55]. **MTFM** combines this **GRC** with a topology-aware weighting scheme ([C.4](#)) and a fuzzy whitenization model ([C.5](#)).

There are three classes of topological trust relationship used; Direct, Recommendation, and Indirect, as discussed in [Subsection 2.1.4](#). Where an observing node  $n_i$  assesses the trust of another target node,  $n_j$ ; the Direct relationship is  $n_i$ ’s own observations  $n_j$ ’s behaviour. In the Recommendation case, a node  $n_k$  which shares Direct relationships

with both  $n_i$  and  $n_j$ , gives its assessment of  $n_j$  to  $n_i$ . In the Indirect case, similar to the Recommendation case, the recommender  $n_k$  does not have a direct link with the observer  $n_i$  but  $n_k$  has a Direct link with the target node,  $n_j$ . These relationships give node sets,  $N_R$  and  $N_I$  containing the nodes that have recommendation or indirect, relationships to the observing node respectively.

Fix equation links on this page after finishing grey stuff

$$\begin{aligned} T_{i,j}^{\text{MTFM}} &= \frac{1}{2} \cdot \max_s \{f_s(T_{i,j})\} T_{i,j} \\ &+ \frac{1}{2} \frac{2|N_R|}{2|N_R| + |N_I|} \sum_{n \in N_R} \max_s \{f_s(T_{i,n})\} T_{i,n} \\ &+ \frac{1}{2} \frac{|N_I|}{2|N_R| + |N_I|} \sum_{n \in N_I} \max_s \{f_s(T_{i,n})\} T_{i,n} \end{aligned} \quad (2.9)$$

Where  $T_{i,n}$  is the subjective trust assessment of  $n_i$  by  $n_n$ , and  $f_s = [f_1, f_2, f_3]$  given as:

$$\begin{aligned} f_1(x) &= -x + 1 \\ f_2(x) &= \begin{cases} 2x & \text{if } x \leq 0.5 \\ -2x + 2 & \text{if } x > 0.5 \end{cases} \\ f_3(x) &= x \end{aligned} \quad (2.10)$$

Plot and explain the point of Whitenization (or move these back to the appendix)

In the case of the terrestrial communications network used in [13], the observed metric set  $X = x_1, \dots, x_M$  representing the measurements taken by each node of its neighbours at least interval, is defined as  $X = [\text{packet loss rate}, \text{signal strength}, \text{data rate}, \text{delay}, \text{throughput}]$ .

Guo et al. [13] demonstrated that when compared against OTMF and Hermes trust assessment, MTFM provided increased variation in trust assessment over time, providing more information about the nodes' behaviours than packet delivery probability alone can.

## 2.4 Conclusion

In the next chapter, the marine communications environment will be studied, as will the current state of the art in the use of autonomy in specifically defence related maritime applications.

Actual Conclusion of Trust Background

# Chapter 3

# Maritime Communications and Operations

## 3.1 Introduction

Introduction to Maritime

## 3.2 Maritime Communications Environment

The key challenges of underwater acoustic communications are centred around the impact of slow and differential propagation of energy (RF, Optical, Acoustic) through water, and its interfaces with the seabed / air. The resultant challenges include; long delays due to propagation, significant inter-symbol interference and Doppler spreading, fast and slow fading due to environmental effects (aquatic flora/fauna; surface weather), carrier-frequency dependent signal attenuation, multipath caused by the medium interfaces at the surface and seabed, variations in propagation speed due to depth dependant effects (salinity, temperature, pressure, gaseous concentrations and bubbling), and subsequent refractive spreading and lensing due to that same propagation variation[56].

### 3.2.1 Mechanics of Acoustic Transmission

Unlike in RF energy transfer (where photons move through space to transmit energy from one place to another), acoustic waves are the result of mechanical perturbation of a medium where localised compressions and extensions pass energy across a medium through that medium's elastic properties. These “compression waves” propagate away from its source, and the rate of this propagation is the sound speed, velocity or  $c$ , measured in  $ms^{-1}$ . This is not to be confused with the fluid velocity corresponding to the instantaneous motion of particles in the medium.

Hydrophones, like their more common microphone equivalent in air, are fundamentally pressure sensors. Acoustic pressure is usually measured in *Pascals* ( $Pa/\mu Pa$ ). In the underwater environment, the dynamic range (difference between instantaneous high

and low pressure values) may be extremely high, often more than 10 orders of magnitude higher. As such, logarithmic notation is justified.

Best to discuss notation here

Useful acoustic signals are generally maintained vibrations rather than instantaneous pulses. They are characterised by their frequency  $f$  expressed in Hertz ( $Hz$ ) or by their Period ( $T$ ) in seconds. In commonly used underwater acoustics, used frequencies range from  $\approx 10Hz - 100kHz$  depending on application.[57].

As with all waves, the relationship between frequency, period and the wavelength is given as in (3.1). As such the generally used upper and lower bounds of wavelength in most applications is from  $1.5m@10Hz$  to  $0.015m@100kHz$ .

$$\lambda = cT = \frac{c}{f} \quad (3.1)$$

This wide range of frequencies and wavelengths allow for a diverse set of constraining factors; (Paraphrased from Lurton [58]).

- *Attenuation* in water; limiting the maximum usable range, which increases very rapidly with frequency
- *Dimensions* of sound source; which increase at lower  $f$  for a given transmission power
- *Spatial Selectivity* of sources and receivers as  $f$  increases, due to similarly increasing directivity of energy propagation.
- *Acoustic Response* of target surfaces (analogous to receiver gain in RF networks).

### 3.2.2 Velocity and density

Air has a baseline density of approximately  $1.3kgm^{-3}$ , and the speed of sound is typically static around  $340ms^{-1}$ . In sea water, acoustic wave velocity is close to  $c = 1500ms^{-1}$  (generally between  $1450ms^{-1}$  and  $1550ms^{-1}$  depending on temperature, pressure, salinity etc.). Similarly variable is sea water density, which is nominally  $\rho = 1030kgm^{-3}$ .

While the sea/air surface is (ideally) a simple refractive interface, the interface between open seawater and marine sediment is graduated, with density ranges between  $1200 - 2000kgm^{-3}$ . This results in refractive and reflective velocities in the sediment interface ranging from  $1500 - 2000ms^{-1}$ [58].

For comparison, the speed of light in air/water is  $2.99 \times 10^8 ms^{-1}$  and  $2.249 \times 10^8 ms^{-1}$  respectively.

this might be better as a table

Mackenzie [59] proposed a more accurate model of acoustic velocity incorporating archival data from 15 worldwide sites that takes Temperature, Salinity and Depth into consideration.

$$10 \log a(f) = 0.11 \cdot \frac{f^2}{1+f^2} + 44 \cdot \frac{f^2}{4100+f^2} + 2.75 \times 10^{-4} f^2 + 0.003 \quad (3.6)$$

**Figure 3.1:** Thorp's Absorption Model[57]

$$\begin{aligned} c = & 1448.96 + 4.591T - 5.304 \times 10^{-2}T^2 + 2.374 \times 10^{-4}T^3 \\ & + 1.340(S - 35) + 1.630 \times 10^{-2}D + 1.675 \times 10^{-7}D^2 \\ & - 1.025 \times 10^{-2}T(S - 25) - 7.139 \times 10^{-13}TD^3 \end{aligned} \quad (3.2)$$

Where  $T$  is the temperature in Celsius,  $S$  the salinity in parts per thousand, and  $D$  is the depth below the surface in meters.

Need to discuss Speed of Sound Profiles

### 3.2.3 Intensity and Power

The energy of an acoustic wave is encapsulated into its kinetic and potential parts; where its kinetic energy corresponds to the active motion energy of the particles in the medium, and the potential energy corresponding to the elastic potential of the medium in displacement/compression.

The acoustic intensity ( $I$ ) is the energy flux mean value per unit of surface and time (3.3) in Watts/ $m^2$  where  $p_0$  is the plane wave amplitude (pressure) and  $P_{rms} = p_0/\sqrt{2}$

$$I = \frac{p_0^2}{2\rho c} = \frac{p_{rms}^2}{\rho c} \quad (3.3)$$

### 3.2.4 Attenuation

The attenuation that occurs in an underwater acoustic channel over a distance  $d$  for a signal about frequency  $f$  in linear (3.4) and dB forms (3.5) is given as;

$$A_{aco}(d, f) = A_0 d^k a(f)^d \quad (3.4)$$

$$10 \log A_{aco}(d, f)/A_0 = k \cdot 10 \log d + d \cdot 10 \log a(f) \quad (3.5)$$

where  $A_0$  is a unit-normalising constant,  $k$  is a geometric spreading factor (commonly taken as 1.5 for practical use, but may be 2 for perfect spherical propagation or 1 for perfect plane-wave propagation), and  $a(f)$  is the absorption coefficient, that may be modelled in a variety of ways.

Thorp's formula (Equation 3.6) is very simple, only depending on  $f$ , and is designed to be most accurate about a temperature of 4°C at a depth of  $\approx 1Km$ . The Ainslie & McColm model is more complex, and incorporates the acidity of the water ( $H^+$ ) as well as temperature ( $T$ ), salinity ( $S$  in parts per trillion) but not depth (Fig. 3.2). The Fisher-Simmons model (Fig. 3.3) is significantly more complex, taking into account the

$$10 \log a(f) = 0.106 \frac{t_1 f^2}{t_1^2 + f^2} e^{\frac{H+8}{0.56}} + 0.52 \left(1 + \frac{T}{43}\right) \left(\frac{S}{35}\right) \frac{t_2 f^2}{t_2^2 + f^2} e^{\frac{-D}{6}} + 4.9 \times 10^{-4} f^2 e^{-(\frac{T}{27} + \frac{D}{17})} \quad (3.7)$$

Where

$$t_1 = 0.78 \sqrt{\frac{S}{35} e^{\frac{T}{26}}} \\ t_2 = 42 e^{\frac{T}{17}}$$

**Figure 3.2:** Ainslie & McColm Absorption Model

$$10 \log a(f) = A_1 P_1 \frac{t_1 f^2}{t_1^2 + f^2} + A_2 P_2 \frac{t_2 f^2}{t_2^2 + f^2} + A_3 P_3 f^2 \quad (3.8)$$

Where

$$A_1 = 1.03 \times 10^{-8} + 2.36 \times 10^{-10} \cdot T - 5.22 \times 10^{-12} \cdot T^2 \\ A_2 = 5.62 \times 10^{-8} + 7.52 \times 10^{-10} \cdot T \\ A_3 = (55.9 - 2.39 \cdot T + 4.77 \times 10^{-2} \cdot T^2 - 3.48 \times 10^{-4} \cdot T^3) \times 10^{-15} \\ t_1 = 1.32 \times 10^3 (T + 273.1) e^{\frac{-1700}{T+273.1}} \\ t_2 = 1.55 \times 10^7 (T + 273.1) e^{\frac{-3052}{T+273.1}} \\ P_1 = 1 \\ P_2 = 10.3 \times 10^{-4} \cdot P + 3.7 \times 10^{-7} \cdot P^2 \\ P_3 = 3.84 \times 10^{-4} \cdot P + 7.57 \times 10^{-8} \cdot P^2$$

**Figure 3.3:** Fisher-Simmons Absorption Model

effects of boric acid concentrations and dissolved magnesium sulphate. While there are several limitations on this model in terms of its being fixed at a salinity of 35 ppt and a pH of 8, as this model incorporates depth, temperature, distance and frequency, it is very attractive for research directed at high variability environments.

Possibly need to switch this with the Francois Garrison model which, depending on your source, is the refined version (or vice versa)

Regardless of the variations of particular attenuation models, comparing  $A_{aco}(d, f)$  with the RF Free-Space Path Loss model ( $A_{RF}(d, f) \approx \left(\frac{4\pi df}{c}\right)^2$ ), the impact of range on signal power is exponential underwater, rather than quadratic in terrestrial RF ( $A_{aco} \propto f^{2d}$  vs  $A_{RF} \propto (df)^2$ ). While both frequency dependant factors are quadratic, approximating the factors in (3.6),  $f \propto A_{aco}$  is at least 4 orders of magnitude higher than  $f \propto A_{RF}$

$$A_{RF}(d, f) \approx \left(\frac{4\pi df}{c}\right)^2 \text{ where } c \approx 3 \times 10^8 \text{ ms}^{-1} \quad (3.9)$$

**Table 3.1:** Contributing factors to Ocean Ambient Acoustic Noise

| Source            | Approximation   |
|-------------------|---|
| Turbulence        | $10 \log N_t(f) = 17 - 30 \log f$   |
| Shipping          | $10 \log N_s(f) = 40 + 20(s - 0.5) + 26 \log f - 60 \log(f + 0.03)$       |
| Wind Driven Waves | $10 \log N_w(f) = 50 + 7.5w^{\frac{1}{2}} + 20 \log f - 40 \log(f + 0.4)$ |
| Thermal Noise     | $10 \log N_{th}(f) = 15 + 20 \log f$                                      |

### 3.2.5 Ambient Noise Model

Ambient ocean noise can be assumed to be Gaussian with a continuous power spectral density in dB re  $\mu\text{Pa}$  per Hz, driven by four major factors, shown in [Table 3.1](#) [60].

check what  $s$  and  $w$  are in this

### 3.2.6 Multipath effects

Refractive lensing and the multi-path nature of the medium result in line of sight propagation being extremely unreliable for estimating distances to targets. The first arriving acoustic signal has as the very least curved in the medium, and commonly has reflected off the surface/seabed before arriving at a receiver, creating secondary paths that are sometimes many times longer than the first arrival path, generating symbol spreading over orders of seconds depending on the ranges and depths involved. Thus, the multi-path channel transfer function can be described by

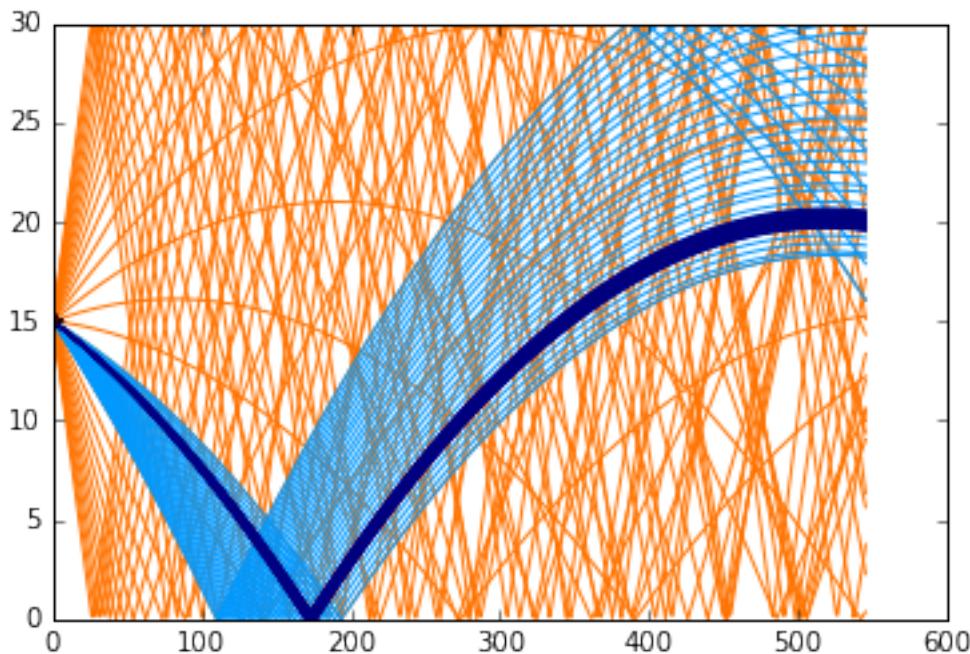
$$H(d, f) = \sum_{p=0}^{P-1} h(p) = \sum_{p=0}^{P-1} \Gamma_p / \sqrt{A(d_p, f)} e^{-j2\pi f \tau_p} \quad (3.10)$$

where  $\tau_p = d_p/c$ ,  $c \approx 1500 \text{ms}^{-1}$

where  $d = d_0$  is the minimal path length between the transmitter and receiver,  $d_p, p = \{1, \dots, P-1\}$  are the secondary path lengths,  $\Gamma_p$  models additional losses incurred on each path such as reflection losses at the surface interface, and  $\tau_p = d_p/c$  is the delay time.

### 3.2.7 Modelling and Simulation of the Acoustic Medium / Channel

Several toolkits exist in a variety of states that perform communications agent simulation, most notably the NS-2 / 3 family of frameworks and their addons. Some of these frameworks, such as SUNSET [61] and AquaTools [62].



**Figure 3.4:** Non-Linear Marine Propagation in an isothermal profile

Vectorise and Label

Beyond the NS family, there are many other communications and simulation modelling systems such as OpNet++[63] and MATLAB toolkits such as the AcTUP interface to the Ocean Acoustics Library.

expand this, justify AUVNetSim, reactive mobility, python compatibility, SimPy  
Etc.

### 3.2.8 Routing and Network Design for UANs

Forward Error Correction coding is used on such channels to minimise packet losses.

Summary of Akyildiz02/05

## 3.3 Marine Operations

### 3.3.1 Typical AUV mission profiles

Typical AUV missions, payloads, and available equipment

### 3.3.2 Potential Future Applications

Future Applications of AUVs

### 3.3.3 Need for Trust in Maritime Networks

As Autonomous Underwater Vehicle (AUV) platforms become more capable and economical, they are being used in many applications requiring trust. These applications are using the collective behaviour of teams or fleets of these AUVs to accomplish tasks [64]. With this use being increasingly isolated from stable communications networks, the establishment of trust between nodes is essential for the reliability and stability of such teams. As such, the use of trust methods developed in the terrestrial MANET space must be re-appraised for application within the challenging underwater communications channel.



## Chapter 4

# Assessment of TMF Performance in Marine Environments

### 4.1 Introduction

As MANETs grow beyond the terrestrial arena, their operation and the protocols designed around them must be reviewed to assess their suitability to different communications environments to ensure their continued security, reliability, and performance. With demand for smaller, more decentralised MANET systems in a range of domains and applications, as well as a drive towards lower per-unit cost in all areas, TMFs are going to be increasingly applied to resource constrained applications, as the benefits and efficiencies they present are significant. This work is primarily concerned with the analytical establishment of hard trust within a topologically dynamic network of autonomous actors. Beyond the constraints of the communications environment, knock on pressures in battery capacity, on-board processing, and locomotion simultaneously present opportunities and incentives for malicious or selfish actors to appear to cooperate while not reciprocating, in order to conserve power for instance. These multiple aspects of potential incentives, trust, and fairness do not directly fall under the scope of single metric trusts discussed above, and this context indicates that a multi-metric approach may be more appropriate. These increasingly decentralised applications present unique threats against trust management [64].

Previous research has established the advantages of implementing TMFs in 802.11 based MANETs, particularly in terms of preventing selfish operation in collaborative systems [11], and maintaining throughput in the presence of malicious actors [12]

To date this work has been limited to terrestrial, RF based networks.

One area of application is the underwater marine environment, where extreme challenges to communications present themselves (propagation delays, frequency dependent attenuation, fast and slow fading, refractive multipath distortion, etc.)(Chapter 3). In addition to the communications challenges, other considerations such as command and control isolation, as well as power and locomotive limitations, drive towards the use of

teams of smaller and cheaper AUV platforms. In underwater environments, communications is both sparse and noisy. Therefore the observations about the communications processes that are used to generate the trust metrics, occur much less frequently, with much greater error (noise) and delay than is experienced in terrestrial RF MANETs.

In addition to the communications challenges, other considerations such as command and control isolation, as well as power and locomotive limitations, drive towards the use of teams of smaller and cheaper AUVs. As such, the use of trust methods developed in the terrestrial MANET space must be re-appraised for application within the underwater context [23]. Many UANs use MANET architectures, however the marine environment presents new challenges for trust management frameworks that have been developed for use in conventional (i.e. Terrestrial RF) MANETs.

It is shown that single metric trust systems are not directly suitable for the marine context in terms of the different threat and cost scenario in that environment.

These single metric TMFs provide malicious actors with a significant advantage if their activity does not impact that metric. In the case where the attacker can subvert the TMF, the metric under assessment by that TMF does not cover the threat mounted by the attacker. This causes a significant negative effect on the efficiency of the network, as the TMF is assumed to have reduced the possible set of attacks when it has actually made it more advantageous to attack a different part of the networks operation. An example of such a situation would be in a TMF focused on PLR where an attacker selectively delays packets going through it, reducing overall throughput but not dropping any packets. Such behaviour would not be detected by the TMF.

For the purposes of this work, from those TMFs discussed in Subsection 2.3.3, Hermes trust establishment, OTMF and MTFM are selected as indicative single and multi metrics frameworks for comparison, as Hermes captures the core operation of a pure single metric assessment methodology and OTMF provides a comparison that combines assessments from across nodes to develop trust opinions.

From the discussion on the nature of the communications environment in Section 3.2, it's clear that before assessing communications metrics a simulated underwater environment, appropriate scaling factors must be found that are realistic from an application perspective but are also comparable in some form to the MANET case.

## 4.2 Modelling of UAN network

### 4.2.1 Mobility, Topology, and Communications

Four mobility patterns are investigated:

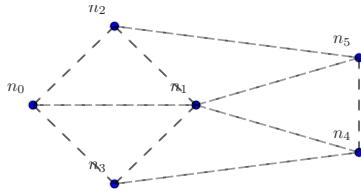
1. All Nodes Static
2. Malicious node mobile
3. Malicious node mobile, all other nodes static

#### 4. All nodes mobile

For this case, the mobility model used is a random walk on the nodes modeled kinematic response, i.e. the node periodically picks a spherically normalised random direction in the XY plane. Maximum node speed (limited by kinematic acceleration/turning constraints) is  $1.5ms^{-1}$ .

The six nodes are initially arranged as per Fig. 4.1 with each node on average 100m from each other as per [13]. The use of six nodes and the particular layout enables the investigation of the three trust relationships based on minimum path topologies, such that the node generating the trust assessments,  $n_0$  has Direct, Recommendation, and Indirect trust assessments of  $n_1$  available to it from itself,  $[n_2, n_3]$ , and  $[n_4, n_5]$  respectively. (See Section 2.1.4)

Collaborations with NATO Centre for Maritime Research and Experimentation (CMRE) in La Spezia, and Defence Science and Technology Laboratorys (DSTLs) Naval Systems Group inform that this is a practical team-size for environmental and defence applications.



**Figure 4.1:** Initial layout with nodes spaced an average of 100m apart

#### 4.2.2 Simulation Background

Simulations were conducted using a Python based simulation framework, SimPy [65], with a network stack built upon AUVNetSim [66], with transmission parameters (Table 4.1) taken from and validated against [57], [67] and [62]

it would be worth while going through this verification explicitly as an appendix

Given the differences in delay and propagation between RF and marine networks, it would not be expected that the same application rates (e.g. packet emission rates or throughput) and node separations are equally stable in this environment. Therefore, a zone of performance is characterised within which the network has stable operation.

### 4.3 Establishing Scale Factors in Communications Rate

In this section the simulated communications environment is characterised to establish an optimal packet emission rate for comparison against [13]. This optimal emission rate is taken to be an emission rate that provides reasonable network stability and protection from network saturation. Network saturation is the point at which a network can no

**Table 4.1:** Comparison of system model constraints as applied between Terrestrial and Marine communications

| Parameter                    | Unit             | Terrestrial       | Marine          |
|------------------------------|------------------|-------------------|-----------------|
| Simulated Duration           | s                | 300               | 18000           |
| Trust Sampling Period        | s                | 1                 | 600             |
| Simulated Area               | km <sup>2</sup>  | 0.7               | 0.7-4           |
| Transmission Range           | km               | 0.25              | 1.5             |
| Physical Layer               |                  | RF(802.11)        | Acoustic        |
| Propagation Speed            | m/s              | $3 \times 10^8$   | 1490            |
| Center Frequency             | Hz               | $2.6 \times 10^9$ | $2 \times 10^4$ |
| Bandwidth                    | Hz               | $22 \times 10^6$  | $1 \times 10^4$ |
| MAC Type                     |                  | CSMA/DCF          | CSMA/CA         |
| Routing Protocol             |                  | DSDV              | FBR             |
| Max Speed                    | ms <sup>-1</sup> | 5                 | 1.5             |
| Max Data Rate                | bps              | $5 \times 10^6$   | $\approx 240$   |
| Packet Size                  | bits             | 4096              | 9600            |
| Single Transmission Duration | s                | 10                | 32              |
| Single Transmission Size     | bits             | $10^7$            | 9600            |

longer successfully deliver the offered load<sup>1</sup> presented to it to the relevant destinations (throughput), and is characterised by a peak and a subsequent decline in the throughput of the network when varying the packet emission rate.

In order to establish the point at which the network becomes saturated due, a range of packet emission rates were explored between 0.01 packets per second (pps), equivalent to 96 bits of offered load per node, up to 0.07 pps (672 bps per node). Initial node separation was set as per Guo at 100m, and each simulation is run 16 times, with each instance modelling a 8 hour mission time.

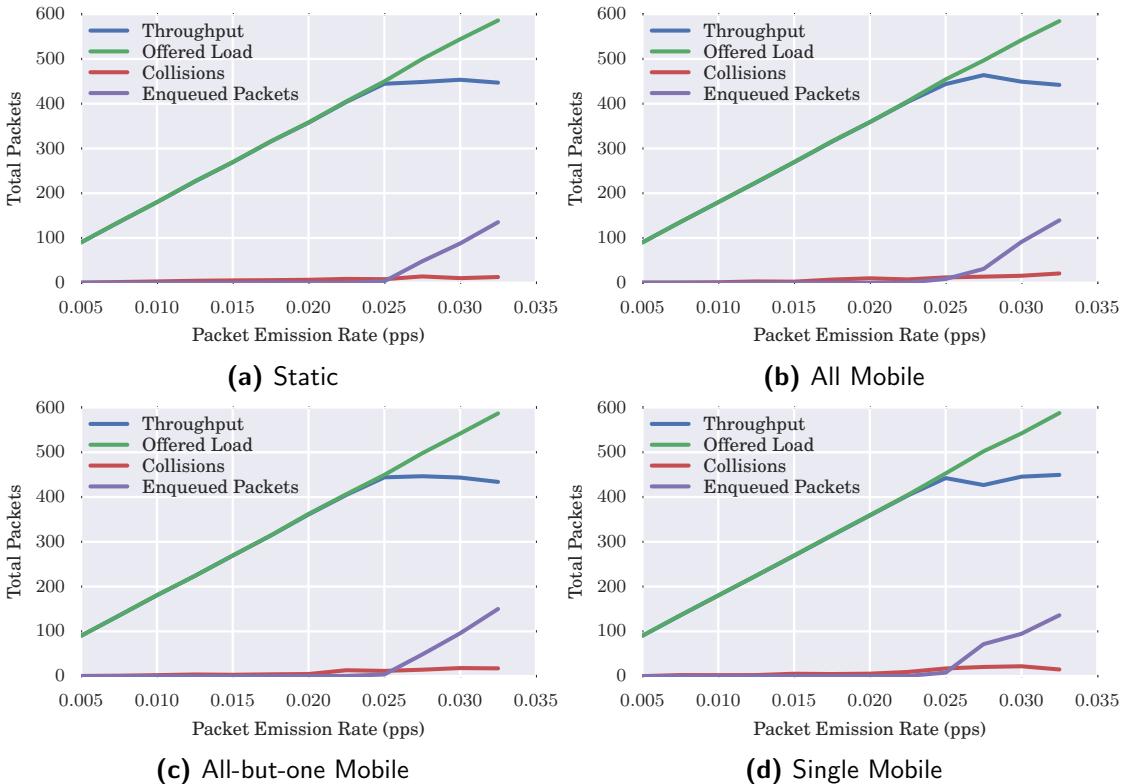
Need to have a discussion about mission configurations at some point

Looking first at the Static mobility case, where all nodes are stationary; from Fig. 4.3a it is already clear that the throughput curve, exhibits a saturation point close to 0.025 pps. Similarly in Fig. 4.3b, the precipitous drop in packet delivery probability beyond 0.025 pps, indicating that this is a strong candidate value for an upper-limit to the safe operating zone in terms of packet emission in the small static case. From Fig. 4.3c, raising packet emissions above 0.25pps results in a significant increase in end-to-end delay. As per Table 4.1, the CSMA based Medium Access Control (MAC) incurs a certain amount of control overhead in the form of Request To Send (RTS) packets,

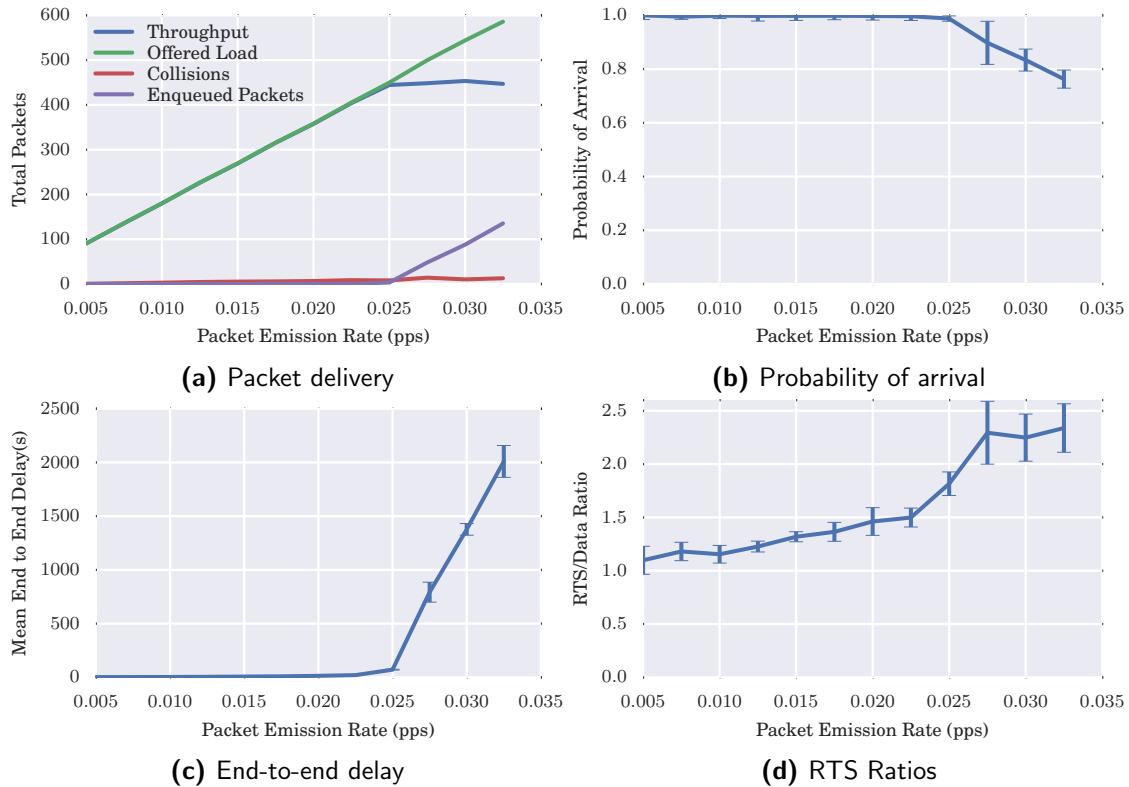
<sup>1</sup>It will become important to note that Offered Load in this case includes packet retransmissions

when a node attempts to acquire time in its neighbourhood. In Fig. 4.3d, the ratio of Control/Data packets increases linearly up to 1.5 until just before 0.025pps, and then accelerates to almost 2.5, further demonstrating that the network has become critically congested. It is worthwhile noting that in Fig. 4.3a that even as the saturation point is passed, packet collisions do not significantly increase, and that the saturation is in fact driven by contention in the medium rather than congestion-collisions.

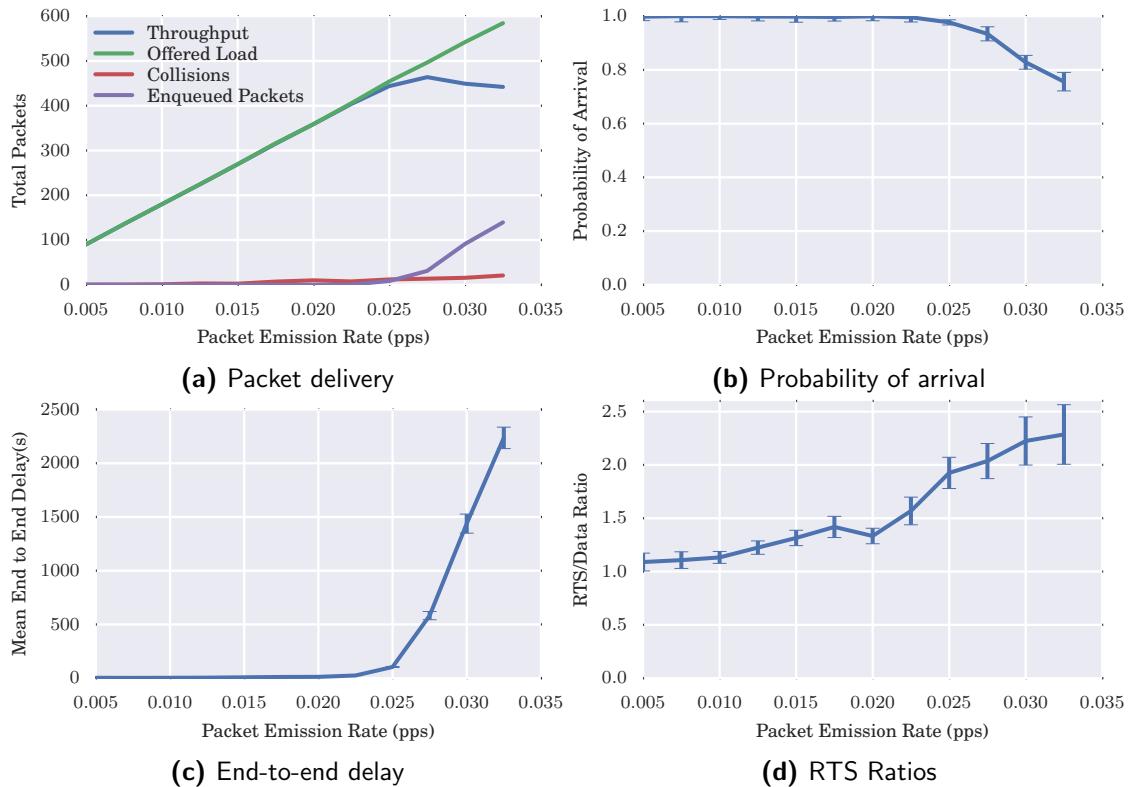
Results are also included from the remaining mobility cases (all nodes mobile; all-but-one node mobile; single mobile node), however from Figs. 4.2, 4.4- 4.6 that the throughput threshold behaviour is qualatitively similar regardless of mobility for this initial node separation.



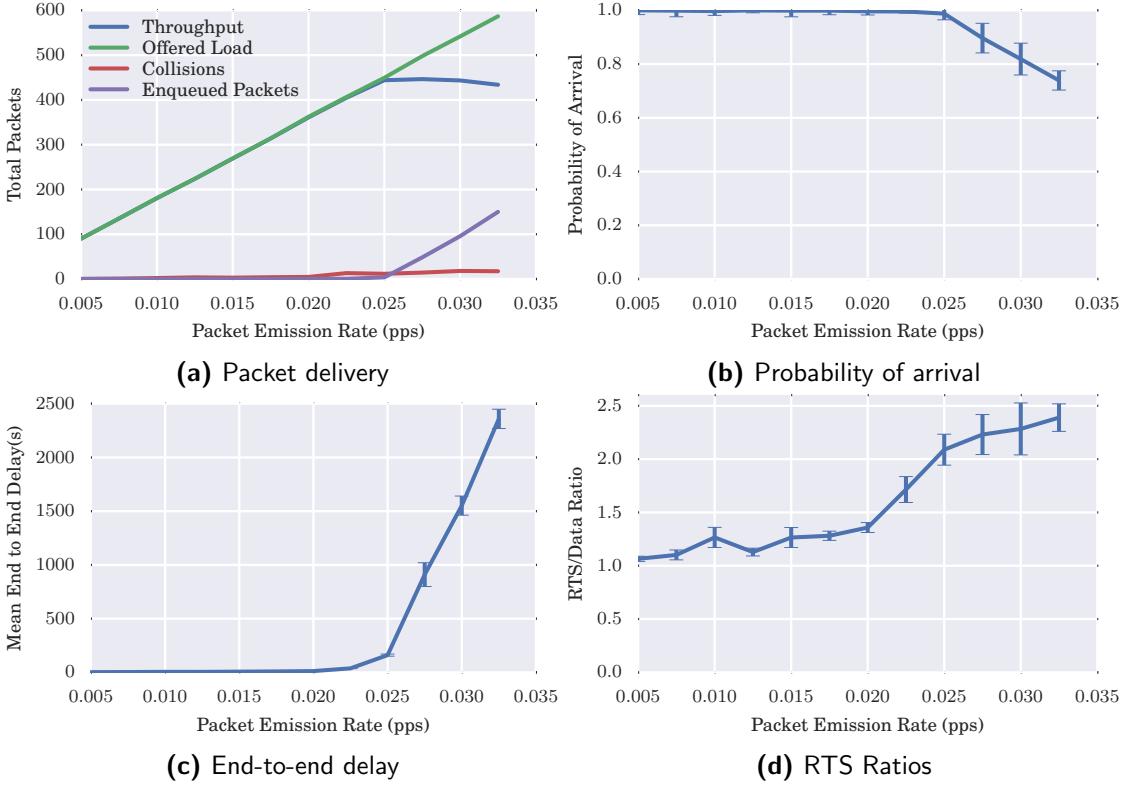
**Figure 4.2:** Throughput performance overview for all mobilities under varying emission rates  
IS THIS ENOUGH?



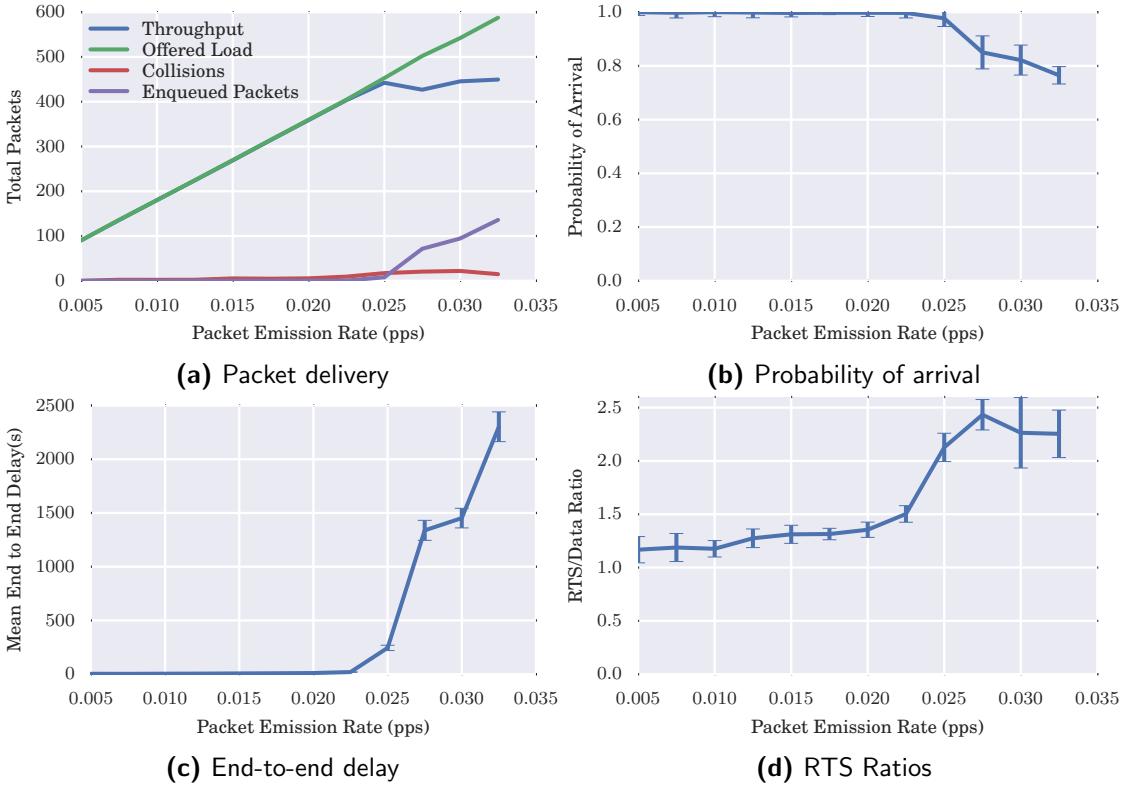
**Figure 4.3:** Network performance varying packet emission rates for the static case



**Figure 4.4:** Network performance varying packet emission rates for the all mobile case



**Figure 4.5:** Network performance varying packet emission rates for the all-but-one mobile case



**Figure 4.6:** Network performance varying packet emission rates for the single mobile case

### 4.3.1 Scale Factors in Physical Node Distribution

In this section the effect of node-separation scaling on communications operation is characterised for comparison against [13]. This is particularly important considering the significant scale factor differences in terms of the speed of propagation in the medium, and the range of potential desired operation.

From [Table 4.1](#), the operating transmission range of acoustic is  $\approx 6$  times further than 802.11, indicating that a suitable operating environment will have an area  $\approx \sqrt{6}$  times the area of the 802.11 case. Therefore, a reasonable experimental range would have an upper bound of performance around this scaling factor, where nodes are approximately 400m apart.

According to Xu, RTS/CTS handshake functionality cannot operate well as interference protection at node separations beyond 0.56 times the transmission range [68]. In the case of marine acoustic transmission at the stated power output, above  $1500m \times 0.56 = 840m$ , handshake overheads should begin to dominate channel access.

redo these graphs with wider separations 1000m

This is due to reduced channel availability due to collisions, which are then due to a much longer potential contention period between nodes.

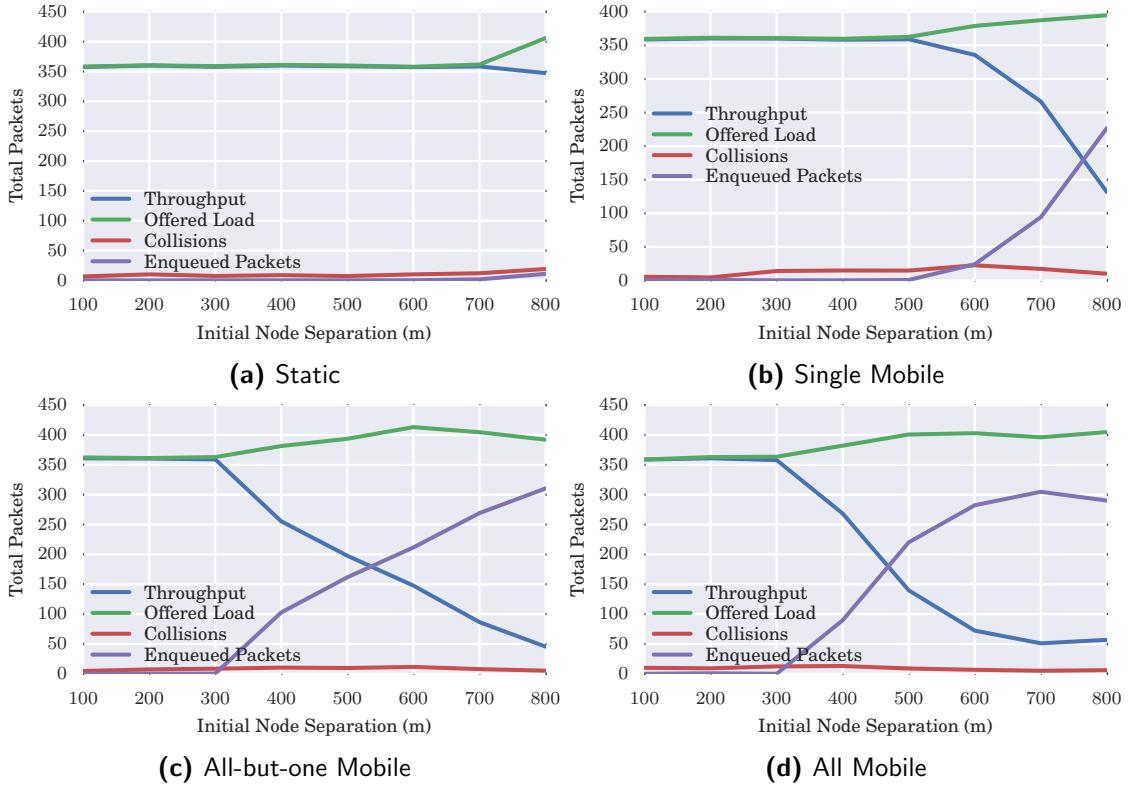
A reasonable range around this is to scale from 100m apart on average to 800m, and from the previous section, a packet emission rate of 0.02pps (slightly below the 0.025pps saturation threshold) is used to explore this space. The “environment” of the simulations is also scaled in accordance with the node scaling, based on an initial environmental “water-box” of 1km for the 100m node separation, i.e. the water-box is consistently ten times larger than the initial node separation.

In the case where all nodes remain static, increasing node separation does not significantly impact throughput, delay, delivery probability or [RTS](#) ratios until rising above 700m ([Fig. 4.8](#)), nearly double our initial estimate of where an appropriate separation zone would be.

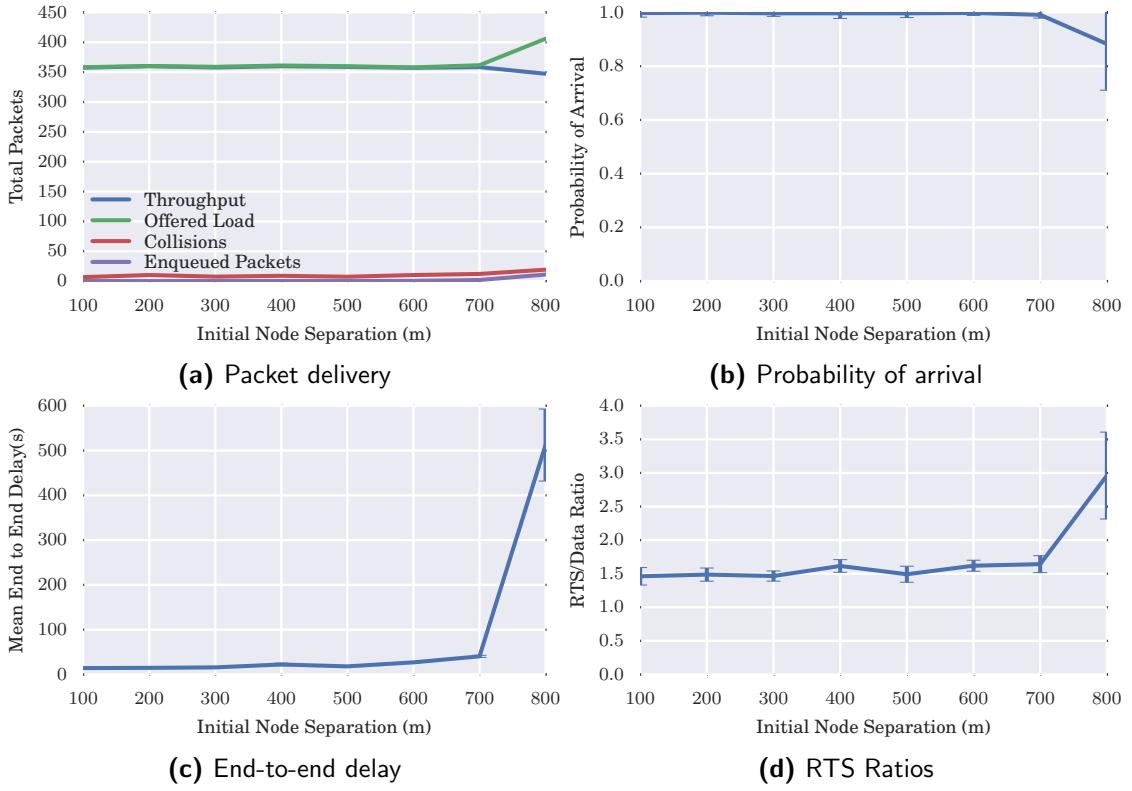
The other mobility cases tell a very different story; as can be seen in [Fig. 4.7b](#), where adding a single mobile node to the network induces a saturation-style response at 500m, and this drops further in [Fig. 4.7c](#) and [Fig. 4.7d](#), reducing the separation of saturation at this emission rate to just 300m.

Another aspect of these results to highlight is that the Offered Load presented to the network *increases* beyond the collapse of the throughput curve. This indicates that there is a subtly different saturation behaviour with respect to separation than the simple congestion argument with respect to packet emission rate; packets are simply taking too long to cross the increasingly sparse network and in-transit packet routes are logically disconnected and require retransmission.

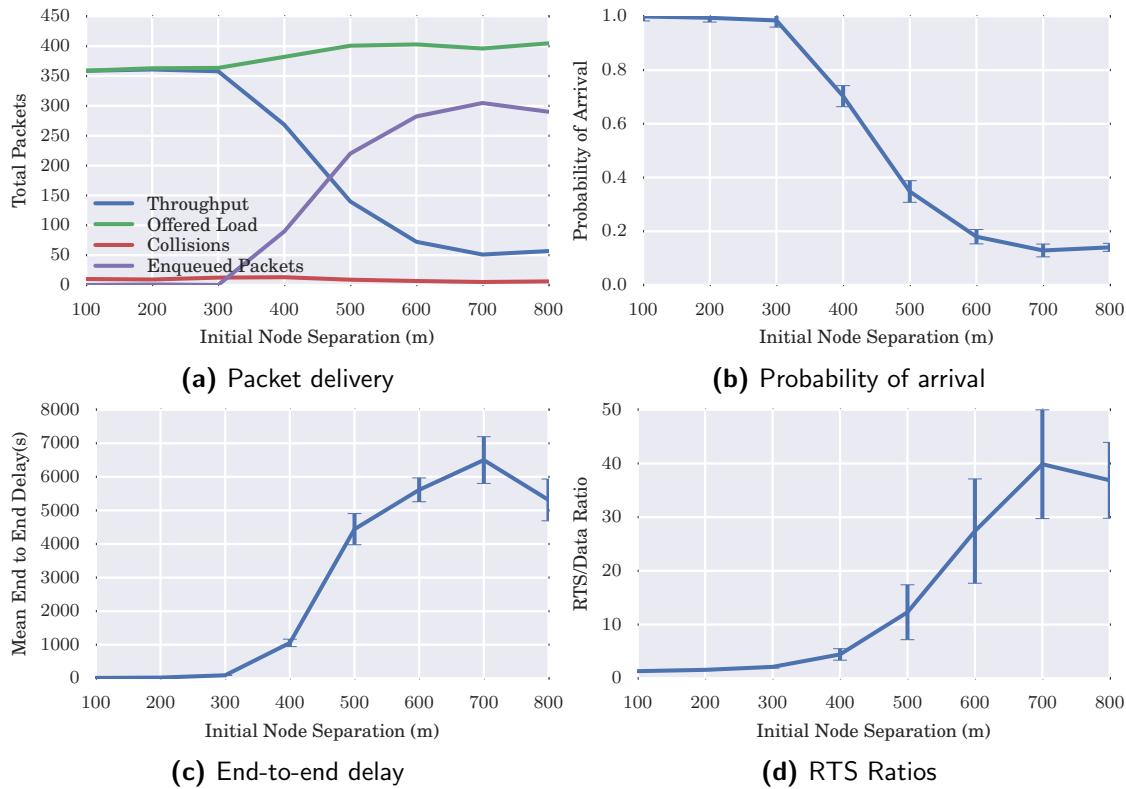
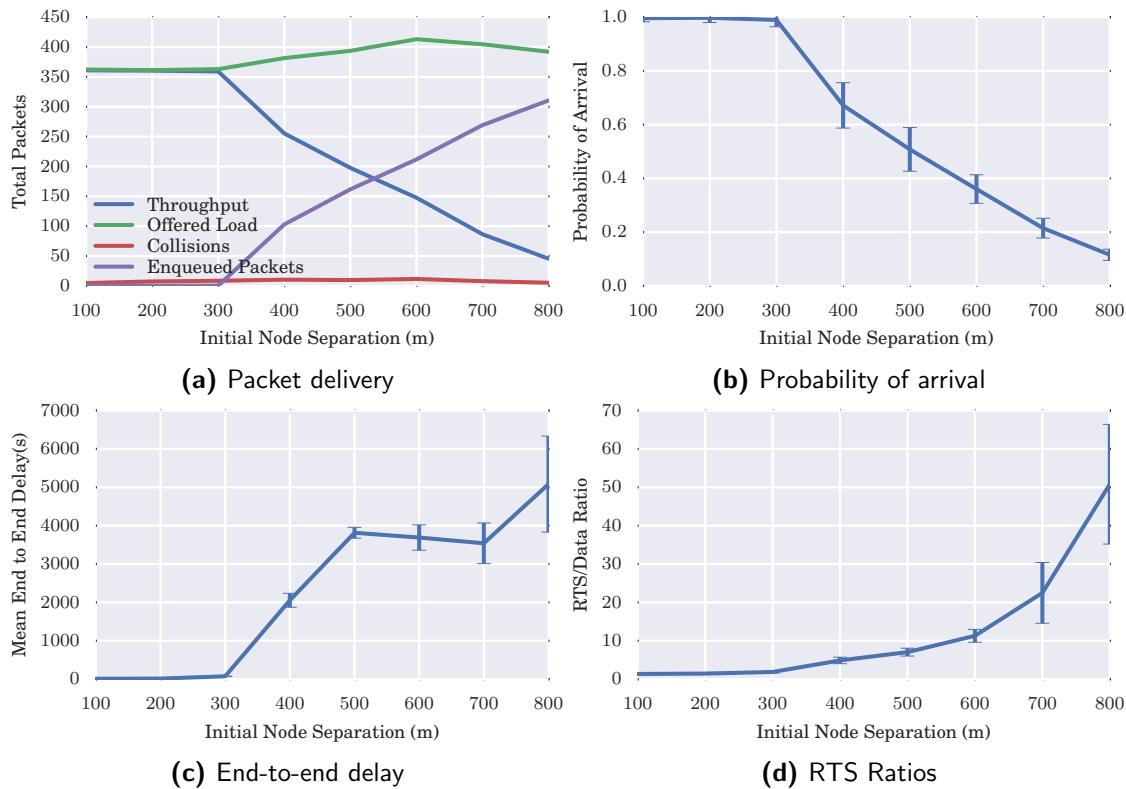
Another interesting aspect is the behaviour of the Enqueued Packet lines and e2e delay lines; They “Bump”; no idea why yet

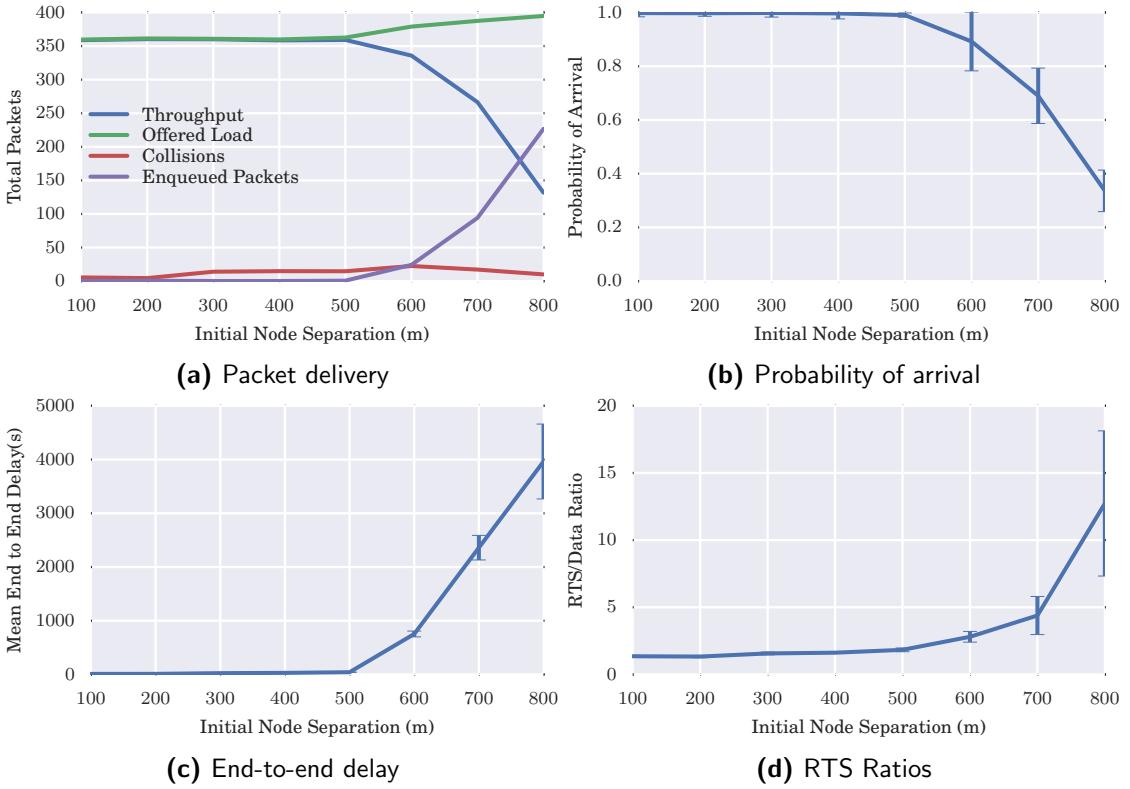


**Figure 4.7:** Throughput performance overview for all mobilities under varying separation *IS THIS ENOUGH?*



**Figure 4.8:** Network performance varying node separation for the static case

**Figure 4.9:** Network performance varying node separation for the all mobile case**Figure 4.10:** Network performance varying node separation for the all-but-one mobile case



**Figure 4.11:** Network performance varying node separation for the single mobile case

**Table 4.2:** Tabular view of data from Fig. 4.11, including ideal propagation time

| Initial Node Separation (m) | Delay(s)  | Probability of Arrival | RTS/Data Ratio | Ideal Delivery Time(s) |
|-----------------------------|-----------|------------------------|----------------|------------------------|
| 100                         | 10.3551   | 0.9977                 | 1.3546         | 1.0314                 |
| 200                         | 11.1631   | 0.9973                 | 1.3322         | 1.1029                 |
| 300                         | 24.2225   | 0.9983                 | 1.5650         | 1.1743                 |
| 400                         | 29.4864   | 0.9965                 | 1.6210         | 1.2457                 |
| 500                         | 41.7093   | 0.9904                 | 1.8331         | 1.3171                 |
| 600                         | 753.4040  | 0.8922                 | 2.8038         | 1.3886                 |
| 700                         | 2360.0826 | 0.6899                 | 4.3889         | 1.4600                 |
| 800                         | 3963.9830 | 0.3360                 | 12.7323        | 1.5314                 |

## 4.4 Combined Scale Factor Analysis

It's clear from the previous results that the relationship between emission rates, separations and mobilities is tightly coupled and not totally clear cut. To arrive at a more optimal operating region, a coupled analysis is performed across both emission rate and initial separation distance.

Given what has been discussed so far; it's clear that in identifying an appropriate operating region, it is important to not only ensure throughput, but that that throughput is timely. For instance, in [Fig. 4.11](#) (tabulated in [Table 4.2](#)), a small increase in separation beyond the apparent throughput-peak at 500m to 600m, which constitutes an increased ideal marine acoustic "time of flight" between nodes by 0.02s, increases the average actual delay by 1800%.

To capture these performance requirements, the feature scaled product of Throughput and Delay is taken and plotted against rate and separation in [Fig. 4.12](#).

This does NOT make for easy comparison between graphs as the scaling is different for each mobility, but I need to think about how to fairly solve this

$$V = |S| \times (1 - |D|) \quad (4.1)$$

For each scenario, the observed Throughput across the network ( $S$  in bytes) is normalised across all observations (i.e. each combination of Node Separation and Emission Rate), as is average end-to-end Delay ( $D$ ). The normalised delay is inverted ( $1 - |D|$ ) and the product of this and the normalised throughput is used as the basis of a two-dimensional linear interpolation shown in [Fig. 4.12](#).

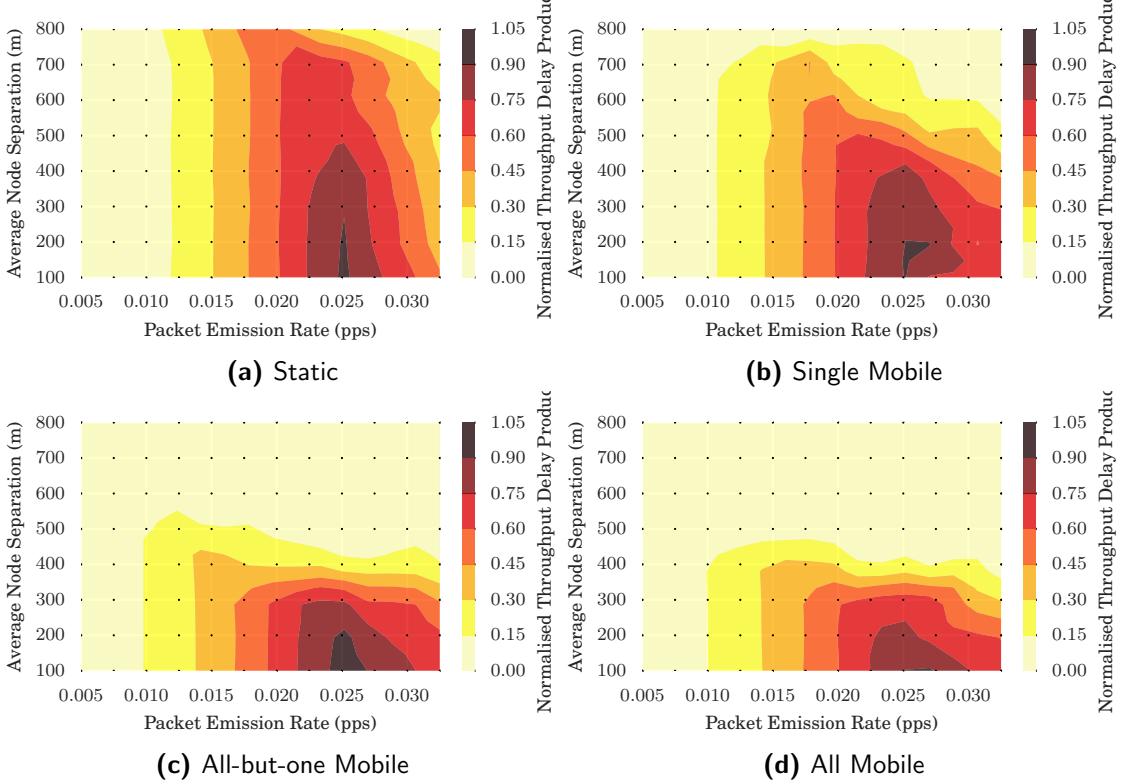
Attempt to Formalise the relationship between separation, offered load, throughput and delay

## 4.5 Conclusions

An appropriate safe operating zone for marine communications has been established by investigating the impact of variations of the communications rate and physical distribution across the mobility scenarios.

These findings can be summarised as that when the separation is increased, the emission rate at which the network becomes saturated decreases, reducing overall throughput. This throughput degradation is tightly coupled with the mobility, as increasing mobility leads to increasing delays as routes are constantly broken, re-advertised and re-established. For instance, where all nodes are static, significant drops in throughput are not seen until node separation approaches 800m, nearly double the initial estimate. However, when all nodes are randomly walking the saturation point collapses from 0.025pps at 300m to 0.015pps at 400m.

Double Check These Numbers Before Release



**Figure 4.12:** Normalised Throughput-Delay Product for all mobilities under varying separation and emission rate

These results indicate that a good area to continue operating in for a range of node separations is at  $0.015\text{pps}$ , and that a reasonable position scaling is from  $100\text{m}$  to  $300\text{m}$ , beyond which communication becomes increasingly unstable, especially in terms of end-to-end delay. These results are similar to work performed in [66], and are expected in such a sparse, noisy, and contentious environment.

expand this section to include discussion and results of single mobility models

The results from ?? and ?? show that the single-node mobility models don't capture the reality of the network

this is a place holder for actual information

The reason for this is that in other mobility combinations, the node targeted for misbehaviour ( $n_1$ ) will already be behaving differently compared to the rest of the network regardless of the misbehaviour.



## Chapter 5

# Use of Physical Behaviours for Trust Assessment

### 5.1 Introduction

This chapter proposes a new approach to trust in resource-constrained networks of autonomous systems based on their physical behaviour, using the motion of nodes within a team to detect and potentially identify malicious or failing operation within a cohort. This is accomplished by looking specifically at operations within the three dimensions of the underwater space. A series of composite metrics based on physical movement are presented and applied to the detection and discrimination of sample physical misbehaviours. This approach opens the possibility of bringing information about both the physical and communications behaviours of autonomous MANETs together to strengthen and expand the application of future Trust Management Frameworks in sparse and/or resource constrained environments.

Early attempts to secure and protect the integrity of Mobile Ad-hoc Networks have relied on various forms of strong-cryptography to protect information being transferred from tampering or malicious inspection. While such approaches protect the integrity of individual pieces of data, the increased computation, and storage requirements of modern, strong, decentralised cryptographic systems presents a clear avenue for Denial of Service (DoS) attacks on MANETs. This threat is particularly relevant in resource-constrained networks, where one or more aspects of the environment are limited, be it available power, mobility, data storage, onboard processing, bandwidth, and channel resources such as capacity and delay. In such networks, where there is a requirement of security and/or integrity monitoring, strong-cryptographic methods present an entirely new opportunity to potential attackers.

This should be moved back

One solution to the trade-off between DoS-protection, and security is the assessment of “trustworthiness” of nodes within a local network. “Trust” in this case is an assessment of capability of a node based on previously observed behaviour. Using this

Trust to make simple routing decisions is significantly simpler and faster than strong-cryptographic methods, particularly in multi-hop networks or resource constrained networks[69]. With Trust being reliant on the near-real-time awareness of some behaviour, and cryptography on the pre-establishment of some entropy store and the repeated reinforcement of that numerical security, they represent two very different approaches to system integrity with very different costs/benefits and in practice, some elements of both methodologies will be used in different contexts and applications.

However, these approaches to operational security have been totally focused on the establishment of trust/security in the communications domain, and ignore other potential threats to the network exploited through physical movement. This threat is particularly evident in collaborative autonomous systems where nodes are tasked to accomplish some survey / exploration / observation objective in a distributed fashion, where individual nodes make decisions based on the actions of their “team”. This collaboration opens the opportunity for a physically-misbehaving actor to selfishly conserve it’s own resources, or maliciously “drain” a given target node. Current security / trust systems applied to MANETs are not concerned with the threat of such physical misbehaviours.

In Section 5.2, we review the current use cases, deployments and mobility patterns of collaborative AUV operations, and the state-of-the-art in underwater localisation techniques. In Section 5.3, we discuss the use of Trust Management Frameworks and their relevance and applicability to marine operations.

Probably do away with this, repeated in a few other places

In Section 5.4, we propose a collection of metrics to characterise the physical behaviours of node, and establish a set of physical “misbehaviours” to test these against. In Section 5.5, we design a series of simulations, and tests to assess the detection and identification capabilities of three potential physical metrics for trust assessment.

## 5.2 AUV Mobility and Localisation

The use and applications of Autonomous Underwater Vehicles (AUVs) has undergone a great expansion in recent years; current applications and considerations are summarised below.

### 5.2.1 AUV operations and deployments

#### Hydrographic Survey

The use of AUVs in the place of manned-surface platforms or tethered undersea platforms enables greatly increased spatial and temporal sampling. Importantly, the separation of AUVs from the noisy sea surface enables much more efficient survey operations. This is particularly important when comparing to classical tow-line based measurements; where the mobility of the AUVs enables for much tighter-turning survey patterns or operation in inaccessible or hard-to-reach locations such as polar survey[? ].

Another significant factor is cost; the daily cost of operating a manned vessel can be considerably higher than the costs of deploying, operating and recovering one or more AUVs with equivalent capabilities[2]. Additionally, the use of low-power “glider” AUVs has lowered the barrier to entry for extended mission types, such as persistent environmental survey, or open-ocean operations. Depth-hardened AUVs have also opened up the deepest parts of the oceans to exploration, with the onboard autonomy, imagery and Simultaneous Location and Mapping (SLAM) techniques allowing deep-dwelling survey AUVs to react to bottom-surface features without the need for a tight craft-to-surface control loop. The natural extension of these kind of applications is the use of AUVs on ice-covered planets such as Europa, where three-dimensional, autonomous navigation without an on-the-loop controller is vital for mission resource efficiency and success.

### **Hull and Infrastructure Inspection**

Ongoing concerns regarding the security, safety and legality of international shipping has driven the application of AUVs to the area of near-surface hull and infrastructure inspections, looking for damage as well as devices such as limpet mines and other contraband. This use case puts a range of unique pressures on the AUV system; requiring highly accurate three-dimensional localisation and path-planning to clearly image the contours of a hull[2]. Similarly, with the increasing use and criticality of intercontinental undersea optical fibre connections, using AUVs for both the laying of and inspection of these cables is an exciting area of work[70][71].

### **Marine Petrochemical**

Oil and Gas industry requirements for high quality, low altitude bathymetry of seabed structures for infrastructure development (pipelines/drill platforms etc.) as well as monitoring of those structures over time (inspection etc.) is another significant application area, and a major driver of research investment. As in Hydrography, the mobility of AUVs is the biggest single advantage over classical platforms[72].

### **Military**

Mine-Countermeasure Operations benefit greatly from, and significantly drive, AUV development; the ability to rapidly explore and covertly survey a potentially dangerous area without risking a human operator is a major benefit. This benefit applies to protection as well as incursion; the ability to have persistent survey of a valuable area such as a forward-operating harbour is increasingly essential, and as AUV technology, autonomy and security practices develop, this use is increasing. This Port Protection capability is particularly complex; teams of AUVs are expected to repeatedly survey an area and remain densely-connected enough to maintain end-to-end communications with all other nodes, in the face of an environment that is possibly not well surveyed

initially, and includes dynamically moving obstacles (i.e. ships). In Sec. 5.5, we use this Port Protection scenario as a baseline for our simplified simulation context.

Look at redoing this with other mobilities (particularly distributed lawnmower)

### 5.2.2 Localisation Technologies

Given the subsurface nature of most AUV operations, terrestrial localisation techniques such as GPS are unavailable (below  $\approx 20\text{cm}$  depth). However, a range of alternative techniques are used to maintain spacial awareness to a high degree of accuracy in the underwater environment.

#### Long baseline (LBL)

Long-baseline localisation systems use a series of static surface/cable networked acoustic transponders to provide coordinated beacons and (usually) GPS-backed relative location information to local subsurface users. Such systems can be accurate to less than  $0.1\text{m}$  or better in ideal deployments and are regularly used in controlled autonomous survey environments such as harbour patrol operations where the deployment area is bounded. However, the initial setup and deployment required in advance of any AUV operation makes LBL difficult to utilise in unbounded or contended areas. LBL systems can also be deployed on mobile surface platforms in the area (ships or buoys for example), but these applications put significant computational pressure on the end-point AUV and have greatly reduced accuracy compared to ideal deployments[73].

#### Doppler Velocity Log (DVL)

Doppler Velocity Logging involves the emission of directed acoustic “pings” that reflect off sea bed/surface interfaces that, when received back on the craft with multi-beam phased array acoustic transducers can measure both the absolute depth/altitude (z-axis) of the craft and through directional Doppler shifting, the relative (xy-translative) motion of the craft since the ping. While classical DVL was highly sensitive to shifting currents in the water column, advances in the development of Acoustic Doppler Current Profiling has turned that situation on its head, enabling the compensation-for and measurement-of water currents down to the sub-meter level[74].

#### Inertial Navigation Systems (INS)

Inertial navigation systems use gyroscopic procession to observe the relative acceleration of a mobile platform. This reference-relative monitoring is particularly useful in the underwater environment, as it detects the motion of AUVs as they are carried by the water itself. Bias Drift is a significant problem for INS systems operating over longer (hundreds of metres) distances, as they usually have some minimal amount of directional bias, that incurs a cumulative effect over time without assistance. Several sensor

synthesis processes have been demonstrated which combine information from INS along with DVL data to improve localisation into the sub-decimeter level[75][76].

### Simultaneous Location and Mapping (SLAM)

Simultaneous Location and Mapping is the process of iteratively developing a feature-based model of an environment, and to use the relative movement within that modeled environment to obtain estimates of absolute positioning. SLAM has been most well developed in the contexts of either visual-based inspection using cameras, or LIDAR-style distance triangulation, however the same principles have been successfully applied using marine sonar readings, providing sub-meter accuracy, real-time, feature-relative localisation information that is (for the most part) environmentally agnostic[77].

In summary, current technology reliably enables AUVs to localise to a sub-metre accuracy in most areas of application.

## 5.3 Trust Management Frameworks

Trust Management Frameworks (TMFs) provide information to assist the estimation of future states and actions of nodes operating as teams, groups or networks. This information is used to optimize the performance of a team against malicious, selfish, or defective misbehaviour by one or more nodes. Previous research has established the advantages of implementing communications-based TMFs in terrestrial, 802.11 based MANETs, particularly in terms of preventing selfish operation in collaborative systems [11], and maintaining throughput in the presence of malicious actors [12]. These observations then inform future decisions of individual nodes, for example, route selection [51].

Recent work has demonstrated the use of a number of metrics to form a “vector” of trust. The Multi-parameter Trust Framework for MANETs (MTFM) [78], uses a range of communications metrics beyond packet delivery/loss rate (PLR) to assess trust. This vectorized trust also allows a system to detect and identify the tactics being used to undermine or subvert trust. This method has been previously applied to the marine space, comparing against a selection of existing communications TMFs [14] showing that MTFM is more effective at detecting misbehaviours in sparse communications environments.

## 5.4 Physical Behaviours for Trust

### 5.4.1 Physical Metrics

Three physical metrics are used to encompass the relative distributions and activities of nodes within the network; Inter-node Distance Deviation (INDD), Inter-node Heading Deviation (INHD), and Node Speed. Conceptually, INDD is a measure of the average spacing of an observed node with respect to its neighbours. INHD is a similar approach

with respect to node orientation. As such, these metrics completely encapsulate and abstract the physical behaviour of any node, potentially performing any misbehaviour. Given that local nodes within the team are aware of the reported positions and velocities of their neighbours, it is believed that this is a reasonable initial set of metrics to establish the usefulness of physical metrics of trust assessment.

Additional metric constructions may be more suitable for certain contexts, platforms or operations, however these were selected in collaboration with UK DSTL and NATO CMRE as suitable, generic, assessments, viable on most current platforms in most current deployment schemes.

$$INDD_{i,j} = \frac{|P_j - \sum_x \frac{P_x}{N}|}{\frac{1}{N} \sum_x \sum_y |P_x - P_y| (\forall x \neq y)} \quad (5.1)$$

$$INHD_{i,j} = \hat{v}|v = V_j - \sum_x \frac{V_x}{N} \quad (5.2)$$

$$V_{i,j} = |V_j| \quad (5.3)$$

Where  $i$  and  $j$  are indices denoting the current observer node and the current observed node respectively;  $x$  is a summation index representing other nodes in the observers region of concern;  $P_j$  is the  $[x, y, z]$  absolute position of the observed node (relative to some coordinated origin point agreed upon at launch) and  $V_j$  is the  $[x, y, z]$  velocity of the observed node.

Thus, the metric vector used for the physical-trust assessment from one observer node to a given target node is;

$$X_{i,j} = \{INDD_{i,j}, INHD_{i,j}, , V_{i,j}\} \quad (5.4)$$

At each time-step, each node will have a separate  $X$  assessment vector for each node it has observed in that time. Ergo the fleet or team as a whole will have  $N \times N$  assessment vectors at each timestep.

#### 5.4.2 Physical Misbehaviours

Misbehaviours in the communications space is heavily investigated area in MANETs [? ][? ][? ][? ], but attacks and misbehaviours in the physical space are far less explored. Both in terrestrial and underwater contexts, as MANET applications expand and become increasingly *de rigueur*, the impacts of physical or operational misbehaviour become increasingly relevant. As in the communications space, the primary drivers of any “misbehaviour” come under two general categories; selfish operation or malicious subterfuge. Autonomous MANETs in general rely (or are at least, most effective) when all nodes operate fairly, be that in terms of their bandwidth sharing, energy usage, routing optimality or other factors. Physically, if a node is being “selfish”, it may preferentially move to the edge of a network to minimise its dynamic work allocation,

or depending on it's intent, may insert itself into the centre of a network to maximise it's ability to capture, monitor, and manipulate traffic going across the network. In the context of a secure operation (or one that's assumed to be secure), the opportunity for capturing a legitimate node and replacing it with a modified clone. Assuming a highly capable outside actor and a multi-channel communications opportunity, there is also the possibility of a node appearing to “play along” with the crowd that occasionally breaks rank to route internal transmissions to a outside agent. In the underwater context this may mean an AUV following the rest of a team along a survey path and occasionally “breaking surface” to communicate to a malicious controller. Alternatively, if an inserted node is not totally aware of a given mission parameter, such as a particular survey or waypointing path, it may simply follow along, hoping not to be noticed.

In all these cases, such behaviour involves some element of behaving differently from the rest of the team, however, there are other cases where such individual “deviance” is observed; where a node is in some kind of mechanical “failure state”. In the underwater context, this could be damage to the drive-train or navigation systems, causing it to lag behind or consistently drift off course. An ideal physical trust management system would be able to differentiate between both “malicious” behaviours and “failing” behaviours.

To investigate this hypothesis, we create two “bad” behaviours; one “malicious”, where a cloned node is unaware of the missions’ survey parameters and attempts to “hide” among the fleet, and a “failing” node, with an impaired drive train, increasing the drag force on the nodes movement. These two behaviours are designated *Shadow* and *SlowCoach* respectively.

## 5.5 Simulation and Validation

### 5.5.1 Simulation Background

Simulations were conducted using a Python based simulation framework, SimPy [65], with a network stack built upon AUVNetSim [66], with transmission parameters taken from and validated against [57] and [67]. For the purposes of this chapter, this network is used for the dissemination of node location information, assuming suitable compression of internally assumed location data compressed into one 4096 bit acoustic data frame, with the network overall emitting approximately 10 frames a minute. Node kinematics are modeled on REMUS 100 AUVs, based on limits and core characteristics given in [79], [?] and [? ].<sup>1</sup>

These limits are given in Table 5.1

### 5.5.2 Node Control Modelling

In our investigation, we use the example of a Port Protection scenario, where a team of six AUVs are tasked with surveying a simplified harbour; in this case a 1kmx1kmx100m

---

<sup>1</sup>While the hydrodynamics of the control surfaces of the AUVs are not modeled in this case, axial drag is modeled as a resistive inertial force on the craft.

**Table 5.1:** REMUS 100 Mobility Constraints as applied in simulation

| Parameter                        | Unit             | Value |
|----------------------------------|------------------|-------|
| Length                           | $m$              | 5.5   |
| Diameter                         | $m$              | 0.5   |
| Mass                             | $kg$             | 37    |
| Max Speed                        | $ms^{-1}$        | 2.5   |
| Cruising Speed                   | $ms^{-1}$        | 1.5   |
| Max X-axis Turn                  | $^{\circ}s^{-1}$ | 4.5   |
| Max Y-axis Turn                  | $^{\circ}s^{-1}$ | 4.5   |
| Max Z-axis Turn                  | $^{\circ}s^{-1}$ | 4.5   |
| Axial Drag Coefficient ( $c_d$ ) | NA               | 3     |
| Cross Section Area               | $m^2$            | 0.13  |

cuboid volume. This is accomplished through a distributed way point system where by the team overall must “check” several points around the exterior and interior of this volume in reasonable time.

This consists of three heuristic rules; Cohesion, Repulsion and Alignment.

$$F_{j,C} = F_+ \left( p_j, \frac{1}{N} \sum_{\forall i \neq j}^N p_i, d_{max} \right) \quad (5.5)$$

$$F_{j,R} = \sum_{\forall i \neq j}^N F_- (p_j, p_i, d_{max}) \mid d_{max} > \|p_i - p_j\| \quad (5.6)$$

$$F_{j,A} = \frac{1}{N} \cdot \left( \sum_{\forall i \neq j}^N \hat{v}_i \right) \quad (5.7)$$

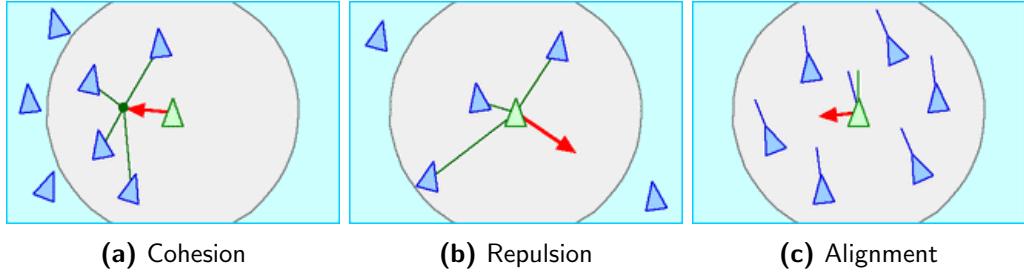
Where  $F$ 's are force-vectors applied to the internal guidance of the AUV,  $F_{j,C}$  representing Cohesion,  $F_{j,R}$  representing Repulsion, and  $F_{j,A}$  as Alignment:  $F_+$  is a scaled vector attraction function, and  $F_-$  is an equivalent repulsion function

$$F_+(p^a, p^i) = (\widehat{p^a - p^i}) \times \frac{|p^a - p^i|}{d} \quad (5.8)$$

$$F_-(p^r, p^i) = (\widehat{p^i - p^r}) \times \frac{|p^r - p^i|}{d} \quad (5.9)$$

### 5.5.3 Standards of Accuracy

The key question of this chapter is to assess the advantages and disadvantages of utilising trust from the physical domain.



**Figure 5.1:** Visual representation of the basic Boidean collision avoidance rules used

It is important to clarify what is meant by “effective” in this case; the “effectiveness” of any trust assessment framework is taken as consisting of several parts, the *accuracy* of detection and identification of a particular misbehaviour, the *complexity* of such analysis, including any specific training required, and the *differentiability* of behaviours using given metrics.

In this case we are particularly interested in the accuracy of detection and identification of malicious / failing behaviours, and as such are looking at three key characteristics of accuracy; true detection accuracy (what percentage of “bad” behaviours are detected at all); false positive rates (what percentage of “control” behaviours are detected as being “bad”); and misidentification rates (how many instances of one bad behaviour are mischaracterised as the other and vice versa).

As such we have three primary questions to answer to establish if these metrics are useful: How accurate are these metrics in being able to easily differentiate between Normal and Abnormal behaviours in terms of True-Positive and False-Positive rates? What differentiation of response, if any, is there between the stated abnormal behaviours? Can a simple classification be built to characterise these differentiations of response, and what is it’s True-Positive/False-Positive accuracy?

#### 5.5.4 Analysis

Having established the metrics under investigation, 64 simulation runs are executed for each scenario (i.e. one node “Maliciously” following the fleet with no mission information (Shadow), one “Failing” node with simulated drive train issues (Shadow), and one baseline control scenario where all nodes are behaving appropriately (Control). Each of these simulated missions last for an hour, matching realistic deployment times based on current MOD/NATO operations[80].

#### Metric Cleaning

In order to assess the viability of using the previously discussed metrics, the raw motion paths recorded by the simulation are fed into an analysis pipeline aimed at abstracting the instantaneous observed values into derived deviations from “normal” behaviour in the team.

$$d_{i,j}^{m,t} = x_{i,j}^{m,t} - \frac{\sum_k x_{i,k}^{m,t}}{|M|} \quad (5.10)$$

$$\alpha_{i,j}^{m,t} = \left| \frac{d_{i,j}^{m,t}}{\sigma d_{i,j}^{m,t}} \right| \quad (5.11)$$

Where  $i$  and  $j$  are indices denoting the current observer node and the current observed node respectively;  $x$  is a summation index representing other nodes in the observers region of concern;  $X$  is the vector of metrics from 5.4;  $d$  is an intermediate value of the distance of a given observation from the mean, and  $\alpha$  is a resulting normalised response value in terms of it's deviation from the mean.

### Behaviour Detection and Classification

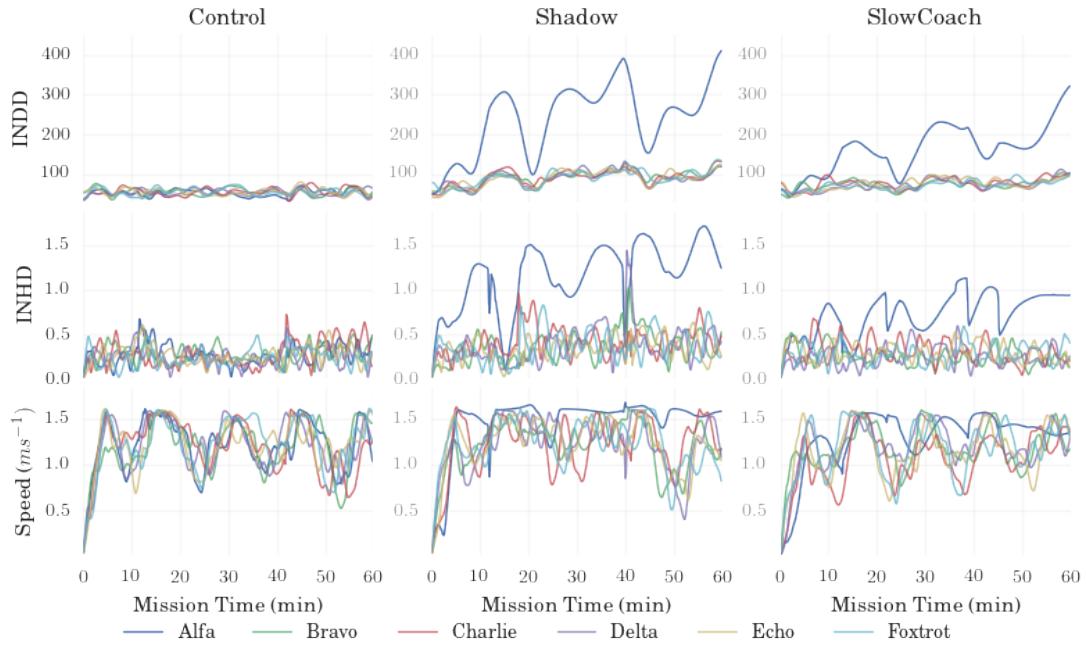
A simple misbehaviour detection is to apply Dixon's Q-test [81] to the resultant  $\sum \alpha$  values for each node for each metric for each run establishing if a “misbehaving node” exists in a given run, and if so, attempt to identify that misbehaving node. For our initial investigation we will use a Confidence Interval of 95%.

Our initial hypothesis is that by using observations of the previously stated physical metrics, that we will be able to detect and identify misbehaviours. Within that context, this Confidence Interval indicates that we would expect only a 5% chance that any run or node identified using the Q-test to *not* be a misbehaving run/node. Further, due to the range of metrics available, by applying the Q-test on a per-metric basis, we can use the “votes” of each metric as a simplified consensus classifier. This classifier may allow us to characterise some aspect of a given misbehaviour in terms of metrics it heavily impacts, and those that are less affected, finding some differentiating-limit between certain behaviours using certain metrics.

$$C_i^m = \Sigma_t \sigma_i^m * \frac{N - 1}{\sum_{x \neq i} \Sigma_t \sigma_x^m} \quad (5.12)$$

### Operational Performance Metrics

While not the focus of this paper, we are also concerned with the impact of these misbehaviours on the mission efficiency of the team overall. We monitor this in three main measurements; the “speed” of the fleet in terms of how many of it's port-protection way points it successfully approaches and passes, the total energy used for communications, and the average end-to-end delay in the acoustic network. We would expect that any misbehaviour in positioning will incur some loss of efficiency, whether it is the fleet being slowed down by a straggler attempting to catch up or of a node moving in an unexpected fashion dragging the team temporarily off course. Given that in acoustic communications, transmission is energetically expensive while reception is not, and while physical misbehaviours will not impact the amount of offered load on the network,



**Figure 5.2:** Observed Metric Values for one simulation of each behaviour ( $x_{i,j}^{m,t}$  from (5.10))

collisions induced by uneven distribution of nodes should have a small but measurable effect on energy used for packet reception.

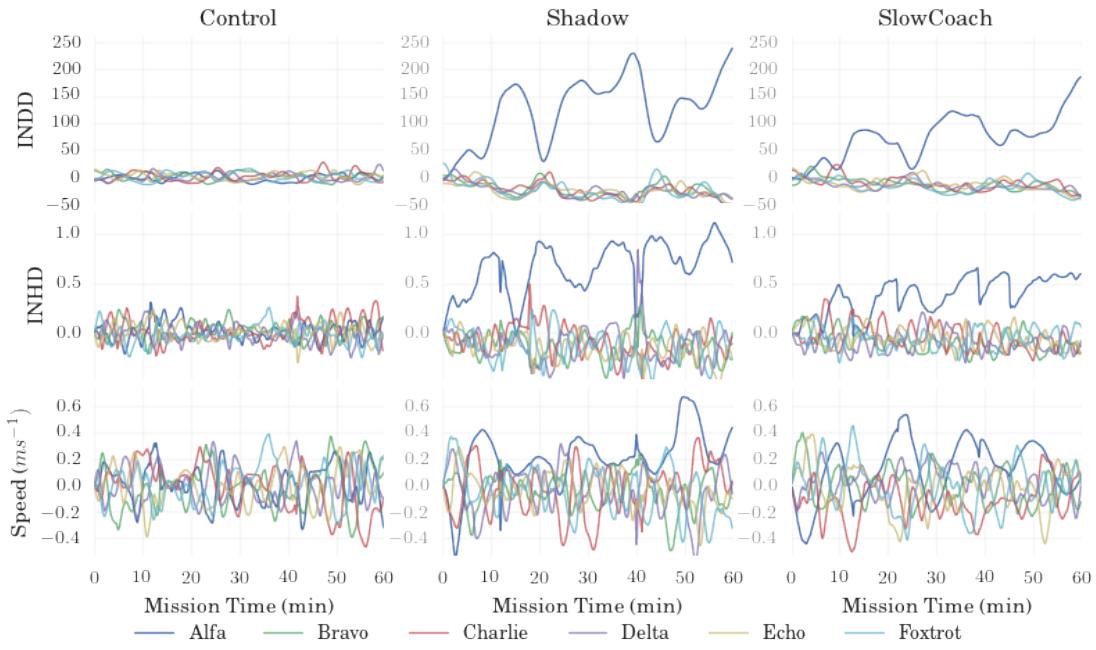
## 5.6 Results and Discussion

Fig. 5.2 shows the raw metric values (vertically) from one run of each behaviour (horizontally), starting with the Control case, where all node are behaving properly. The misbehaving node in the remaining cases. It clear that using the (unitless) INDD and INHD metrics, Alfa is the outlier and other, fairly behaving, nodes are all consistent in their metric values. This outlier-response is not nearly as clear in the Speed metric case (bottom row of Fig. 5.2).

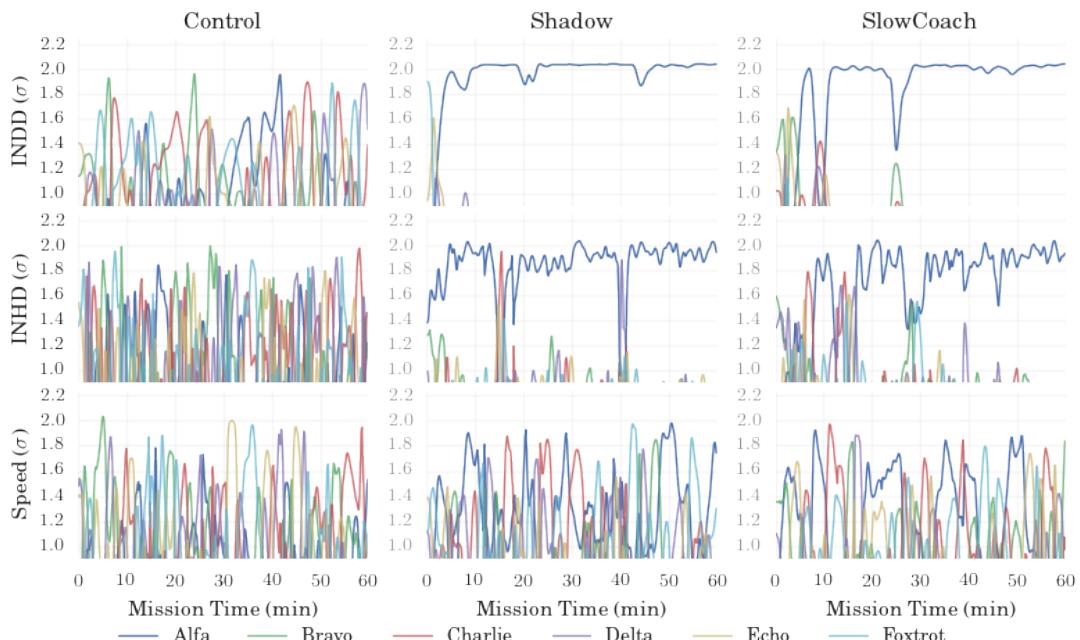
From a behaviour-perspective, it appears that the Shadow behaviour is creating the largest, most obvious deviations.

In Fig. 5.4 the metric values are normalised as per (5.11). This has highlighted the outlying-characteristic of INHD and INDD; largely eliminating the other nodes-responses. In the Speed response of Fig. 5.4, the Speed metric is not obviously highlighting any significant misbehaviours in that metric.

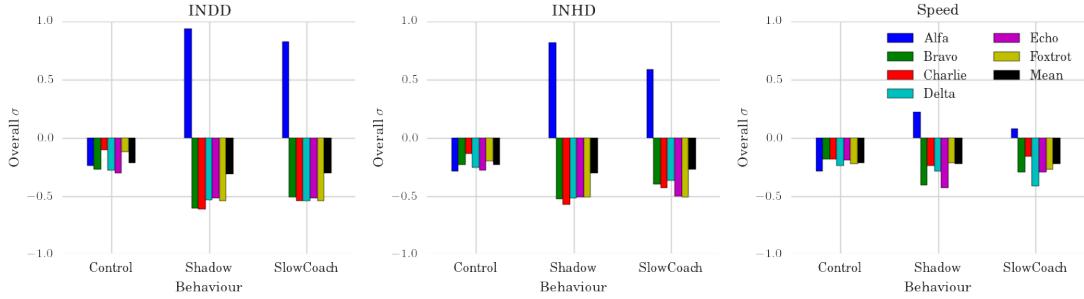
From Fig. 5.5, it appears that Speed is being significantly affected by the differing behaviours, but much less so than INHD/INDD.



**Figure 5.3:** Unnecessary but included for draft discussion Observed Metric Values for one simulation of each behaviour ( $d_{i,j}^{m,t}$  from Fig. ??)



**Figure 5.4:** Normalised Deviance values from one simulation of each behaviour ( $\alpha_{i,j}^{m,t}$  from (5.11))



**Figure 5.5:** Per-Node-Per-Run deviance for each metric, normalised in time ( $\sum \alpha / T$ )

**Table 5.2:** Overall Q-Test Outlier Correct Detection Accuracy

| Behaviour | Mean  | Std   |
|-----------|-------|-------|
| Control   | 0.927 | 0.261 |
| Shadow    | 0.979 | 0.144 |
| SlowCoach | 0.792 | 0.408 |

### 5.6.1 Detection of Misbehaviours

We have demonstrated by graphical result that from our initial metrics, INHD and INDD do appear to accurately and obviously identify the malicious node in the case that there is one. Using the deviance normalisation presented in (5.11), clear, almost contiguous areas under the Alfa-values are observed in Fig. 5.4 in the Shadow and SlowCoach misbehaviours. Further, from Fig. 5.5, it is shown that while it is nowhere near as “clear” as the deviance in INHD and INDD, that the Speed metric is still registering a statistically significant deviation in both misbehaviours, and that the difference between the deviances in Speed may indicate a way to analytically differentiate between the two misbehaviours.

To investigate how this would relate to the ability to blindly detect misbehaviours, the Q-test is applied to  $\Sigma\alpha$  results as used in Fig. 5.5, to attempt to correctly establish:

1. if a node is misbehaving and
2. which node is misbehaving

As such, the “correctness” rule for assessing this strategy is that, in misbehaving cases, the Q-tests should return Alfa (otherwise a “Fail” is recorded), and in the Control case, the Q-test should assert that there are no obvious outliers, (otherwise a “Fail” is recorded again). In Table 5.2, the Null case (Control behaviour) is correctly identified 92% of the time. The “malicious”, Shadow misbehaviour is detected and identified 98% of the time, and the “failing”, SlowCoach misbehaviour is identified just 79% of the time. These values match our intuition from Figs. 5.2 and 5.4.

**Table 5.3:** Per-Metric Q-Test Outlier Detection Accuracy

|      | Behaviour | INDD  | INHD  | Speed |
|------|-----------|-------|-------|-------|
| Mean | Control   | 0.875 | 0.938 | 0.969 |
|      | Shadow    | 1.000 | 1.000 | 0.938 |
|      | SlowCoach | 1.000 | 1.000 | 0.375 |
| Std  | Control   | 0.336 | 0.246 | 0.177 |
|      | Shadow    | 0.000 | 0.000 | 0.246 |
|      | SlowCoach | 0.000 | 0.000 | 0.492 |

We can investigate this further by looking at the “correctness” of the assessments of each metric individually (Table 5.3). In both misbehaviours, INHD and INDD correctly identify Alfa as the misbehaver 100% of the time. However, they misidentify a potential misbehaviour in the Control case 13% and 7% of the time respectively. Meanwhile, Speed correctly identified the Control case 97% of the time, and the Shadow case 94% of the time, but missed the SlowCoach behaviour 63% of the time. This result is surprising on the face of it, as SlowCoach is a misbehaviour that is exclusively about individual node speed and conceptually should have had a much larger impact on the simple Speed metric. However, the collaborative nature of the collision avoidance system, and the existing limits on node kinematics from Table 5.1 appear to be hiding this impact.

### 5.6.2 Identification of Misbehaviours

Having established the ability of INDD, INHD and Speed to all detect physical misbehaviour to a statistically significant level, and having shown that there is a demonstrable difference in response to different misbehaviours, we return to the last question from Sec. 5.5.3; can a simple classifier based on a subset of our results be constructed, and can it be blindly applied to a new set of results successfully?

From (5.12), the per-metric-per-behaviour “Confidence” in the relationship between a given metric deviance and each behaviour is established. It is hypothesised that this confidence can be used as a signature for that metric.

**Table 5.4:** Metric Confidence Responses for known behaviours (5.12)

|      | Behaviour | INDD  | INHD  | Speed |
|------|-----------|-------|-------|-------|
| Mean | Control   | 1.064 | 0.966 | 1.010 |
|      | Shadow    | 4.059 | 3.374 | 2.098 |
|      | SlowCoach | 4.246 | 3.352 | 1.491 |
| Std  | Control   | 0.262 | 0.113 | 0.132 |
|      | Shadow    | 0.398 | 0.436 | 0.206 |
|      | SlowCoach | 0.198 | 0.288 | 0.180 |

Having demonstrated that the Null case (All nodes behaving fairly) can be identified to a strong degree of accuracy, our classifier will continue to use the Q-test across all metrics for that case and concentrate on differentiating the Shadow and SlowCoach behaviours where they exist.

From Table 5.4 it is clear that INHD and INDD have similar responses to both misbehaviours, with significant standard deviations, but the response of the Speed metric is much more stable and discernible; across the range of training simulation runs. In the SlowCoach behaviour, this Speed response centres around 1.5, while the Shadow behaviour centres around 2.0, with these centres being at least one standard deviation away from each other respectively.

Our generated classifier is formalised in (5.13).

$$C \rightarrow \begin{cases} Q^{95}(X) = \emptyset, & \text{Control} \\ Q^{95}(X) \neq \emptyset \wedge \text{Speed}^X \leq 1.75, & \text{Shadow} \\ Q^{95}(X) \neq \emptyset \wedge \text{Speed}^X > 1.75, & \text{SlowCoach} \end{cases} \quad (5.13)$$

Applying this simplified classifier to a blind test set of simulations (of the same scale) gives surprisingly positive results as shown in Table 5.5, with greater than 90% identification rates for both misbehaviours. However, in the Null (Control) case we experience a false-positive rate of nearly 30%, that is to say that in the case where there is no misbehaviour, 30% of the time a node will be misidentified as misbehaving when it is not.

Given the simplicity of the applied classifier, these are strongly positive results for the use of physical metrics for behaviour discrimination; with INHD and INDD proving as strong and obvious “canaries” of misbehaviour, and Speed in this case proving a capable differentiator between conceptually close misbehaviours.

**Table 5.5:** Successful Identification rates on untrained results using (5.13)

| True Behaviour | Probability of Correct Blind Identification |
|----------------|---|
| Control        | 0.719                                       |
| Shadow         | 0.906                                       |
| SlowCoach      | 0.938                                       |

Explain the Minority Classifier

### 5.6.3 Impacts of Misbehaviour on operational performance

The anticipated “small but measurable” effects to communications performance and energy usage are indeed extremely small and within the bounds of statistical uncertainty. One observation of merit was an observed 10% increase in end-to-end delay in the case of the Shadow behaviour; this is due to the misbehaving node “overshooting” the mission

**Table 5.6:** Successful Identification rates on untrained results using (5.13), with outlier consensus checks

| True Behaviour | Probability of Correct Blind Identification |
|----------------|---|
|                |   |
| Control        | 1.000                                       |
| Shadow         | 0.906                                       |
| SlowCoach      | 0.938                                       |

way points and thus temporarily looking local connection to nodes on the opposite side of the fleet from it, causing retransmissions thus, delays. As for physical efficiency, achievement rates were identical to within 2% error on each run across all behaviours, and fleet distance varied by a similar margin. It's possible that our selected behaviours were too unambitious in our impacts, and future work will have to investigate the impact of "heavy-handed" or destructive behaviours on the operational efficiency of autonomous networks.

## 5.7 Conclusion

In this paper we have demonstrated that with current and on-the-horizon underwater localisation techniques, that in certain mobility models, that a set of relatively simple geometric abstractions (INHD, INDD, and Speed), between nodes as part of an Underwater MANET can be used as a Trust Assessment and Establishment metric.

These metrics are application-agnostic and could potentially be applied in other areas of mobile autonomy such as UAV operations and Autonomous Vehicular Networks.

We show, using a Port-Protection way point scenario built upon a Boidian collision prevention behaviour that in a simulated underwater environment, the outputs of these metrics can be used to detect and differentiate between exemplar malicious behaviour and potential failure states.

This verification further supports the assertions the authors have made previously that it is practical to extend Trust protocols such as Multi-parameter Trust Framework for MANETS (MTFM)[78] to include metrics and observations from the physical domain as well as those from the communication domain[82]. This combination of physical and "logical" information would further support the decentralised and distributed establishment of observation based Trust.

# Chapter 6

# Communications Trust Assessment in Underwater MANETs

## 6.1 Introduction

### 6.1.1 Selected Misbehaviours

We are primarily concerned with the direct trust relationship between  $n_0$  and  $n_1$ , i.e.  $n_0$ 's assessment of the trustworthiness of  $n_1$ , or  $T_{1,0}$ .

Guo et al. introduce a range of misbehaviours, including modification of the packet loss rate of routing nodes and limiting throughput on a per-link basis as well as a selection of combined misbehaviours. Given that the established links are already heavily constrained, such attacks would severely impact the general performance of the network beyond the scope of simple selfishness. These direct malicious behaviours effectively trigger saturation collapses in operating regions of the network that should be stable.

Therefore, two more subtle misbehaviours to investigate are;

1. **Malicious Power Control (MPC)**, where  $n_1$  increases its transmit and forwarding power by 20% for all nodes *except* communications from  $n_0$  in order to make  $n_0$  appear to be selfishly conserving energy to the rest of the team, while  $n_1$  itself appears to be performing very well.
2. **Selfish Target Selection (STS)**, where  $n_1$  preferentially communicates, forwards and advertises to nodes that are physically close to it in effort to reduce its own power consumption.

## 6.2 Simulation Results and Discussion

Having established a safe operating range for comparison at 300m average separation and an emission rate of 0.015pps, each of the three selected behaviours (Fair, MPC, STS) are

performed in both the static and mobile scenarios. We select a trust assessment period of 10 mins for a five hour mission to scale in comparison to relative bitrates experienced (1Mbps vs  $\approx$  15bps).

The six metrics used for grey assessment are; transmitted and received throughput and power, delay, and **PLR** as calculated by aborted and unacknowledged, transmissions. Compared to [13], this metric set lacks a data rate quantity as the network is not dynamically adjusting bandwidth. In context of **GRC** generation (C.2), the best sequence  $g$  was selected using the lowest PLR, delay, and powers, and the highest throughputs, and the worst sequence,  $b$  the inverse of these metrics, reflecting the observations made in Section 3.2.

The particular factors under discussion are the relative performance of **MTFM** against **OTMF** and Beta with respect to statistical stability across mobilities and in responsiveness to changing network behaviour. We establish a similar result set by initially tracking the resultant trust values established by **MTFM** in the pair of mobility scenarios, shown in Fig. A.2. We are also concerned with the opinions of  $n_1$  provided to  $n_0$  by other nodes, where  $[T_{1,2}, T_{1,3}]$  and  $[T_{1,4}, T_{1,5}]$  denote the sets of recommendation and indirect trust assessment respectively.

We also include aggregate assessments;  $T_{1,\text{Avg}}$ , the unweighted mean of direct trust assessments of  $n_1$  from all nodes and  $T_{1,\text{MTFM}}$ , the final **MTFM** trust assessment value based on both network topology and whitenization from (C.5).

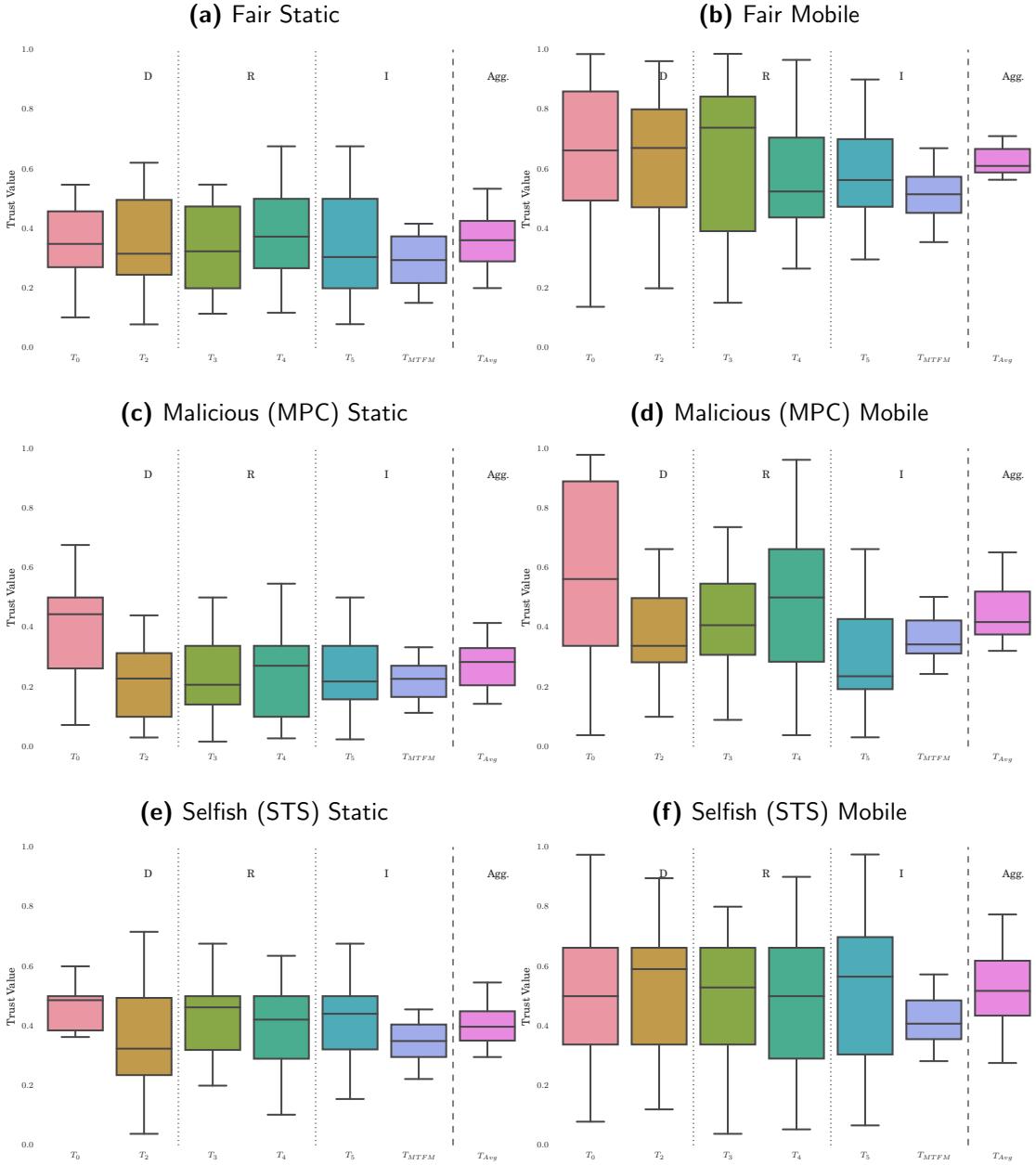
The variability in assessment is coupled to mobility; in the static case (Fig. 6.1a), the nodes exhibit relatively consistent distributions. In the full mobility case, shown in Fig. 6.1b, this subjective variability is greatly increased. As the topology is highly dynamic, delays due to re-establishing routes can be very large, perturbing the trust value. The  $T_{1,\text{MTFM}}$  displays a significantly reduced variation than those of the individual subjective observations in all cases, even when compared to the unweighted average,  $T_{1,\text{Avg}}$ . This demonstrates  $T_{\text{MTFM}}$ 's value as an aggregating trust assessment in such sparse and noisy environments. Further, in Fig. 6.1d a much higher variability in assessment is observed in  $T_0$ , correctly indicating that there is something wrong with the relationship between  $n_0$  and  $n_1$ .

### 6.2.1 Comparison between **MTFM**, **Hermes** and **OTMF**

As per [13], “fair” scenarios were also performed with no malicious behaviour, applying **OTMF** and **Hermes** assessment as well as **MTFM**, providing like-for-like comparison of assessment. For simplicity of presentation, only the fully-mobile scenario is considered, as we are concerned with the establishment of trust in mobile networks.

In the thesis, we're concerned about a lot more than just the all mobile results

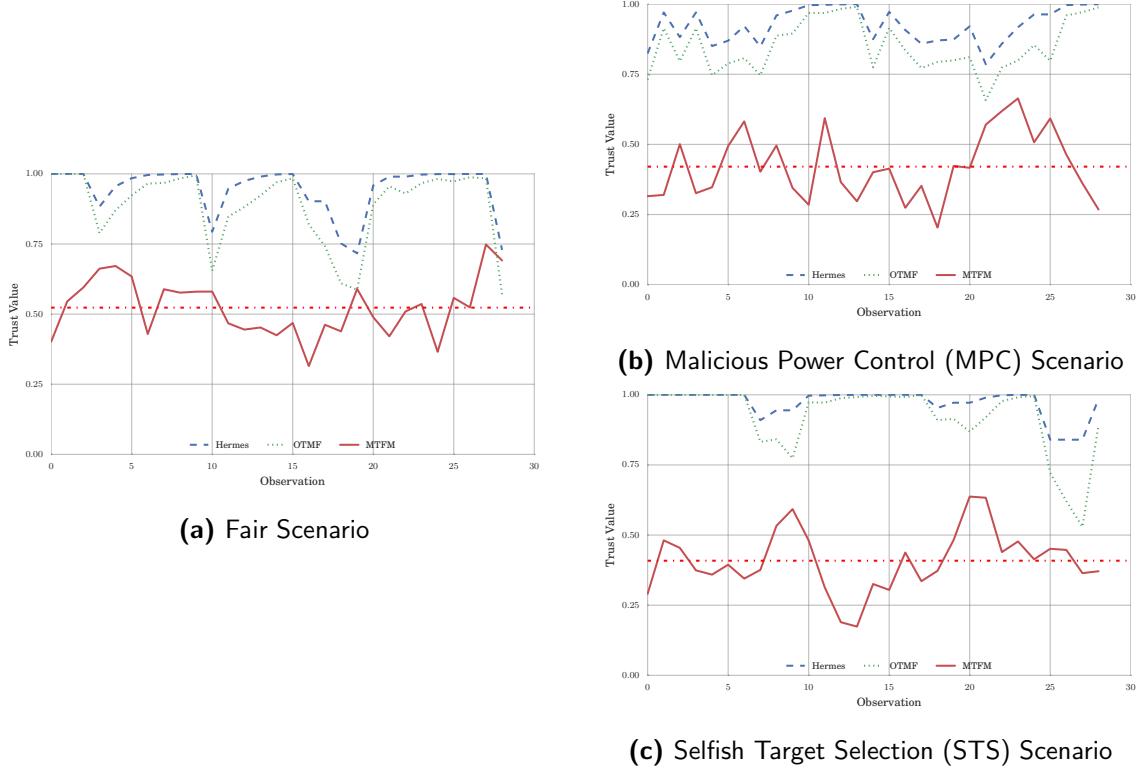
The use of Forward Beam Routing and a CSMA/CA MAC scheme from AUVNetSim [66] in our simulation mitigates a significant number of packet losses through collision avoidance and contention handling, leading to the situation that the only genuinely lost packets occur when a node moves completely out of range of any other node and time out



**Figure 6.1:** MTFM Trust assessments of  $n_1$  ( $T_{1,X}$ ), showing Direct, Recommender and Indirect relationships, as well as the Aggregate trust assessments from combining these

occurs in route discovery rather than transmission. As such, confirmed packet losses are relatively rare and in a delaying network like this, it is difficult to set a differentiating time out between packets that are in the network but queued, and packets that are actually “lost”.

The single metric TMFs used in conventional MANETs require regular and constant input to shape and adjust their evaluations, which for a network with significant and irregular delays such as this, is not practical. This renders OTMF and Hermes assessment at best uninformative and at worst misleading; consistently providing nodes a high



**Figure 6.2:**  $T_{1,0}$  for Hermes, OTMF and MTFM assessment values for fair and malicious behaviours in the fully mobile scenario (mean of MTFM also shown)

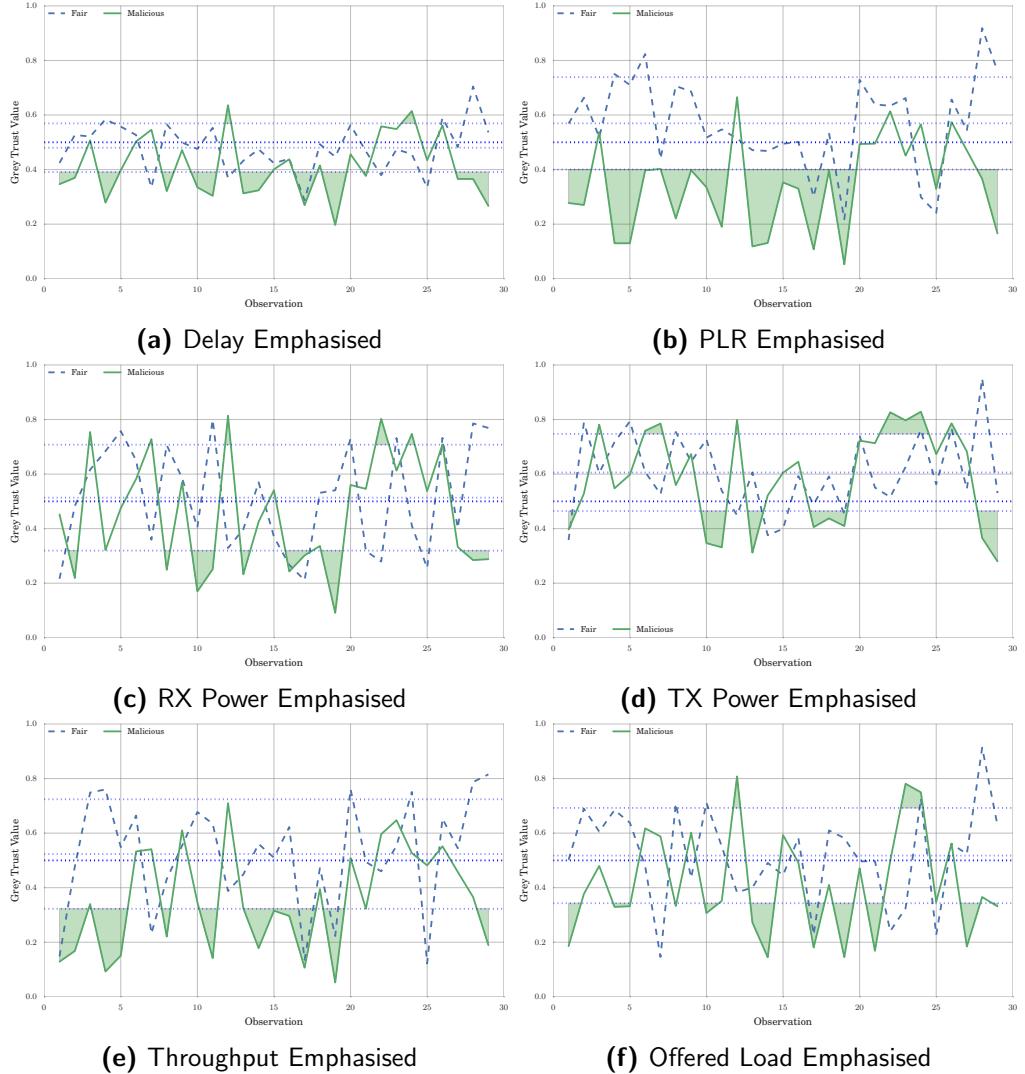
trust assessment as they have very little information to extract trust from.

Fig. 6.2 shows a comparison between the unweighted response of MTFM compared to OTMF and Hermes assessment functions on the same data for the fair, malicious and selfish behaviours respectively. It is important to note a distinction between the expectations of MTFM compared to other TMFs; MTFM is primarily concerned with the identification of differences in the behaviours of nodes in a network, and is relative rather than absolute. That is to say that under MTFM, nodes are compared against the worst current performances across metrics of other observed nodes and graded against them, rather than the absolute (objective) approach taken by many TMFs. In these cases, particularly since the methods of attack were not directly related to PLR, OTMF and Hermes have not registered significant activity in either misbehaviour when compared to the fair scenario. The difference between the MTFM trust assessments under “fair” and “malicious” behaviour is lowered by  $\approx 10\%$  in both cases, in terms of the mean values returned. At run time, similar results could be attained by an exponentially weighted moving average filter (EWMA).

On their own, neither OTMF, Hermes, or unbiased MTFM appear to be effective in detecting or identifying malicious behaviour in this environment, in fact OTMF and Hermes don’t appear to differentiate between fair and selfish scenarios at all.

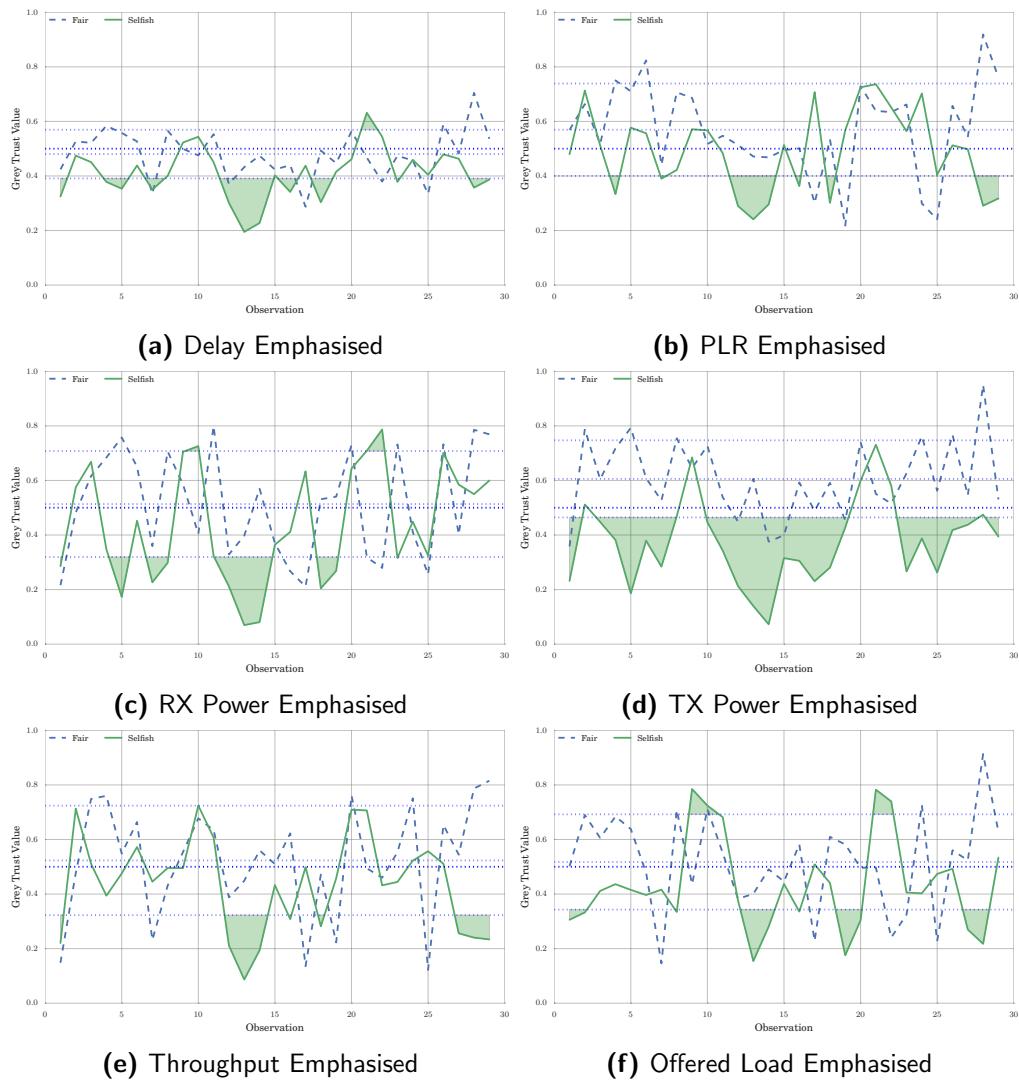
### 6.2.2 Metric Weighting

A sequence of vectors that preferentially weight each metric in Eq. (C.3) to each of the three simulation runs. For a metric weight vector  $H$ , where the metric  $m_j$  is emphasised as being twice as important as the other metrics, forming an initial weighting vector  $H' = [h_1 \dots h_M]$  such that  $h_i = 1 \forall i \neq j; h_j = 2$ . That vector  $H'$  is normalised such that  $\sum H = 1$  by  $H = \frac{H'}{\sum H'}$ . Using this process the primary aspects of an attack can be extracted and highlighted by comparing against the deviation from the “fair” result set.



**Figure 6.3:**  $T_{1,MTFM}$  in the All Mobile case for the Malicious Power Control behaviour, including dashed  $\pm\sigma$  envelope about the fair scenario

Fig. 6.3 shows that the malicious node is consistently outside the  $\pm\sigma$  (one standard deviation above and below the mean) envelope of the fair scenario it's being compared to. This is particularly true for PLR, with smaller impacts on delay, received power and offered load. This weighted delta in received throughput is minimal to insignificant compared to the width of the detection envelope, occasionally breaching the envelope for a short period.



**Figure 6.4:**  $T_{1,MTFM}$  in the All Mobile case for the Selfish Target Selection behaviour, including dashed  $\pm\sigma$  envelope about the fair scenario

In the selfish case (Fig. 6.4) a much lower weighted delta in PLR and delay is observed, with greatly increased impact on transmission power. In comparison to [13], these results are qualitatively similar, however here the differences between the fair case and the misbehaviours are less clear than in the comparable terrestrial space. Guo et al. show similar types of behaviour but report a weighted delta from  $\approx 0.4$  to  $\approx 0.9$  across the simulation period, compared to our maximum delta in  $P_{TX}$  in selfish behaviour (Fig. 6.4d) of  $\approx 0.3$  for an inconsistent interval.

### 6.2.3 Weight Significance Analysis for Behaviour Classification

For a more quantitative assessment of the viability of multi-metric trust assessment methods, taking the qualitative analysis above and apply a Random Forest regression [83] to assess the relative importance of the selected metrics on relative detectability of malicious behaviour. Random Forest accomplishes this by generating a large number of

random regression trees and prune these trees to fit incoming data. The target function for this regression was the area between the target behaviours weighted  $T_{MTFM}$  curve and the  $\pm\sigma$  envelope of the base behaviour as shaded in Figs. 6.3 and 6.4. From this training process, the relative importance of each input feature (metric) can be inferred in terms of how good it is to differentiate between the fair case and a given misbehaviour. Additionally a cross correlation analysis is performed to establish the correlations between given metric weighting emphasis and the output of the target function. Our intention is to establish the metrics that not only differentiate both misbehaviours from the fair case, but also what metrics differentiate the two misbehaviours from each other.

Applying this target regression to 729 different metric weight vector emphasis combinations reveals that each of the three combinations (i.e. comparing fair to misbehaviours, and comparing the misbehaviours) present distinct patterns of significance in three primary metrics; received throughput, transmitted power, and PLR, with delay, received power and transmitted throughput playing a lesser role. Practically this means that in order to accurately distinguish between these scenarios, these primary metrics should be higher-weighted in the generation of  $T_{1,MTFM}$  in (C.4).

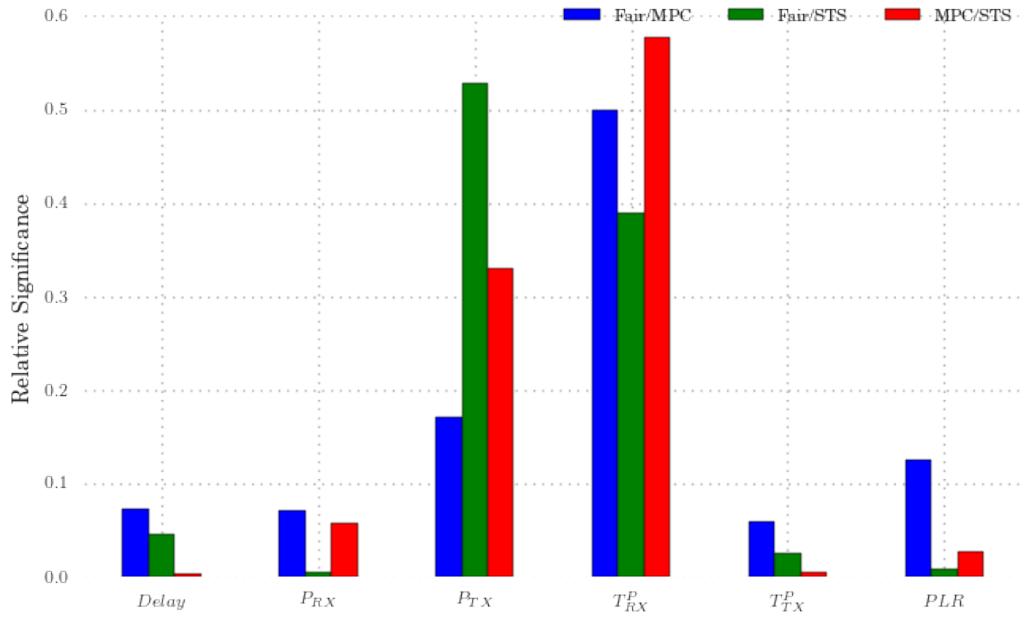
It may initially appear odd that the relative significance of the received throughput is similar between all three scenario combinations, however a correlation analysis shows that in the MPC attack; the received throughput is positively correlated with successful classification against the fair case ( $R = +0.71, p \approx 10^{-100}$ ), while the inverse is the case for the STS attack ( $R = -0.70, p \approx 10^{-100}$ ). It is expected that Transmitted power should be the defining characteristic of STS ( $R = +0.72, p < 10^{-100}$ ) as the node is acting fairly from a protocol perspective but is acting unfairly at a higher (incentive) level; it is performing fairly in terms of its communications with other nodes, however it is preferring to communicate with nodes that it can expend less energy communicating with. A summary of these correlations is shown in Table. 6.1.

Comparing Figs. 6.2, 6.3b, and 6.4b, while it is possible that in a cleaner, less sparse, and less noisy environment, OTMF would be able to detect the MPC behaviour, Fig. 6.5 shows that PLR plays almost no part at all in detecting the STS behaviour, and so OTMF would not detect the attack.

**Table 6.1:** Correlation Coefficients between metric weights and behaviour detection targets

| Correlation | Delay | $P_{RX}$ | $P_{TX}$ | $G$    | $S$    | PLR    |
|-------------|-------|----------|----------|--------|--------|--------|
| Fair / MPC  | 0.199 | 0.159    | -0.416   | 0.708  | -0.238 | -0.401 |
| Fair / STS  | 0.179 | -0.009   | 0.724    | -0.697 | -0.145 | -0.052 |
| MPC / STS   | 0.058 | -0.134   | 0.146    | -0.768 | 0.052  | 0.146  |

As such this presents the open opportunity to develop a heuristic weight search scheme to detect malicious behaviour without the comparison to the fair scenario. This



**Figure 6.5:** Random Forest Factor Analysis of Malicious (MPC, Selfish (STS) and Fair behaviours compared against each-other

would be accomplished by assessing the impact of differential metric weighting on the mean trust assessment rather than comparing co-weighted valuations across scenarios.

### 6.3 Conclusions

It has been demonstrated that existing MANET Trust Management Frameworks are not directly suitable to the sparse, noisy, and dynamic underwater medium. By comparing the operation and performance of trust establishment in MANETs in a simulated underwater environment has demonstrated that in order to have any reasonable expectation of performance, that throughput and delay responses must be characterised before implementing trust. While the MTFM value does not display any immediate difference between the two behaviours, it has been shown that by exploring the metric space by weight variation, the existence and nature of the malicious behaviour can be discovered. Another difference is that MTFM is significantly more computationally intensive than the relatively simple Hermes / OTMF algorithms. The repeated metric re-weighting required for real time behaviour detection is therefore an area that requires optimization. With significant delays (from seconds to many minutes), in a fading, refractive medium with varying propagation characteristics, the environment is not as predictable or performant as classical MANET TMF deployment environments.

It is shown that, without significant adaptation, single metric probabilistic estimation based TMFs are ineffective in such an environment. Additionally, it's clear that existing frameworks are overly optimistic about the nature and stability of the communications channel, and can overlook characteristics that are useful for assessing the

behaviour of nodes in the network. This indicates that there is a good case, particularly within constrained **MANETs** as this, for multi-vector, and even multi-domain trust assessment, where metrics about the communications network and topology would be brought together with information about the physical behaviours and operations of nodes to assess trust.

A significant additional factor of trust assessment in such a constrained environment is that there may be long periods where two edge nodes (for instance,  $n_0 \rightarrow n_5$ ) may not interact at all. This can be due to a range of factors beyond malicious behaviour, including simple random scheduling coincidence and intermediate or neighbouring nodes collectively causing long back-off or contention periods. This disconnection hinders trust assessment in two ways; assessing nodes that do not receive timely recommendations may make decisions based on very old data, and malicious nodes have a long dwelling time where they can operate under a reasonable certainty that the **TMF** will not detect it (especially if the node itself is behaving disruptively).

### 6.3.1 Future Work

One solution to this would be to move from a stepping-window of trust observations to a continuous trust log, updated on packet reception rather than waiting regular periods for packets to be analysed. Future work will investigate the improvement of weight-based detection algorithms, the stability of **GRC** under multi-node collusion, the development of real-time outlier detection, and the introduction of physical behavioural metrics into the trust assessment context.



# Chapter 7

# Multi-Domain Trust Assessment in Collaborative Marine **MANETs**

## 7.1 Introduction

In this chapter, a multi-domain trust management framework (MD-TMF) is demonstrated in collaborative marine MANETs. A methodology is demonstrated that applies Grey Sequence operations and Grey Generators to provide continuous trust assessment in a sparse, asynchronous metric space across multiple domains of trust. By utilising information from multiple domains, it is demonstrated that trust assessment can be more accurate and consistent in identifying misbehaviour than in single-domain assessment. Further, a methodology for assessing the usefulness of individual metrics in this cross-domain space is demonstrated, allowing for the elimination of redundant metrics, simplifying the runtime assessment process.

## 7.2 Initial Optimisation of Multi-Domain Trust with Pre-defined Domains

A key question in this chapter is to assess the advantages and disadvantages of utilising trust from across domains. This includes a secondary question as to how trust assessments from these domains are most effectively combined.

It is important to clarify what is meant by “effective” in this case; the “effectiveness” of any trust assessment framework is taken as consisting of several parts.

1. the *accuracy* of detection and identification of a particular misbehaviour
2. the *timeliness* of such detections
3. the *complexity* of such analysis, including any specific training required

4. the *commonality* of the results of any detections between perspectives (also termed “isomorphism” of results)

### 7.2.1 Communications Trust Metrics

The metric vector is constructed using those trust metrics that are applicable to the marine environment from [78], as the simulated marine acoustic modem stack does not operate on the same tiered data-rate approach as used in the 802.11 stack, the data rate metric was not included. Remaining metrics are; Delay, Received and Transmitted power, Throughput ( $S$ ), Offered Load ( $G$ ) and **PLR**.

Thus, the metric vector used for communications-trust assessment is;

$$X_{comms} = \{D, P_{RX}, P_{TX}, S, G, PLR\} \quad (7.1)$$

### 7.2.2 Physical Trust Metrics

Three physical metrics are selected to encompass the relative distributions and activities of nodes within the network; **Inter-Node Distance Deviation (INDD)**, **Inter-Node Heading Deviation (INHD)**, and Node Speed. These metrics encapsulate the relative distributions of position and velocity within the fleet, optimising for the detection of outlying or deviant behaviour within the fleet.

Conceptually, **INDD** is a measure of the average spacing of an observed node with respect to its neighbours. **INHD** is a similar approach with respect to node orientation.

$$INDD_{i,j} = \frac{|P_j - \sum_x \frac{P_x}{N}|}{\frac{1}{N} \sum_x \sum_y |P_x - P_y| (\forall x \neq y)} \quad (7.2)$$

$$INHD_{i,j} = \hat{v}|v = V_j - \sum_x \frac{V_x}{N} \quad (7.3)$$

$$V_{i,j} = |V_j| \quad (7.4)$$

Thus, the metric vector used for physical-trust assessment is;

$$X_{phy} = \{INDD, INHD, V\} \quad (7.5)$$

Need to actually show physical only trust measurements

### 7.2.3 Cross Domain Trust Metrics

This simplest possible combination is a vector concatenation across domain metric vectors; in this case;

$$X_{merge} = (X_{comms}|X_{phy}) = \{D, P_{RX}, P_{TX}, S, G, PLR, INDD, INHD, V\} \quad (7.6)$$

### 7.2.4 Metric Weight Analysis Scheme

From (C.3), the final trust values arrived at are dependent on metric values, the weights assigned to each metric, and the structure of the  $g$ ,  $b$  comparison vectors.

referencing the right equ in the wrong place

This permits the assessment of the significance of different metrics in the detection and identification of different behaviours. The primary aspects of a (mis)behaviour can be detected and assessed by comparing a weighted trust assessment against the deviation from a “fair” result set using the same weight, i.e. we are interested in the weight schemes that create the largest difference between fair and misbehaving cases.

For a metric weight vector  $H$ , where the metric  $m_j$  is emphasised as being twice as important as the other metrics, an initial weighting vector  $H' = [h_1 \dots h_M]$  is formed such that  $h_i = 1 \forall i \neq j; h_j = 2$ . That vector  $H'$  is then scaled such that  $\sum H = 1$  by  $H = \frac{H'}{\sum H'}$ .

The construction of the  $g$  and  $b$  vectors from Equation C.2 depends on the particular metric, e.g. Throughput ( $S$ ) on a link is assumed to be positively correlated to trustworthiness and so follows the default construction ( $g(S) \mapsto \max, b(S) \mapsto \min$ ), whereas in the case of a metric such as delay, this relationship is inverted, i.e. longer delays indicate less trustworthy activity ( $g(D) \mapsto \min, b(D) \mapsto \max$ ). This inversion relationship (i.e. those with the construction  $g(x) \mapsto \min, b(x) \mapsto \max$ ) is signified by a negative weight.

In complex environments, the relationship between metrics trustworthiness correlations is not always as obvious as the throughput / delay examples. This phenomenon was mentioned by Guo [78], but was manually configured for each metric for each behaviour and no analytical method for quantitatively establishing such relationships has been presented since.

With the nine selected metrics from across communications and physical behaviours, we can explore this metric space by varying the weights associated with each metric, and choose to emphasise across three levels; i.e. metrics can be ignored or over-emphasised. Naively this results in  $3^9 = 19683$  combinations, however as these weights are being normalised, redundant duplicates can be eliminated, e.g.  $[0, 0, 0, 0, 1, 0, 0, 0, 0] \equiv [0, 0, 0, 0, 2, 0, 0, 0, 0]$  leaving 18661 unique weights for analysis.

To assess the performance of a given weight combination (i.e. an optimisation factor), we are initially interested in the metric weight vector that consistently provides the largest deviation in the final trust value  $T$  across the cohort, i.e. producing the most clear detection of a node misbehaving in that particular fashion. This is approached as an inverse outlier filtering problem, and the range outside a  $\pm\sigma$  envelope compared to the equivalent weighting in a known “fair” behaviour is selected to assess detection (or comparing to other misbehaviours to assess discrimination). See Subsection 6.2.2. Note

that at this point we establish “signatures” of different behaviours rather than optimal detection weights.

#### Duplicating C6 Metric Weighting Section

We apply a Random Forest regression [83] to assess the relative importance of the selected metrics on relative detectability of malicious behaviour. Random Forest accomplishes this by generating a large number of random regression trees and prune these trees based on how accurate they are in correctly matching the input data. In this case that data is the deviation in trust observed ( $\Delta T$ ) between a two behaviours, i.e. maximising the ability to tell the difference between two given behaviours (i.e. “Fair” and “Malicious”). A major advantage of Random Forest in this case is that by walking the most successful regression trees, we can acquire an already normalised maximal activation weight for the particular behaviour comparison being tested.

After establishing the importance of weights in particular behaviours, a final weight is arrived at by algorithmically those few metrics that are important, rather than having to further explore the computationally expensive weight-space.

Using this approach, the results of these simulations can be explored, condensing the multi-dimensional problem (target / observer / behaviour / metric / time) down to a more manageable level for analysis.

#### 7.2.5 Significance Analysis

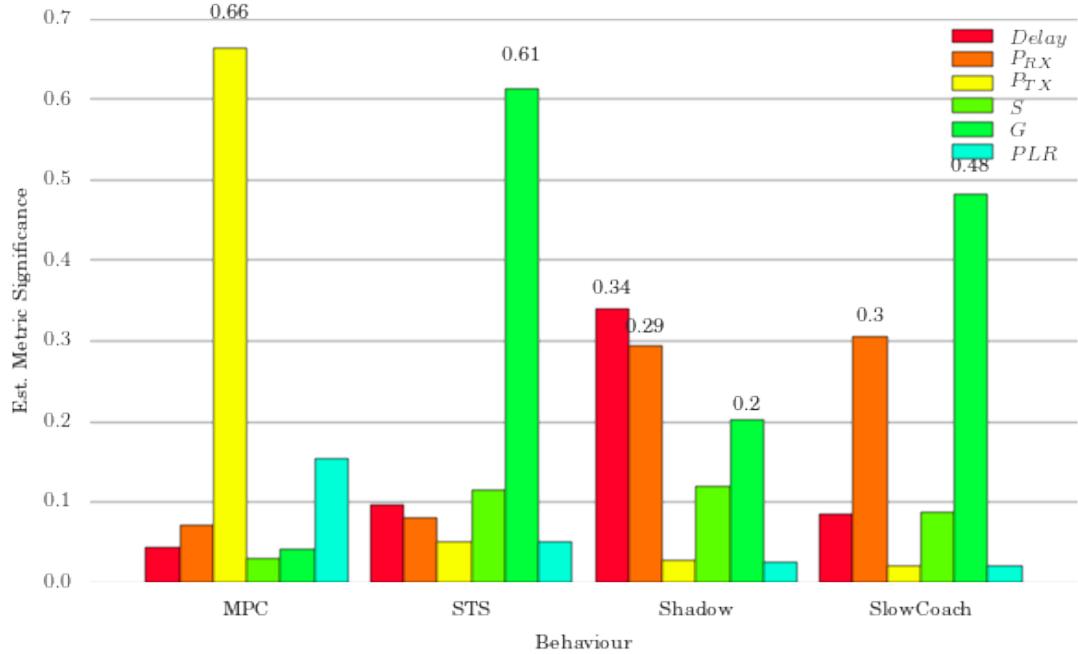
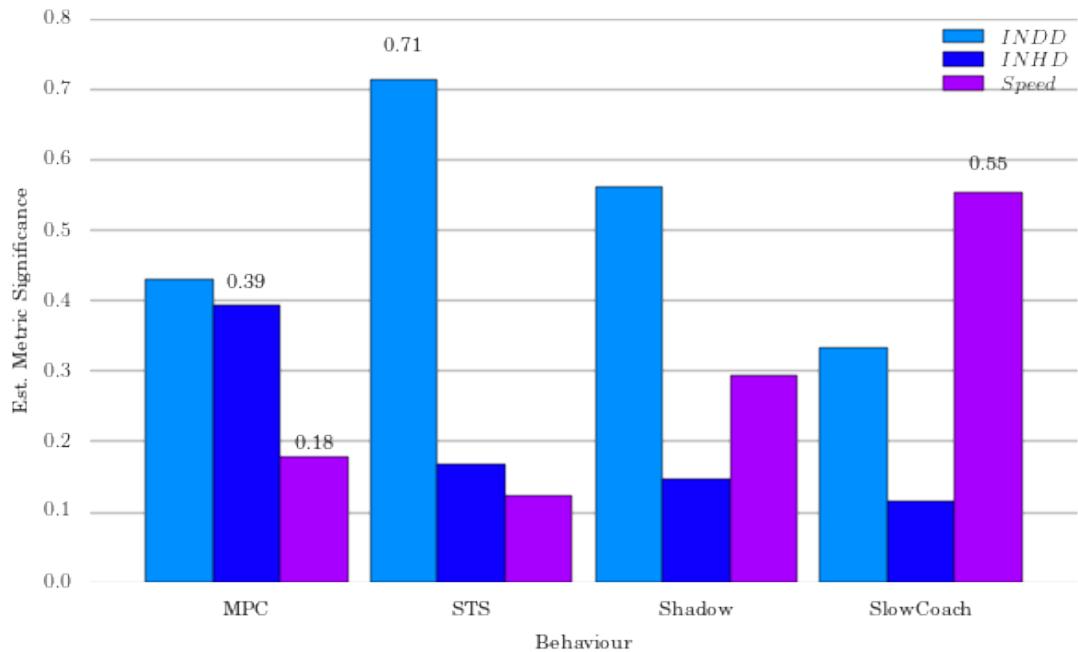
First the results of the Random Forest regression assessment are discussed; Figs 7.1 and 7.2, show the resultant feature extraction signatures for Comms-only and Physical-only metric selections respectively, and in Fig 7.3, these metric spaces are brought together and reassessed.

In both single-domain cases, there are clear “signatures” in misbehaviours that don’t directly target that domain ( $P_{RX}$  in the Physical Shadow and Slowcoach behaviours in Fig 7.1 and  $INDD$  in the Selfish Target Selection behaviour in Fig 7.2). This inter-domain activity is to be expected in MANETs in general, where the physical reality of the network (i.e. distance between nodes) directly impacts the behaviour of the logical communications network (i.e. delay between nodes), and is a useful characteristic for differentiating potential misbehaviours.

#### Come back to this and talk about redundancy

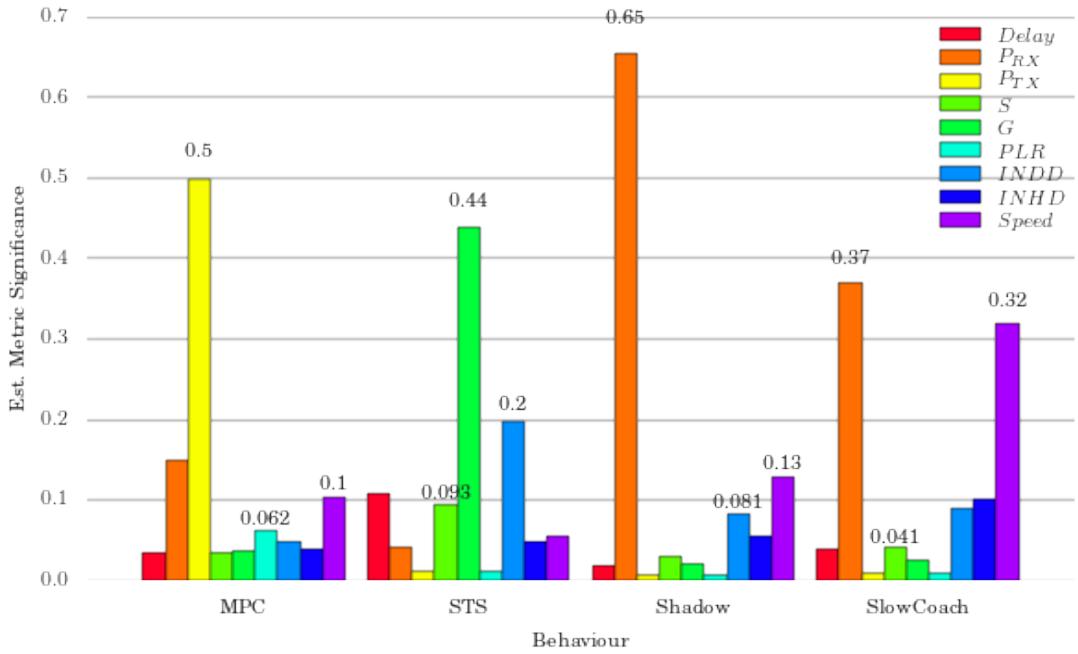
#### 7.2.6 Weight Assessment

From this significance information, a “estimated” signature for each behaviour can be inferred, which can then be fed back into the assessment framework. The aim of this iteration is to minimise the number of weight permutations required to come to a conclusion about the behaviour under observation.

**Figure 7.1:** Plot of Communications Metric Feature Extraction ( $X_{comms}$ )**Figure 7.2:** Plot of Physical Metric Feature Extraction ( $X_{phys}$ )

However, these approximated signatures have no information regarding the “sign” of the  $g,b$  comparison vectors from (C.3), i.e. there is no hint as to whether the relationship is  $g(x) \mapsto \max, b(x) \mapsto \min$  or  $g(x) \mapsto \min, b(x) \mapsto \max$

One option would be to go back to the regression point and expand the combination options to include negative values, signifying inverted  $g,b$  relationships, however this

**Figure 7.3:** Multi Domain Metric Features Extraction ( $X_{merge}$ )**Table 7.1:** Multi Domain Metric Feature Correlation ( $X_{merge}$ )

|              | <i>Delay</i> | <i>P<sub>RX</sub></i> | <i>P<sub>TX</sub></i> | <i>S</i> | <i>PLR</i> | <i>G</i> | <i>INDD</i> | <i>INHD</i> | <i>Speed</i> |
|--------------|--------------|-----------------------|-----------------------|----------|------------|----------|-------------|-------------|--------------|
| Misbehaviour |              |                       |                       |          |            |          |             |             |              |
| MPC          | -0.187       | 0.129                 | 0.579                 | 0.006    | 0.069      | -0.146   | 0.040       | -0.190      | -0.297       |
| STS          | -0.195       | -0.035                | 0.019                 | -0.100   | 0.019      | 0.381    | -0.209      | 0.057       | 0.062        |
| Shadow       | 0.004        | -0.654                | 0.030                 | -0.016   | 0.030      | 0.063    | 0.120       | 0.158       | 0.266        |
| SlowCoach    | -0.157       | -0.533                | 0.013                 | -0.132   | 0.013      | -0.028   | 0.159       | 0.206       | 0.460        |

is combinatorially explosive.<sup>1</sup> Instead, the “significance” weight is permuted against it’s possible combinations of “flips”, i.e. for  $X_s = [0.3, 0.4, 0.01, 0.02, 0.27]$  could also be  $X_s^p = [0.3, -0.4, 0.01, 0.02, 0.27]$  and so on. This sign permutation is filtered based on a threshold value (0.01), so for all indices below that threshold will not be permuted on, halving the number of combinations required for each indices eliminated. This reduces the number of additional assessments required from  $1.9 \times 10^6$  to approximately 500 (when applied to all nine metrics).

The best of these permutations is selected to both maximise the (correct) deviation between each nodes trust perspectives and to minimise the trust value reported for the misbehaving nodes;  $\Delta T \rightarrow \max^+$  (Equation 7.7, results summarised in Table 7.2).

<sup>1</sup>The current version of this analysis uses three metric weights; ignored, standard, emphasised, giving  $3^9 = 19683$  combinations. Expanding this to include inverted standard and inverted emphasised weights would raise that to  $5^9 = 1.9 \times 10^6$

Additionally, a “False Positive” assessment,  $\Delta T^-$  [Equation 7.8](#) is shown in [Table 7.3](#) which encapsulates the average false positive selection rate.

$$\Delta T_{ix} = \frac{\sum_{j \neq x} (\overline{T_{i,j}}^{\forall t})}{N - 1} \quad (7.7)$$

$$\Delta T_{ix}^- = \frac{\sum_{j \neq x} \Delta T_{ij}}{N - 1} \quad (7.8)$$

This isn't right. DT doesn't include it's own value!

Where  $i$  is a given observer,  $x$  is the known misbehaving node,  $\overline{T_{i,j}}^{\forall t}$  is the average weighted trust assessment of node  $j$  observed by node  $i$  across time and  $N$  is the number of nodes.

Conceptually,  $\Delta T_{ix}$  represents the “Distrust” of the target node  $x$ , as the difference in trust value from  $0 \rightarrow 1$ , the higher the better.  $\Delta T_{ix}^-$  represents the average  $\Delta T_{ij}$  for all other nodes, representing the likelihood of another node being as highly distrusted as  $x$ , where positive values indicate that  $x$  is not the obvious outlier, negative values indicate that  $x$  is a very clear outlier, and near-zero values indicate a difficulty in selection of any outlier from the cohort.

Could do with a conceptual graphic showing what these look like, although it'd be messy as all hell

Could also do with a investigation into the deviation of T's; so far most of this analysis averages everything, which is almost certainly not the best approach alone

The “best” weight permutations, as shown in [Table 7.4](#), are applied to untrained datasets for these results.

An exemplar subset of the results is shows in Figs [7.4](#)- [7.15](#), with the “misbehaving node” highlighted with heavier lines, with any observations about the rest of the cohort faded and dashed. For each node assessment, the mean for that assessment over that time period is also included as a solid / dashed line respectively for clarity.

The most intuitively “Communications” behaviour, [MPC](#), scores comfortably in the 90th percentile range in both Communications Domain ([Fig. 7.4](#)) and Full Domain ([Fig. 7.6](#)) trust assessments. As seen in [Table 7.4](#), both the “Full” and “Comms” metric optimisations heavily weigh  $P_{TX}$ , and as this is the metric directly modified by the misbehaviour, it is expected that this is easily discernable using these domain weights. However when this communications information is unavailable, as is the case in the use of Physical Domain metrics alone in [Fig. 7.5](#), the misbehaving node (Alfa) is completely undiscernible compared to the other nodes, with all nodes in the cohort tending to a trust value of 0.5. How this discernibility would fare under varying emphasis of behaviours is an open question

Answer what happens when you vary MPC power variation

**Table 7.2:**  $\Delta T$  across domains and “proposed” behaviours targeting known misbehaving node

| Domain \ Behaviour | MPC  | STS  | Shadow | SlowCoach | Avg. |
|--------------------|------|------|--------|-----------|------|
| Full               | 0.90 | 0.10 | 0.50   | 0.63      | 0.53 |
| Comms              | 0.95 | 0.17 | 0.28   | 0.27      | 0.42 |
| Phys               | 0.02 | 0.02 | 0.43   | 0.76      | 0.31 |
| Comms Alt.         | 0.96 | 0.20 | 0.41   | 0.49      | 0.51 |
| Phys Alt.          | 0.51 | 0.12 | 0.45   | 0.66      | 0.43 |
| Avg.               | 0.67 | 0.12 | 0.41   | 0.56      | 0.44 |

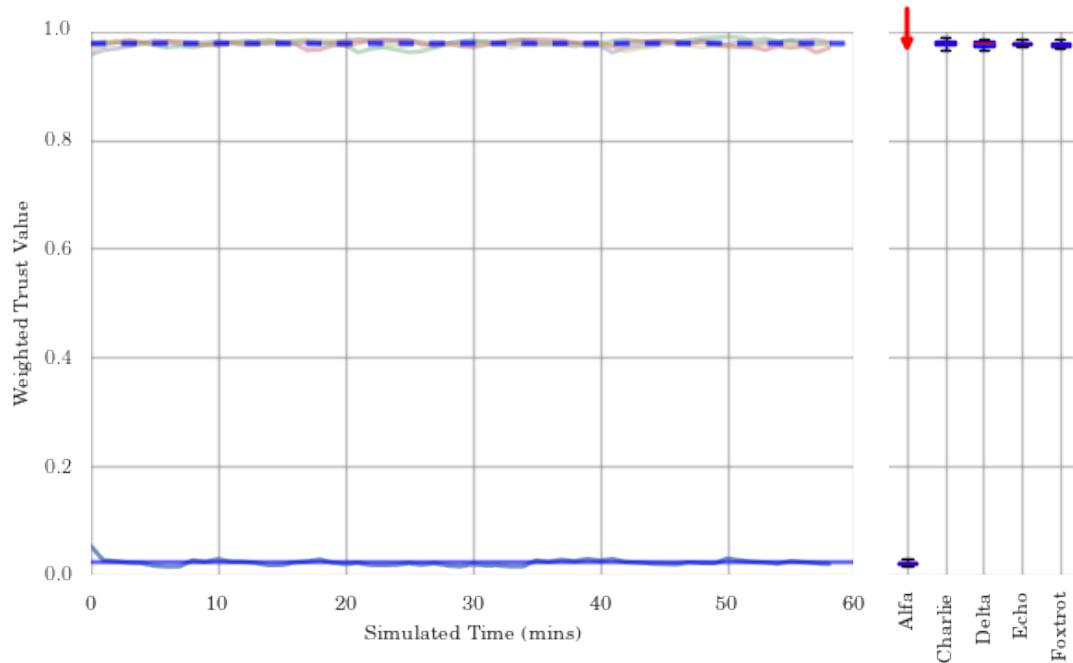
Under the most “subtle” behaviour; **STS**, where no direct metric is being modified in operation, but where the behaviour is effectively in the “Application layer” of the networking stack, the picture is far more murky. Comparing Figs 7.7 and 7.9, while there is a reasonable dip in the misbehavior’s trust assessment, the high level of variance across the cohort is such that this “mistrust” triggering is neither consistent or obvious. From Table 7.4, the metric of import is  $G$ , the Offered Load on the network, and given it’s negative weighting, this matches the expectation that the node doing “less than it’s fair share” is potentially misbehaving. Unfortunately this is the case across the **STS** responses, where in Table 7.2 we have summarized out general results, **STS** has by far and away the lowest average  $\Delta T$  in all domains. Interestingly however is the observation that Comms-only trust performs slightly better than Full trust weighting.

Referring to Figs 7.1 and Fig. 7.3, it’s clear that the offered load ( $G$ ) is the almost singular feature of this behaviour, due to it’s almost completely logical behaviour that is only loosely coupled to the state of the environment. The massive emphasis placed on load could only be diminished by putting it together in a larger ensemble. In Figs 7.10 and 7.12, the misbehaving node is much more obvious than in **STS**, which is moderately surprising for a physically-focused behaviour. Further, there is a roughly 20% improvement when incorporating the full metric space.

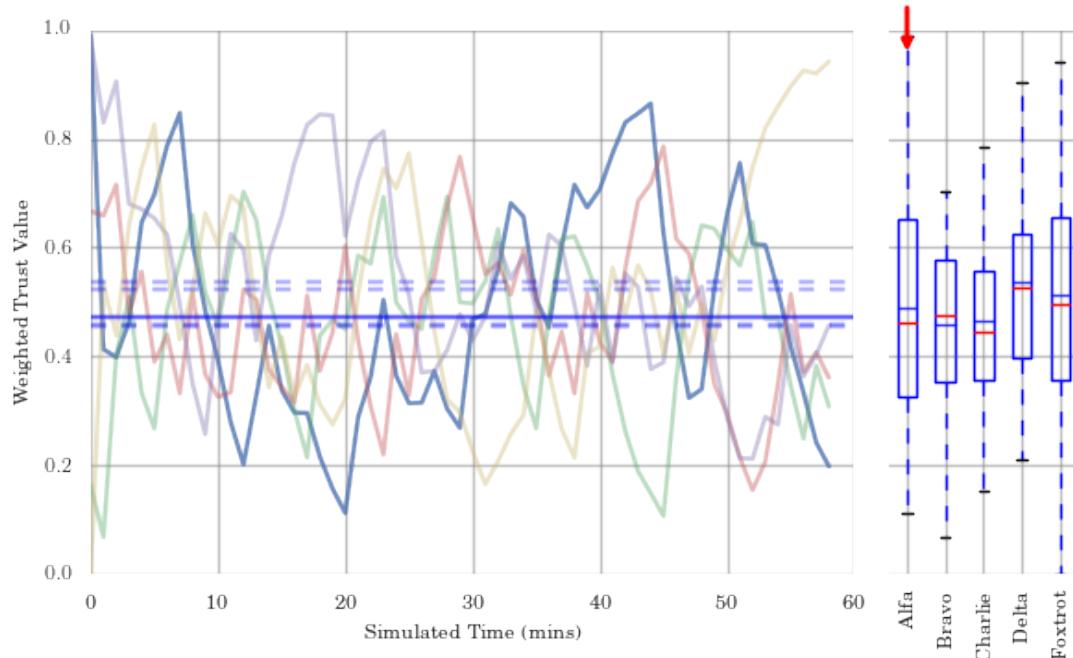
From Table 7.2, the Shadow behavior is the most consistently detectable behaviour across selected metric domains.

See Section D.1 for a full collection of graphs showing the comparison of the malicious nodes trust value against the instantenous mean of the remaining cohort. See Subsection D.2.1 for a full collection of graphs showing the comparison of non-malicious nodes trust value against the individual values of the remaining cohort. See Subsection D.2.2 for a full collection of graphs showing the comparison of non-malicious nodes trust value against the instantenous mean of the remaining cohort.

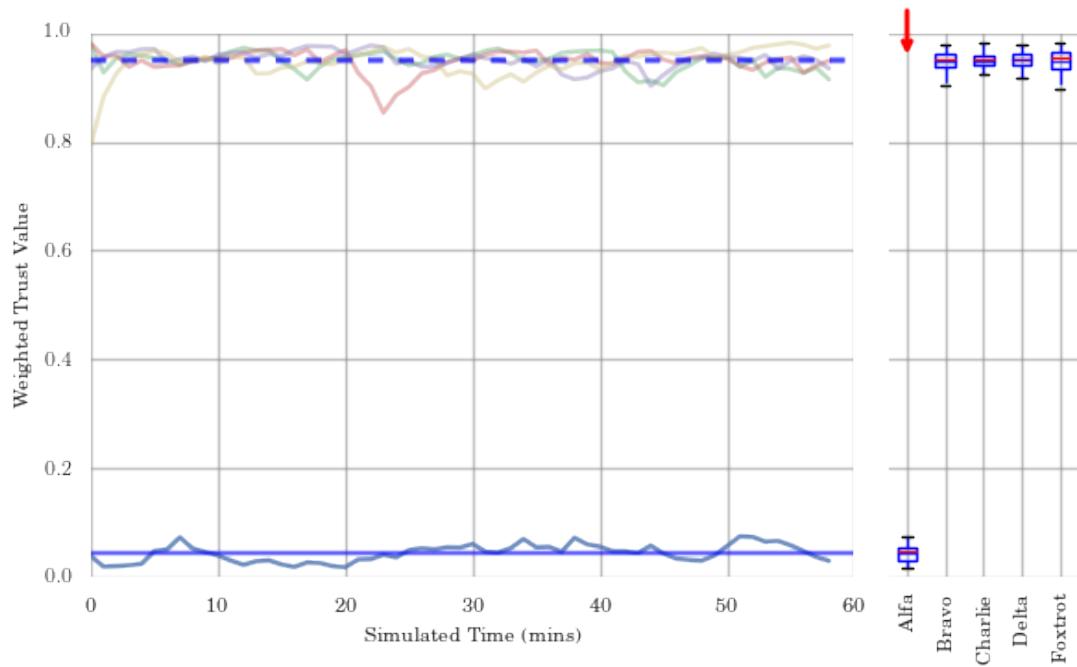
Haven’t worked out a clever way of automatically generating both the basic domain and alternate domain texes easily



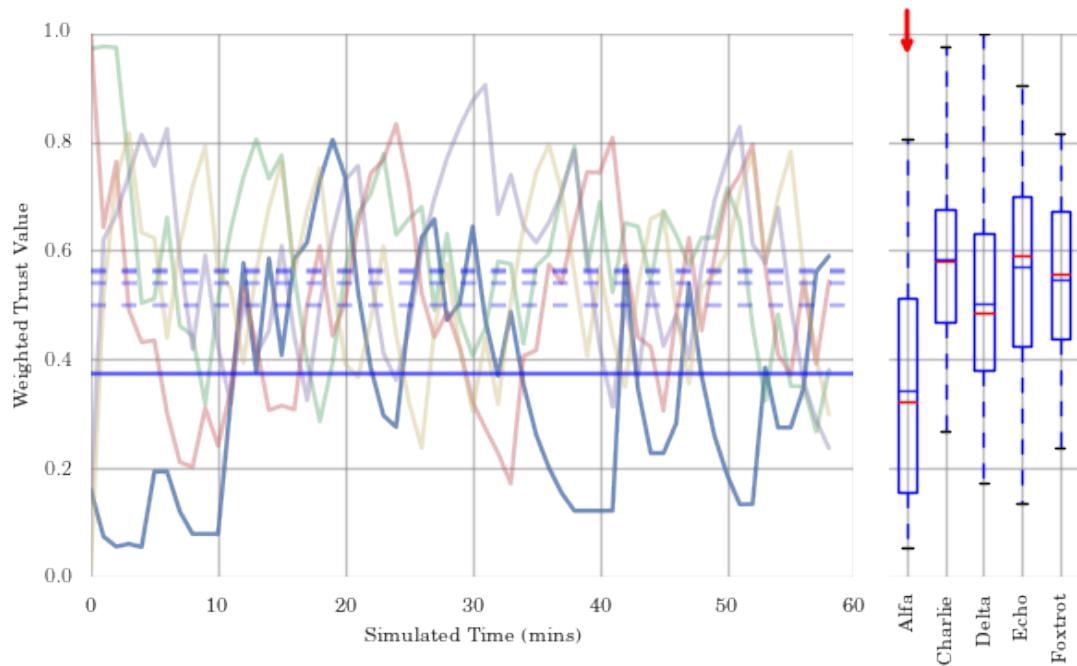
**Figure 7.4:** MPC Comms Metric Trust (showing cohort trust assessments)



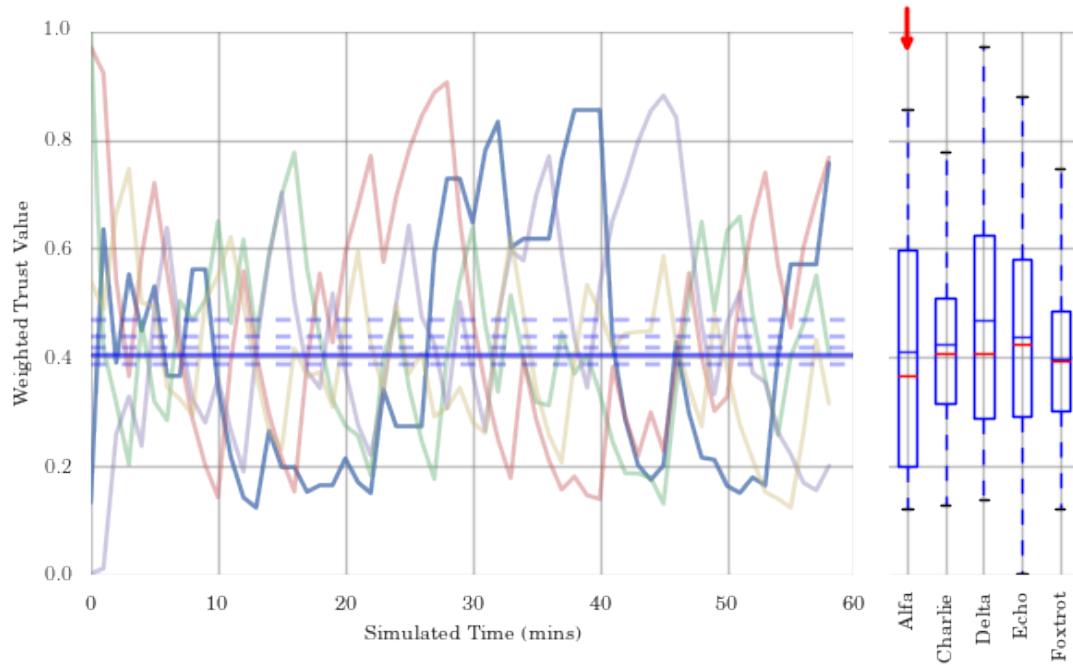
**Figure 7.5:** MPC Physical Metric Trust (showing cohort trust assessments)



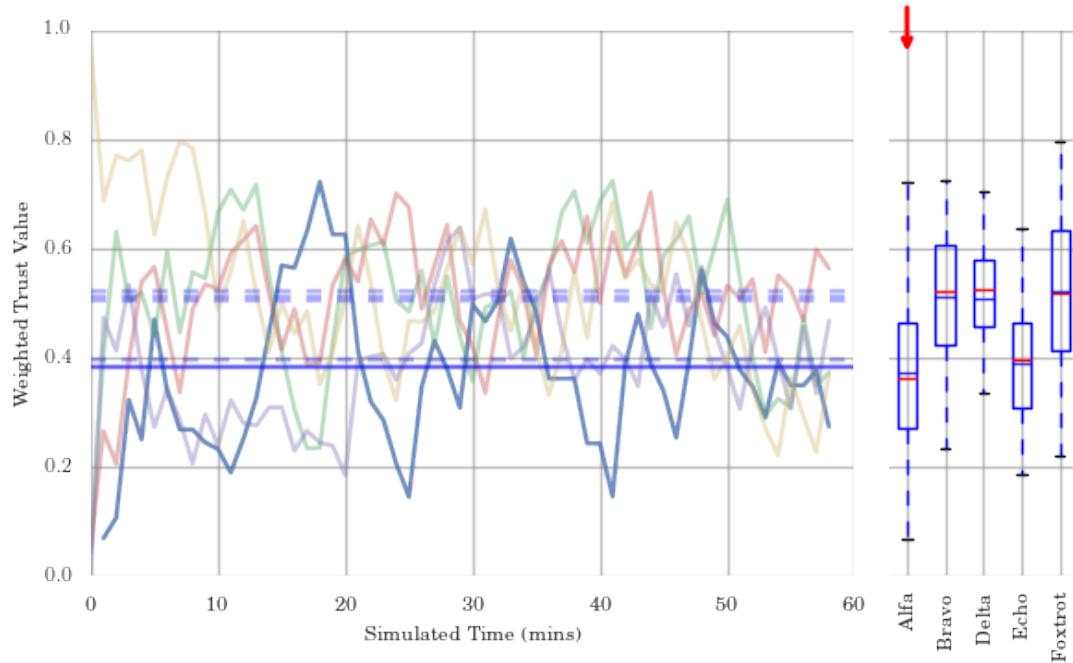
**Figure 7.6:** MPC Full Metric Trust (showing cohort trust assessments)



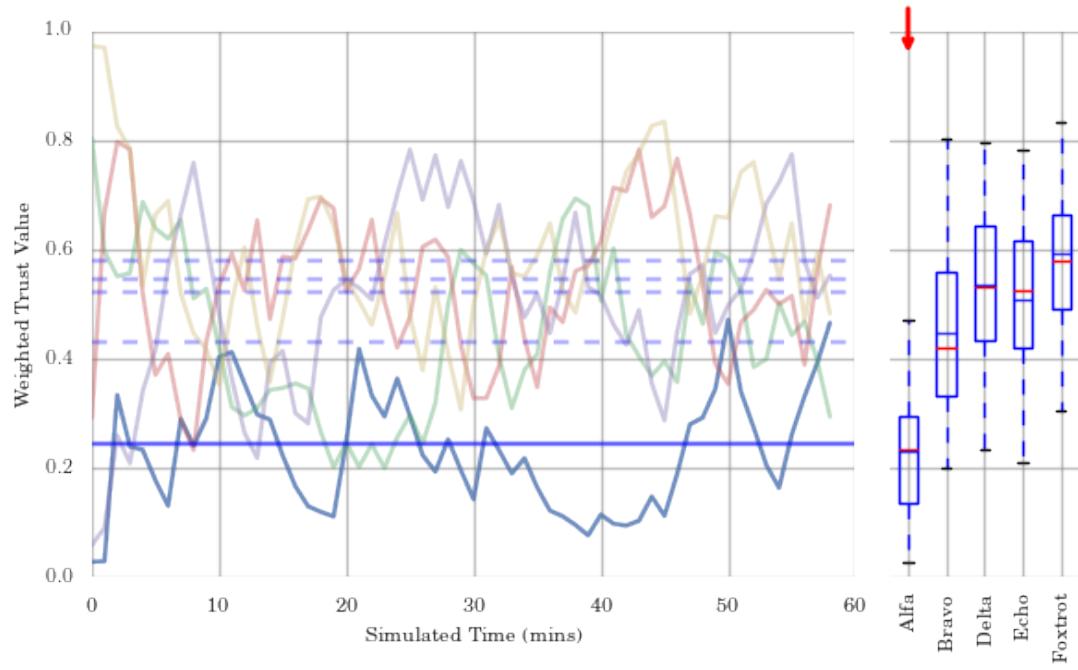
**Figure 7.7:** STS Comms Metric Trust (showing cohort trust assessments)



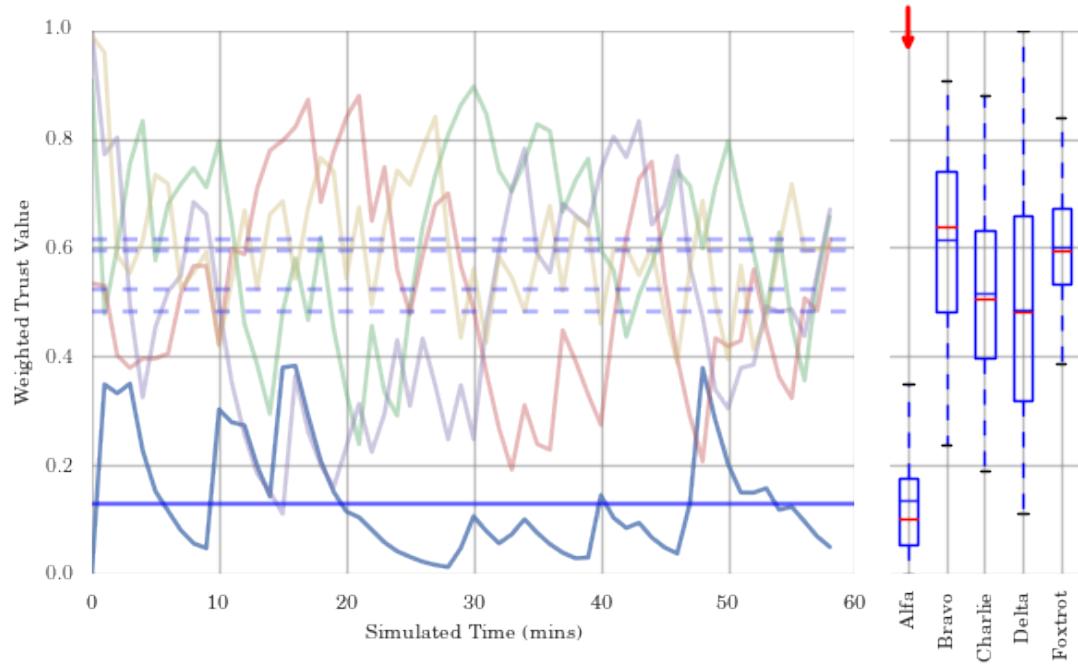
**Figure 7.8:** STS Physical Metric Trust (showing cohort trust assessments)



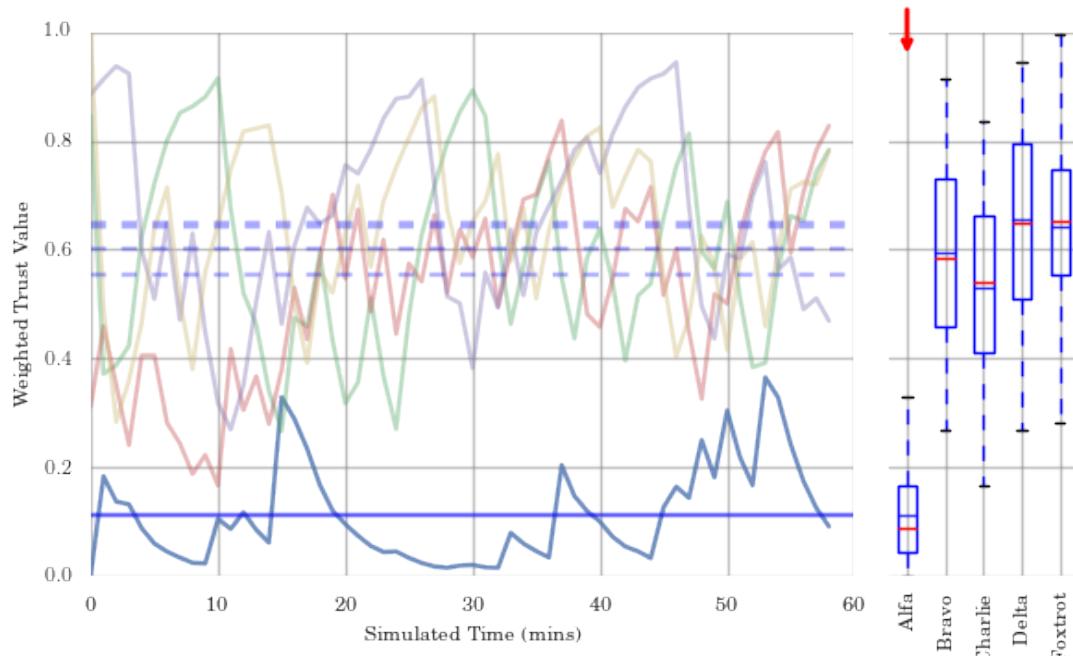
**Figure 7.9:** STS Full Metric Trust (showing cohort trust assessments)



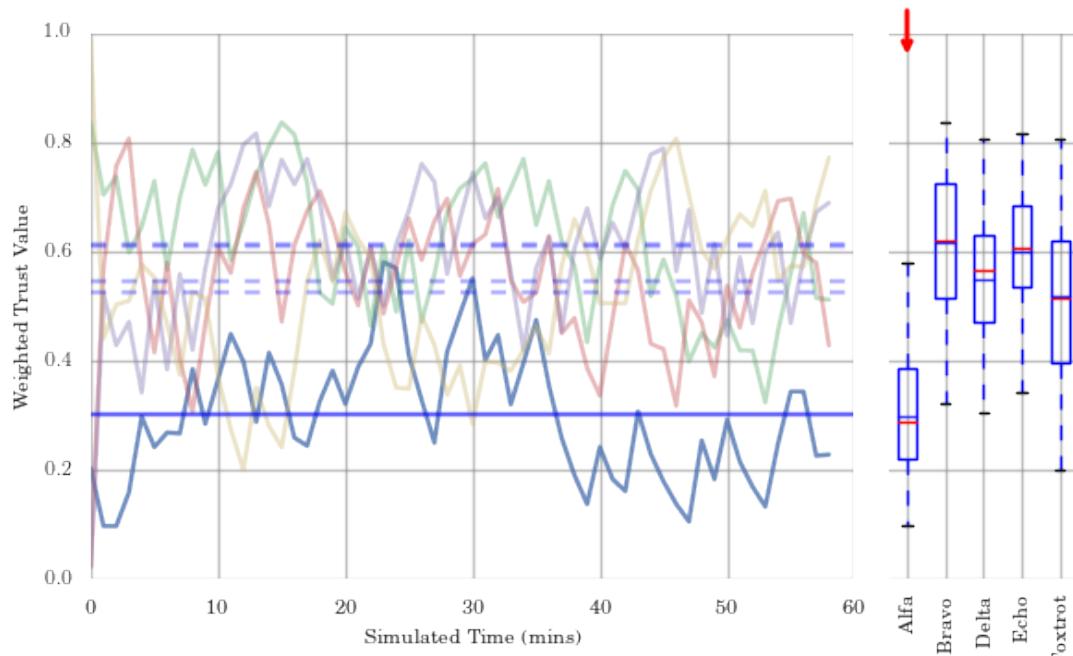
**Figure 7.10:** Shadow Comms Metric Trust (showing cohort trust assessments)



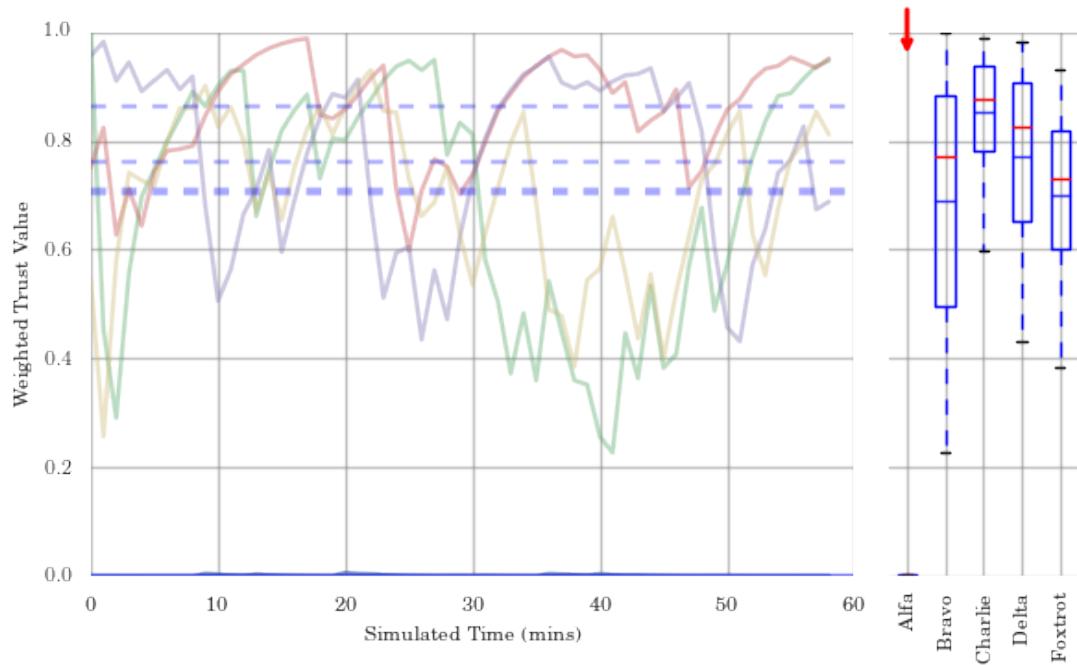
**Figure 7.11:** Shadow Physical Metric Trust (showing cohort trust assessments)



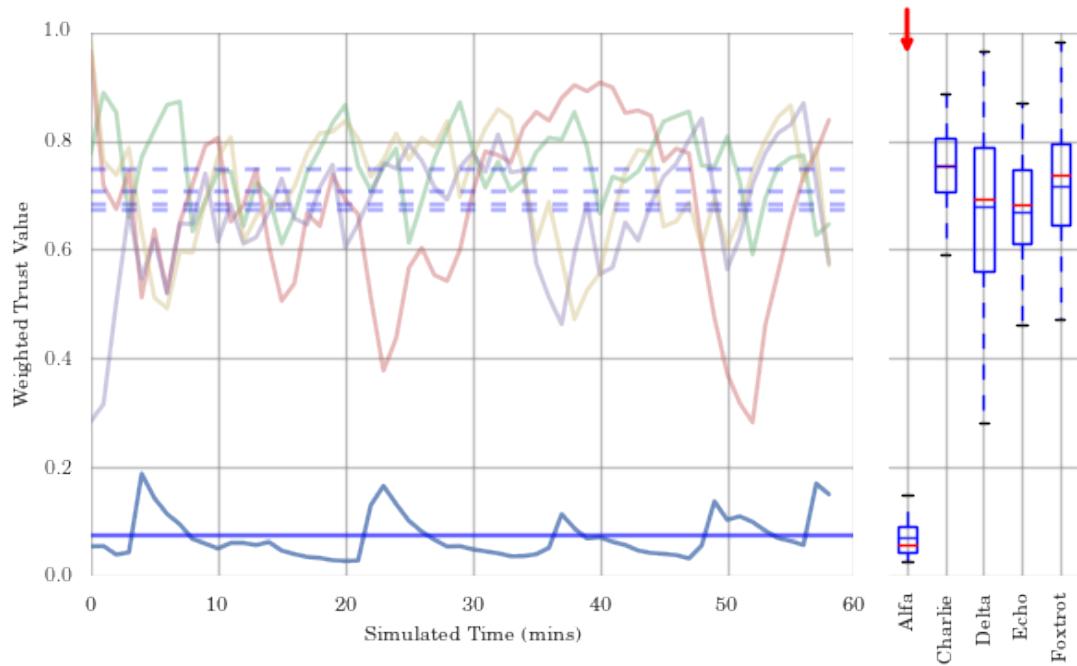
**Figure 7.12:** Shadow Full Metric Trust (showing cohort trust assessments)



**Figure 7.13:** SlowCoach Comms Metric Trust (showing cohort trust assessments)



**Figure 7.14:** SlowCoach Physical Metric Trust (showing cohort trust assessments)



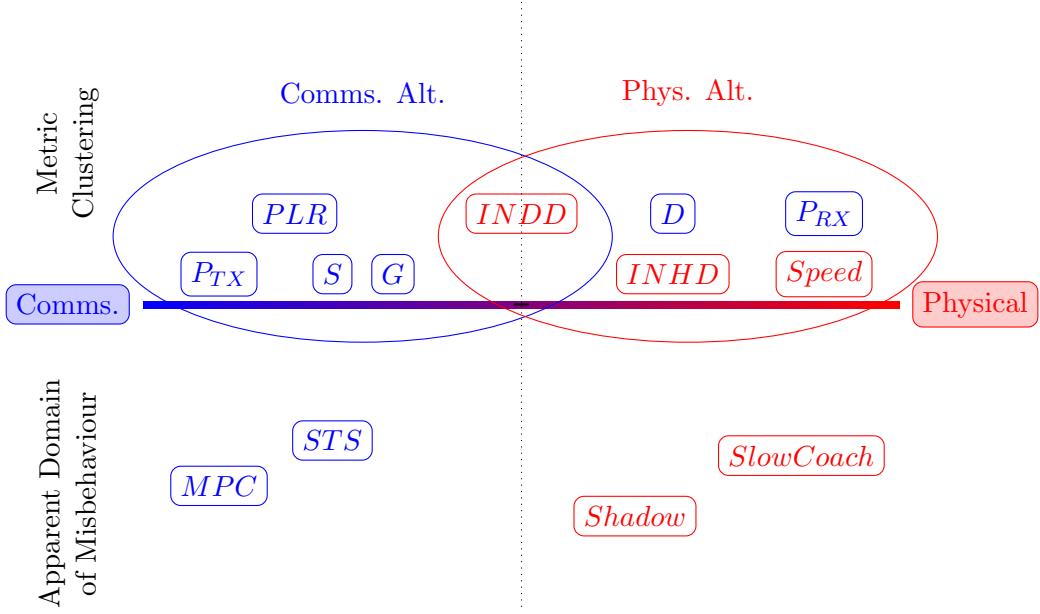
**Figure 7.15:** SlowCoach Full Metric Trust (showing cohort trust assessments)

**Table 7.3:**  $\Delta T^-$  False Positive assessments across domains and “proposed” behaviours across non-misbehaving nodes

| Domain \ Behaviour | MPC   | STS   | Shadow | SlowCoach | Avg.  |
|--------------------|-------|-------|--------|-----------|-------|
| Full               | -0.23 | -0.03 | -0.12  | -0.16     | -0.13 |
| Comms              | -0.24 | -0.04 | -0.07  | -0.07     | -0.10 |
| Phys               | -0.01 | -0.01 | -0.11  | -0.19     | -0.08 |
| Comms Alt.         | -0.24 | -0.05 | -0.10  | -0.12     | -0.13 |
| Phys Alt.          | -0.13 | -0.03 | -0.11  | -0.17     | -0.11 |
| Avg.               | -0.17 | -0.03 | -0.10  | -0.14     | -0.11 |

**Table 7.4:** Optimised metric vector weights per domain trained upon and behaviour targeted

| Domain, Behaviour \ Metric |           | <i>Delay</i> | <i>P<sub>RX</sub></i> | <i>P<sub>TX</sub></i> | <i>S</i> | <i>G</i> | <i>PLR</i> | <i>INDD</i> | <i>INHD</i> | <i>Speed</i> |
|----------------------------|-----------|--------------|-----------------------|-----------------------|----------|----------|------------|-------------|-------------|--------------|
| Full                       | MPC       | -0.033       | 0.154                 | 0.495                 | 0.034    | -0.035   | 0.062      | -0.047      | -0.039      | -0.101       |
|                            | STS       | -0.106       | 0.042                 | 0.010                 | 0.095    | 0.438    | 0.010      | -0.194      | -0.049      | -0.055       |
|                            | Shadow    | 0.019        | 0.656                 | 0.007                 | -0.030   | -0.021   | 0.007      | -0.081      | -0.054      | -0.125       |
|                            | SlowCoach | 0.040        | 0.373                 | 0.009                 | -0.042   | -0.025   | 0.009      | -0.087      | 0.099       | -0.316       |
| Comms                      | MPC       | 0.045        | 0.068                 | 0.665                 | 0.029    | -0.043   | 0.150      |             |             |              |
|                            | STS       | 0.098        | 0.083                 | 0.047                 | 0.118    | -0.608   | 0.046      |             |             |              |
|                            | Shadow    | -0.358       | 0.279                 | 0.025                 | 0.119    | 0.193    | 0.024      |             |             |              |
|                            | SlowCoach | -0.082       | 0.309                 | 0.021                 | 0.090    | 0.478    | 0.020      |             |             |              |
| Phys                       | MPC       |              |                       |                       |          |          |            | -0.439      | -0.383      | -0.178       |
|                            | STS       |              |                       |                       |          |          |            | -0.729      | -0.164      | -0.108       |
|                            | Shadow    |              |                       |                       |          |          |            | -0.555      | -0.142      | -0.304       |
|                            | SlowCoach |              |                       |                       |          |          |            | -0.285      | -0.118      | -0.597       |
| Comms Alt.                 | MPC       |              |                       | 0.731                 | 0.019    | -0.024   | 0.211      | -0.014      |             |              |
|                            | STS       |              |                       | 0.040                 | -0.131   | -0.444   | 0.038      | -0.348      |             |              |
|                            | Shadow    |              |                       | 0.033                 | -0.124   | -0.104   | 0.032      | -0.707      |             |              |
|                            | SlowCoach |              |                       | 0.029                 | -0.164   | -0.184   | 0.028      | -0.595      |             |              |
| Phys Alt.                  | MPC       | 0.043        | 0.389                 |                       |          |          |            | -0.311      | -0.075      | -0.183       |
|                            | STS       | -0.356       | 0.095                 |                       |          |          |            | -0.235      | -0.135      | -0.179       |
|                            | Shadow    | 0.081        | 0.577                 |                       |          |          |            | -0.097      | 0.070       | -0.175       |
|                            | SlowCoach | -0.106       | 0.309                 |                       |          |          |            | -0.067      | 0.099       | -0.420       |



**Figure 7.16:** Assumptions made about the relevant domains of impact / detectability of misbehaviours, and domain relevance of metrics, may not be optimal

### 7.3 Conclusion

In this chapter we demonstrate that in harsh environments, multi-domain trust assessment can perform better on average than single-domain counterparts, both in terms of robustness and sensitivity, but also covering a wider region of the potential behaviour space,

The extension of the methodologies of multi-vector trust into the marine space are already demonstrated, however including information from physical observations of actors in a network enables the detection and identification of a much wider range of behaviours. We also demonstrate a method for assessing trust metrics in harsh environments in terms of their relative significance, and a method for establishing classification signatures for misbehaviours.

It is to be noted that this presented method is significantly more computationally intensive than the relatively simple Hermes / OTMF algorithms communications only algorithms, and is exponential in complexity as metrics and/or domains are added. The repeated metric re-weighting required for real time behaviour detection is therefore an area that requires optimization. More work needs to be done to characterise how worthwhile this approach is compared to a separate synthesis approach where by MTFM-style trust is generated and assessed on a per-domain basis and subsequently fuzzed.

For greater fidelity and more optimal results, a wider range of weights can be used in the initial regression step; however this is computationally expensive given that weighting is applied to each perspective (i.e. observer/target node pair) for each trust assessment time step, presenting 15 perspectives at each time interval in the 6 node case.

Every effort has been made to avoid over-training the dataset, using cross validating sampling for regression and "best weight" generation, however more meta-analysis is required to further demonstrate the functionality of this process.

# Appendix A

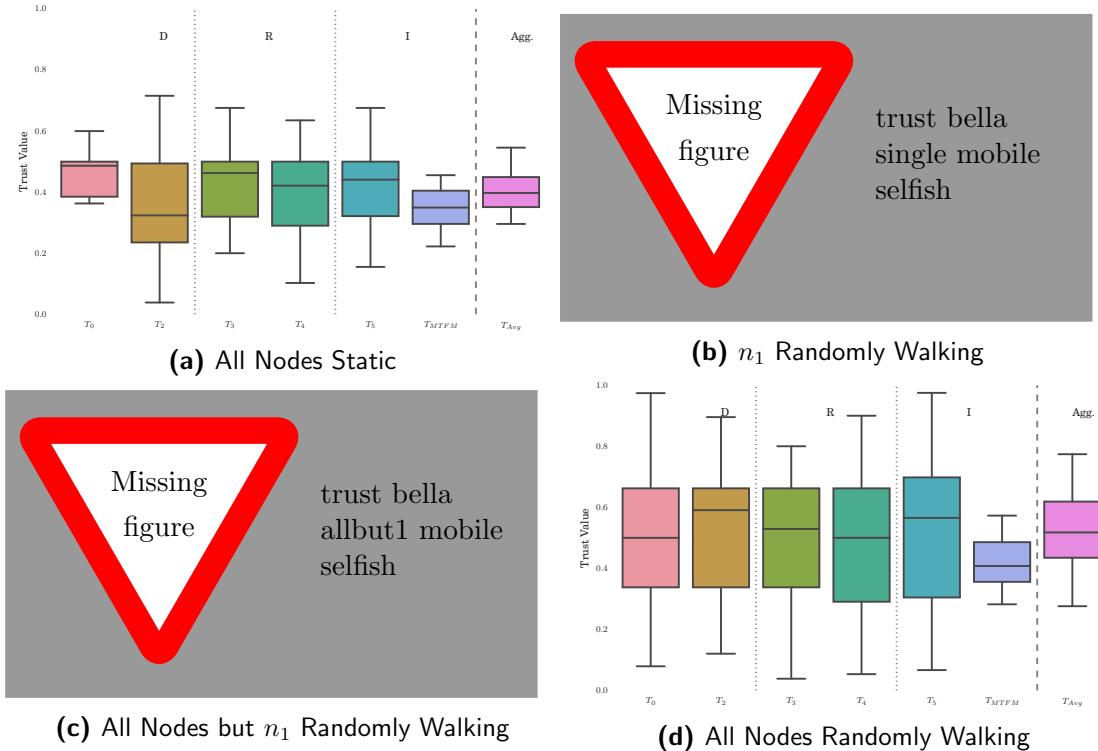
## Orphan Sections

### A.1 Metric Weighting

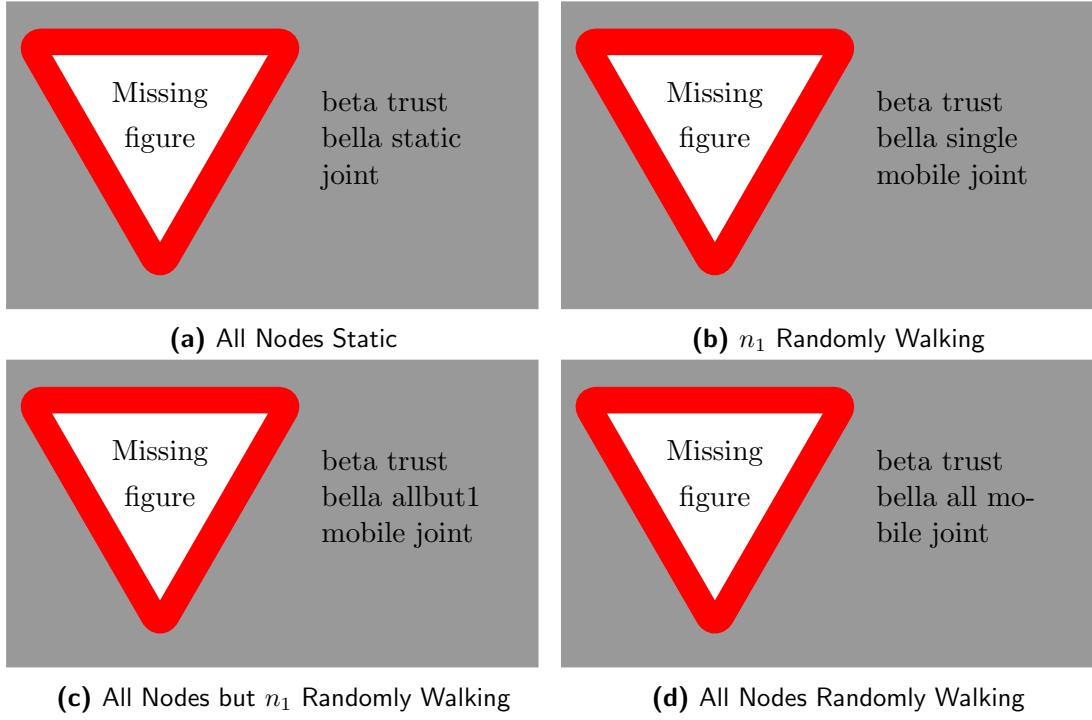
### A.2 UNEDITED PROSE: Real Time Grey Systems

#### *Incoming Train of Consciousness*

For a given metric set  $X$  such that  $X = x_1, \dots, x_M$  representing the  $M$  different types of measurement generated by an observer. If these metrics are not synchronised, for instance if they are interrupt driven such as communications-based observations, generating more abstract measurements requires inherent assumptions about “how to



**Figure A.1:** MTFM Trust assessments for varying mobility options in the selfish case



**Figure A.2:** Beta Trust time varying assessments for of  $n_1$  varying mobility options

accumulate the data while you wait”. For instance, in [14], we demonstrated a periodic trust assessment framework for autonomous marine environments, in such an environment, to establish useful, generalised, data, it was necessary to wait for a relatively long time to accumulate enough data to make assessments. However, this left many ‘smells’; data was being left in-buffer for a long time before being used to make decisions, and by the time the data was collated and processed, it could be wildly different from the reality. Further, while some periods could be extremely sparse or even empty, others could be extremely busy with many records having to be averaged down to provide a ‘single period’ response. Therefore, the implementation of a suitable sequence buffer version of the framework would be beneficial.

Such a sequence buffer framework would involve a tracking predictor that would provide best-guess estimates of an interpolated value for a metric between value updates, and a back-propagation algorithm to retroactively update historical assessments of that metrics so as to better inform any abstracted trust value predictor.

I had initially thought that such a back-propagator would be a total mess as I’d imagined that significant-model-breaking would potentially indicate untrustworthy behaviour, but this is stupid since the per-metric-model has the least information of anyone and is simply there to provide better intermediate values and has no / limited direct impact on the overall trust behaviour.

This back propagation will probably be a pain to implement as it’d require a retroactive reassessment of trust and could get really messy if it was interrupt driven, but it’s better not to prematurely optimise.

### A.3 From end of Defense Trust Conclusions

In order to contextualise the discussions on trust in mixed and hybrid networks, an exemplar scenario is considered. That scenario builds on existing Maritime Autonomy Framework (MAF) investigations(Mollet, J. et al., 2012. Osprey Task 37 Activity 8 - Unmanned Systems Operations: Technical Assurance Work Package - Security Issues and Mitigations - Final Report,)

While the initial assessment does not cover the MHPC PT CONUSE recommendations, it provides a starting point for future trust research in UxV operations. In order to constrain the scope of this project, a single operational scenario will be analysed within documented MCHP CONUSE(Rudge, A., Chapman, K. & Goddard, N., 2012. Information Management for MHPC: Research Strategy,), of Route/Area Survey within both peacetime and wartime contexts, with a Beyond Line of Sight (BLOS) operator. This scenario will be a minimal MCM operation in a littoral area. In field assets will consist of:

- Two squads consisting of Three UUVs, (tacitly modelled on the in-service REMUS 100 UUV), and a USV providing acoustic-RF relay capabilities per-squad
- an UAV providing BLOS Comms
- A remote human operator (MCMV / PJHQ / etc)



The differential between the peacetime and wartime contexts will be an attempted capture of a UUV by a manned surface-based FIS asset. Clearly, this paper has a limited scope and does not attempt to cover every aspect of a trustworthy system.



## Appendix B

# Human Factors related to Trusted Operation of Autonomous Systems

This work has largely considered autonomous systems as entities of wider systems, implicitly involving human operators/agents in some part of the desired operation. We refer to these systems as [Autonomous Collaborative System \(ACS\)](#). As described in [Chapter 3](#), Operational Trust has two main aspects, trust in the system to behave as expected and trust in the interfaces between systems (human/machine and machine/machine). Of all of the interfaces in an Autonomous Collaborative System, the most problematic is that arguably that between the [ACS](#) and the human operator / team of operators. Cummings identified the main challenges to [HSC](#), summarised below:[41]

### B.1 Information Overload

Operator efficiency exhibits an optimum at moderate levels of cognitive engagement, above which cognitive ability is overloaded and performance drops (Otherwise known as the Yerkes-Dodson Law). Additionally, in the case of under-engagement, operators can fall foul of boredom, and become desensitised to changing factors. *However, predicting this point of over-saturation is an open psychophysiological research problem.*

### B.2 Adaptive Automation

Automation is well tailored to consistent levels of activity. This is quite simply not the case many domains. Particularly in defence and military applications, activity is characterised by long periods of “routine” punctuated by high intensity, usually unpredictable, activity. At those interfaces between “calm” and “storm”, where real time situational awareness is imperative, temporary Information Overload is highly probable. Adaptive Automation enables autonomous systems to increase their [LOA](#) based on specific events

in the task environment, changes in operator performance or task loading, or physiological methods. It is taken as given that for routine operations, and increased LOA reduces operator workload, and vice versa. However, this relationship is highly task dependent and can create severe problems in cases of LOA being greater, or indeed lesser, than is required. In the cases of overly-high LOA, operator skill is degraded, situational awareness is reduced as the operator is not as engaged, and the automated system may not be able to handle unexpected events, requiring the operator to take over, which, given the previous points, is a difficult prospect. Alternatively, in sub-optimal LOA, Information Overload can result in the case of high intensity situations, but also the system can fall foul of overly-sensitive human cognitive biases, false positive pattern detection, boredom, and complacency in the case where less is going on. Therefore, as a corollary to Information Overload challenges, there is a need to define the interrelationship between levels of situational activity (or risk) and appropriate levels of automation. *Under what circumstances can AA be used to change the LOA of a system? Does the autonomous system or the human decide to change LOA? What LOAs are appropriate for what circumstances?*

### B.3 Distributed Decision Making

In a modern, non-hierarchical, often distributed or cellular military management system (Network Centric Warfare doctrine for example), tools are increasingly being used to mitigate information asymmetry within command and control. A simple example of this is shared watch-logs in Naval operations, providing temporal collaboration between watch-teams separated in time. The DoD Global Information Grid is another example of a spatial collaborative framework. Recent work has demonstrated the power of collaborative analysis and human-machine shared sensing technologies even with low levels of training on the part of the operators providing superior results and resource efficiencies than either humans or machines alone in survey and search-and-rescue scenarios (Ahmed et al.2014)

Check Security

. As these temporal and spatial collaboration tools increase in complexity and ability, decisions that previously required SA that was only available at higher echelons within the standard hierarchy are available to commanders on the ground, or even to individual team members, enabling the potential for informed decisions to be taken faster and more effectively, enabled by automated strategies to present relevant information to teams based on the operational context. However there are a range of operational, legal, psychological and technical challenges that need to be addressed before confidence in these distributed management structures can be established. Studies into situational awareness sharing techniques (telepresent table-top environments, video conferencing,

and interactive whiteboards) have generally yielded positive results, however investigations into interruptive-communications (such as instant messaging chat) have demonstrated a negative impact on operational efficiency. In short, the biggest problem with distributed decision making in the context of supervisory systems is that *there is no consensus on whether it is advantageous or not, and what magnitude of operational delta is introduced, if any.*

## B.4 Complexity

Beyond simple Information Overload, increasing complexity of information presented to operators is having a negative effect on operational efficiency. In HSC, displays are designed to reduce complexity, introducing abstractions with an aim to presenting the minimum amount of information to the operator required to maintain an accurate and up-to-date mental model of the environmental and operational state. This has led to the development of many domain specific decision support interfaces, however, in academic research, there has been nothing but mixed results. One commonly raised negative is the general bias on the cool factor of interfaces. Immersive 3D visual, aural, or haptic interfaces that at first appraisal seem to provide more approachable information to the operator, and are indeed tacitly preferred by operators in use. However, there has not been any evidence to demonstrate performance improvement when using these tools, and in-fact, *improving the “fidelity” of the interfaces has led to operators overly-relying on these representations of the environment rather than remaining engaged in the environment.*

## B.5 Cognitive Biases and Failing Heuristics

In many areas, operators and commanders are required to make rapid decisions with imperfect information, driven by massively increased information availability and rates of change in areas such as battlefield tactics and global finance markets. However, Human decision making isn't always rational (especially under pressure), and operators use personally derived heuristics to make “rational shortcuts”. This is a double edged sword, where these heuristics can be employed to greatly reduce the normative cognitive load in a stressful situation, but also introduce destructive biases, where these shortcuts make assumptions that don't bear out in reality.

For example, in the context of decision support systems, “Autonomy Bias” has been observed as a complement to the already well known “Confirmation Bias”<sup>1</sup> and “Assimilation Bias”<sup>2</sup>, where operators that have been provided with a “correct” answer by a

<sup>1</sup>Confirmation Bias is the tendency for people to preferentially select from available information that supports pre-existing beliefs or hypotheses.

<sup>2</sup>Assimilation Bias is often thought of as a subset of Confirmation Bias, whereby it specifies that instead of seeking out information supporting of current views, any incoming data is interpreted as being supportive of a particular view without questioning that view, even if it appears contradictory.

decision support system do not look (or see, depending on perspective) for any contradictory information, and will unquestionably follow, increasing error rates significantly.

This behaviour isn't only the reserve of decision support systems, but also in the generic allocation of operator attention; scheduling heuristics are used to decide how much time tasks should be worked on, and time and again, humans are found to be far from optimal in this regard, especially in time-pressured scenarios where these heuristics are in even more demand. Even when operators are given optimal scheduling rules, these quickly fall apart, often due to primary task efficiency degradation after interruption. This highlights a critical interface in the adoption of complex autonomous systems that still demand Man in the loop functionality; if a system is required to have full-time concentrated supervision (e.g. flying a UCAV), but also event-based reactive decision making (e.g. alerts from non-critical subsystems), both tasks are negatively impacted. In an assessment of factors influencing trust in autonomous vehicles and medical diagnosis support systems, Carlson et al also identified that a major factor in an operator or users trust in a system was not only dependant on past performance and current accuracy but also on "soft factors" such as the branding and reputation of the manufacture / designer. (Carlson et al. 2014)

Check Security

Further, autonomous decision support / detection / classification systems have an "uncanny valley" to overcome in terms of accuracy, in that there is a dangerous period when such systems are used but not perfect, but operators become complacent, causing an increased error rate, until such a time that those autonomous systems can match or exceed the detection rates of their human counterparts.

## Appendix C

# Grey System Theory and Grey Trust Assessment

### C.1 Grey numbers, operators and terminology

Grey numbers are used to represent values where their discrete value is unknown, where that number may take its possible value within an interval of potential values, generally written using the symbol  $\oplus$ . Taking  $a$  and  $b$  as the lower and upper bounds of the grey interval respectively, such that  $\oplus \in [a, b] | a < b$ . The “field” of  $\oplus$  is the value space  $[a, b]$ . There are several classifications of grey numbers based on the relationships between these bounds.

don't think classification is the right word here

Black and White numbers are the extremes of this classification; such that  $\dot{\oplus} \in [-\infty, +\infty]$  and  $\ddot{\oplus} \in [x, x] | x \in \mathbb{R}$  or  $\oplus(x)$ . It is clear that white numbers such as  $\ddot{\oplus}$  have a field of zero while black numbers have an infinite field.

Grey numbers may represent partial knowledge about a system or metric, and as such can represent half-open concepts, by only defining a single bound; for example  $\underline{\oplus} = \oplus(\underline{x}) \in [x, +\infty]$  and  $\overline{\oplus} = \oplus(\overline{x}) \in [-\infty, x]$ .

Primary operations within this number system are as follows;

$$\oplus_1 + \oplus_2 \in [a_1 + a_2, b_1 + b_2] \quad (\text{C.1a})$$

$$-\oplus \in [-b, -a] \quad (\text{C.1b})$$

$$\oplus_1 - \oplus_2 = \oplus_1 + (-\oplus) \quad (\text{C.1c})$$

$$\oplus_1 \times \oplus_2 \in [\min(a_1 a_2, a_1 b_2, b_1 a_2, b_2 a_2), \quad (\text{C.1d})$$

$$\max(a_1 a_2, a_1 b_2, b_1 a_2, b_2 a_2)]$$

$$\oplus^{-1} \in [b^{-1}, a^{-1}] \quad (\text{C.1e})$$

$$\oplus_1 / \oplus_2 = \oplus_1 \times \oplus_2^{-1} \quad (\text{C.1f})$$

$$\oplus \times k \in [ka, kb] \quad (\text{C.1g})$$

$$\oplus^k \in [a^k, b^k] \quad (\text{C.1h})$$

where  $k$  is a scalar quantity.

## C.2 Whitenisation and the Grey Core

The characterisation of grey numbers is based on the encapsulation of information in a grey system in terms of the grey numbers core ( $\hat{\oplus}$ ) and its degree of greyness ( $g^\circ$ ). If the distribution of a grey number field is unknown and continuous,  $\hat{\oplus} = \frac{a+b}{2}$ .

Non-essential grey numbers are those that can be represented by a white number obtained either through experience or particular method. [84] This white value is represented by  $\tilde{\oplus}$  or  $\oplus(x)$  to represent grey numbers with  $x$  as their whitenisation. In some cases depending on the context of application, particular gray numbers may temporarily have no reasonable whitenisation value (for instance, a black number). Such numbers are said to be Essential grey numbers.

## C.3 Grey Sequence Buffers and Generators

eqs of sequence buffers and partial derivs

Given a fully populated value space, sequence buffer operations are used to provide abstractions over the dataspace. These abstractions can be *weakening* or *strengthening*. In the weakening case, these operations perform a level of smoothing on the volatility of a given input space, and strengthening buffers serve to highlight and A powerful tool in grey system theory is the use of grey incidence factors, comparing the “likeness” of one value against a cohort of values. This usefulness applies particularly well in the case of multi-agent trust networks, where the aim is to detect and identify malicious or maladaptive behaviour, rather than an absolute assessment of “trustworthiness”.

## C.4 Grey Trust

Grey Theory performs cohort based normalization of metrics at runtime. This creates a more stable contextual assessment of trust, providing a “grade” of trust compared to other observed entities in that interval, while maintaining the ability to reduce trust values to a stable assessment range for decision support without requiring every environment entered into to be characterised. Grey assessments are relative in both fairly and unfairly operating cohorts. Entities will receive mid-range trust assessments if there are no malicious actors as there is no-one else “bad” to compare against.

Guo[13] demonstrated the ability of Grey Relational Analysis (GRA)[85] to normalise and combine disparate traits of a communications link such as instantaneous throughput, received signal strength, etc. into a Grey Relational Coefficient, or a “trust vector”.

In [13], the observed metric set  $X = x_1, \dots, x_M$  representing the measurements taken by each node of its neighbours at least interval, is defined as  $X = [\text{packet loss rate}, \text{signal}$

strength, data rate, delay, throughput]. The trust vector is given as

$$\begin{aligned}\theta_{k,j}^t &= \frac{\min_k |a_{k,j}^t - g_j^t| + \rho \max_k |a_{k,j}^t - g_j^t|}{|a_{k,j}^t - g_j^t| + \rho \max_k |a_{k,j}^t - g_j^t|} \\ \phi_{k,j}^t &= \frac{\min_k |a_{k,j}^t - b_j^t| + \rho \max_k |a_{k,j}^t - b_j^t|}{|a_{k,j}^t - b_j^t| + \rho \max_k |a_{k,j}^t - b_j^t|}\end{aligned}\quad (\text{C.2})$$

where  $a_{k,j}^t$  is the value of a observed metric  $x_j$  for a given node  $k$  at time  $t$ ,  $\rho$  is a distinguishing coefficient set to 0.5,  $g$  and  $b$  are respectively the ‘‘good’’ and ‘‘bad’’ reference metric sequences from  $\{a_{k,j}^t, k = 1, 2 \dots K\}$ , e.g.  $g_j = \max_k(a_{k,j}^t)$ ,  $b_j = \min_k(a_{k,j}^t)$  (where each metric is selected to be monotonically positive for trust assessment, e.g. higher throughput is always better).

Weighting can be applied before generating a scalar value which allows the identification and classification of untrustworthy behaviours.

$$[\theta_k^t, \phi_k^t] = \left[ \sum_{j=0}^M h_j \theta_{k,j}^t, \sum_{j=0}^M h_j \phi_{k,j}^t \right] \quad (\text{C.3})$$

Where  $H = [h_0 \dots h_M]$  is a metric weighting vector such that  $\sum h_j = 1$ , and in the basic case,  $H = [\frac{1}{M}, \frac{1}{M} \dots \frac{1}{M}]$  to treat all metrics evenly.  $\theta$  and  $\phi$  are then scaled to  $[0, 1]$  using the mapping  $y = 1.5x - 0.5$ . The  $[\theta, \phi]$  values are reduced into a scalar trust value by  $T_k^t = (1 + (\phi_k^t)^2 / (\theta_k^t)^2)^{-1}$ . This trust value minimises the uncertainties of belonging to either best ( $g$ ) or worst ( $b$ ) sequences in (C.2).

**MTFM** combines this GRA with a topology-aware weighting scheme(C.4) and a fuzzy whitenization model(C.5). There are three classes of topological trust relationship used; Direct, Recommendation, and Indirect. Where an observing node,  $n_i$ , assesses the trust of another, target, node,  $n_j$ ; the Direct relationship is  $n_i$ ’s own observations  $n_j$ ’s behaviour. In the Recommendation case, a node  $n_k$ , which shares Direct relationships with both  $n_i$  and  $n_j$ , gives its assessment of  $n_j$  to  $n_i$ . The Indirect case, similar to the Recommendation case, the recommender  $n_k$ , does not have a direct link with the observer  $n_i$  but  $n_k$  has a Direct link with the target node,  $n_j$ . These relationships give us node sets,  $N_R$  and  $N_I$  containing the nodes that have recommendation or indirect, relationships to the observing node respectively.

$$\begin{aligned}T_{i,j}^{\text{MTFM}} &= \frac{1}{2} \cdot \max_s \{f_s(T_{i,j})\} T_{i,j} + \frac{1}{2} \frac{2|N_R|}{2|N_R| + |N_I|} \sum_{n \in N_R} \max_s \{f_s(T_{i,n})\} T_{i,n} \\ &\quad + \frac{1}{2} \frac{|N_I|}{2|N_R| + |N_I|} \sum_{n \in N_I} \max_s \{f_s(T_{i,n})\} T_{i,n}\end{aligned}\quad (\text{C.4})$$

Where  $T_{i,n}$  is the subjective trust assessment of  $n_i$  by  $n_n$ , and  $f_s = [f_1, f_2, f_3]$  given as:

$$\begin{aligned} f_1(x) &= -x + 1 \\ f_2(x) &= \begin{cases} 2x & \text{if } x \leq 0.5 \\ -2x + 2 & \text{if } x > 0.5 \end{cases} \\ f_3(x) &= x \end{aligned} \tag{C.5}$$

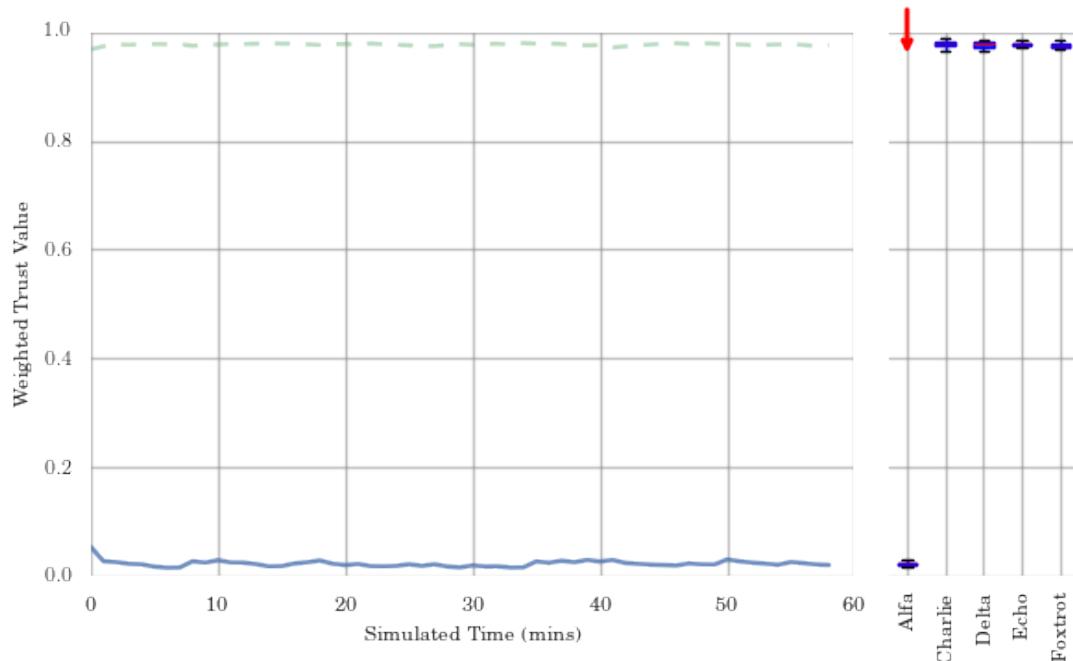
Grey System Theory, by it's own authors admission, hasn't taken root in it's originally intended area of system modelling [84]. However, given it's tentative application to MANET trust, taking a Grey approach on a per metric benefit has qualitative benefits that require investigation; the algebraic approach to uncertainty and the application of "essential and non essential greyness", whiteisation, and particularly grey buffer sequencing allow for the opportunity to generate continuous trust assessments from multiple domains asynchronously.

## Appendix D

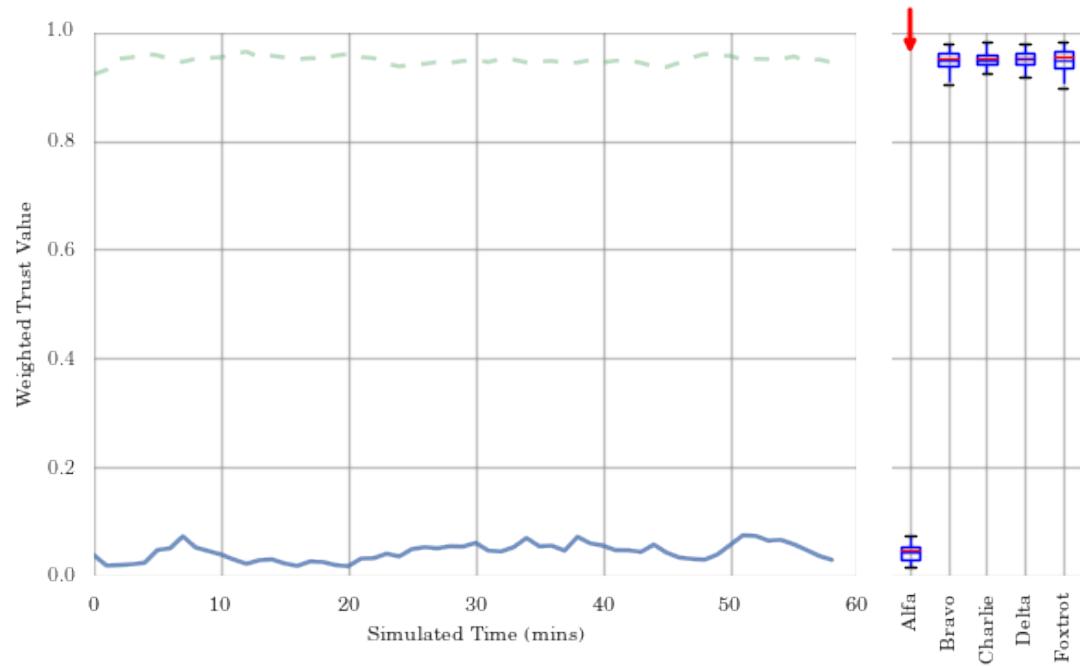
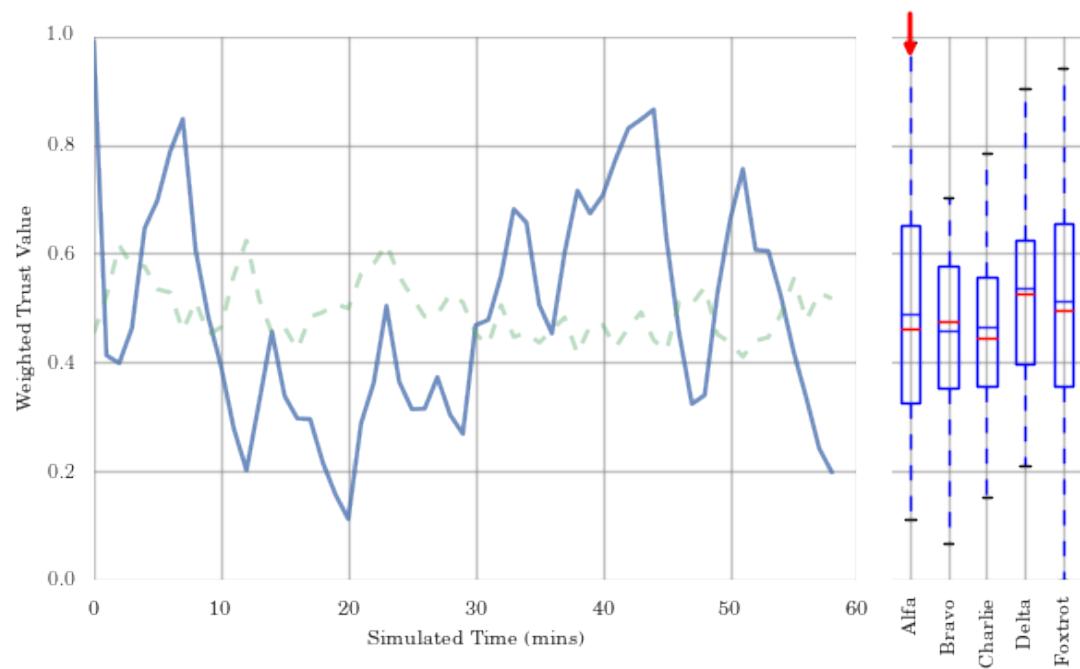
# Additional Graphs

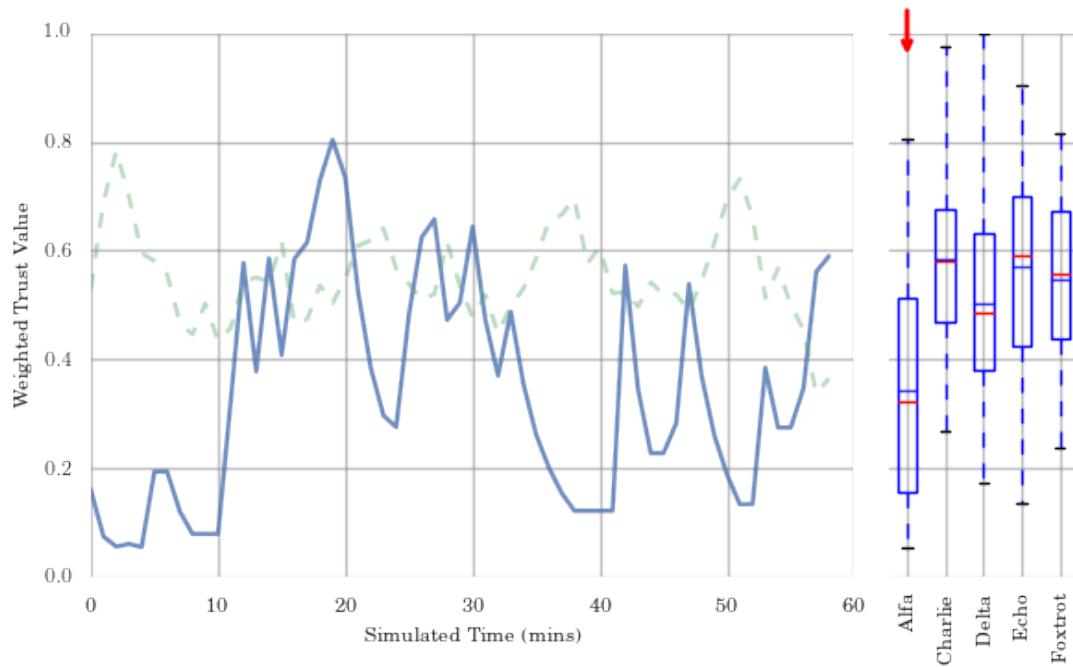
### D.1 From Subsection 7.2.6: Mean-Weighted Multi-Domain Trust Results

These graphs show the trust variability of an individual against a mean-averaged trust response of the remaining cohort of nodes.

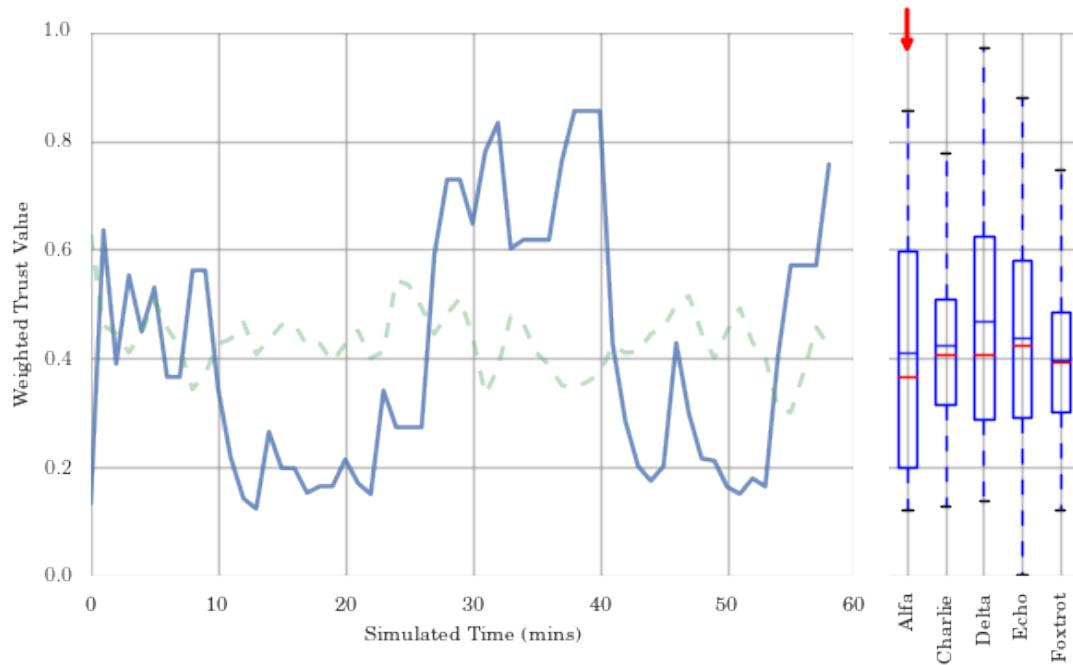


**Figure D.1:** MPC Comms Metric Trust (showing mean of non-misbehaving nodes)

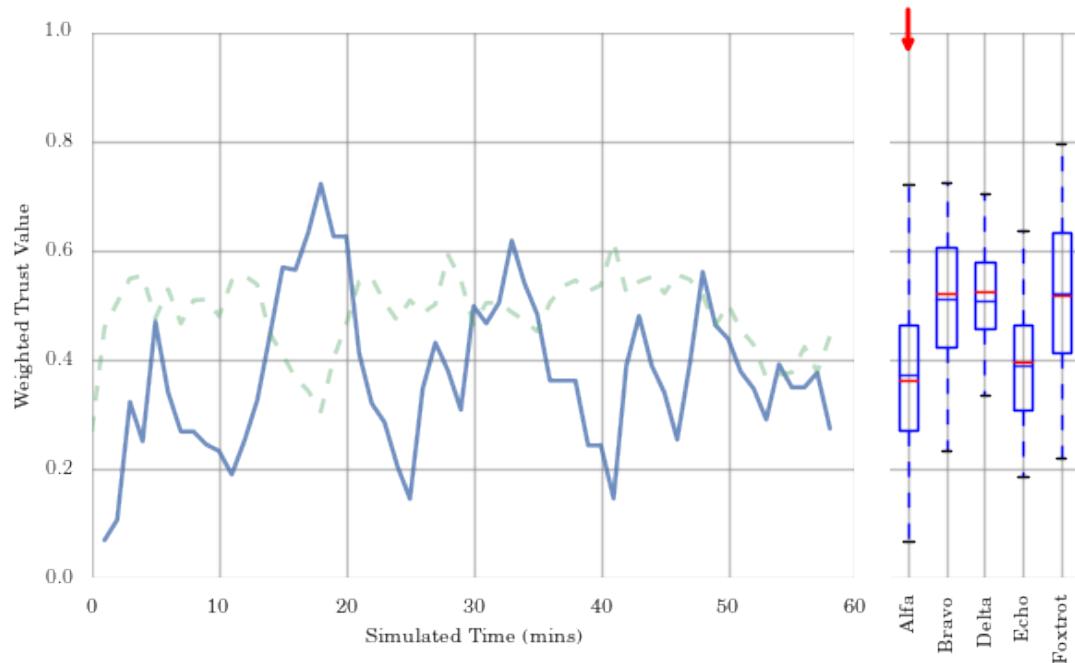




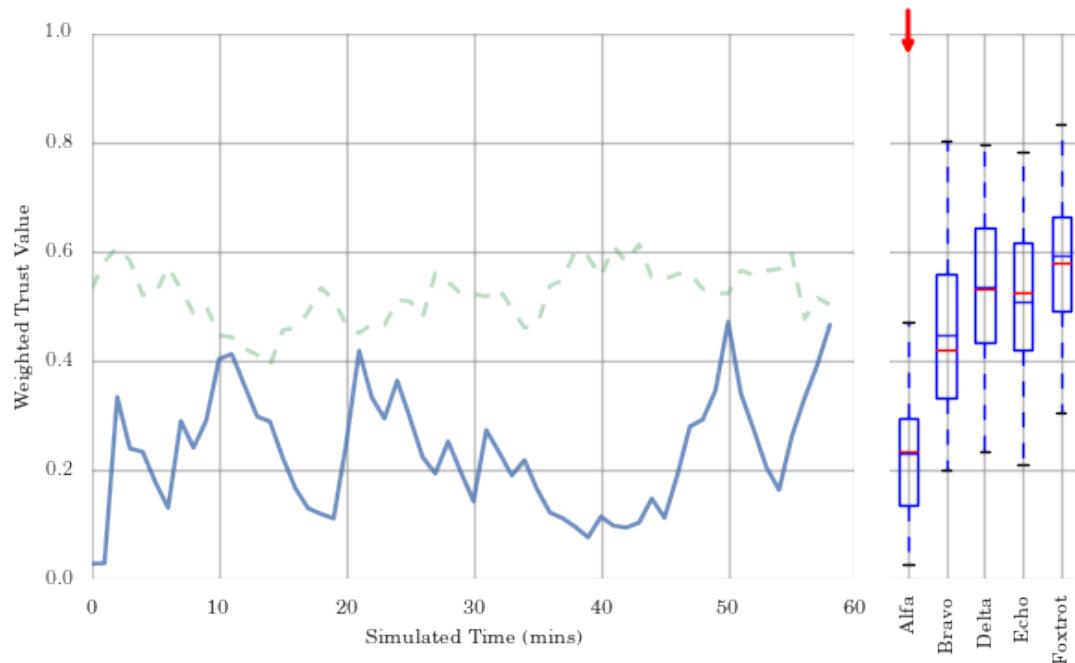
**Figure D.4:** STS Comms Metric Trust (showing mean of non-misbehaving nodes)



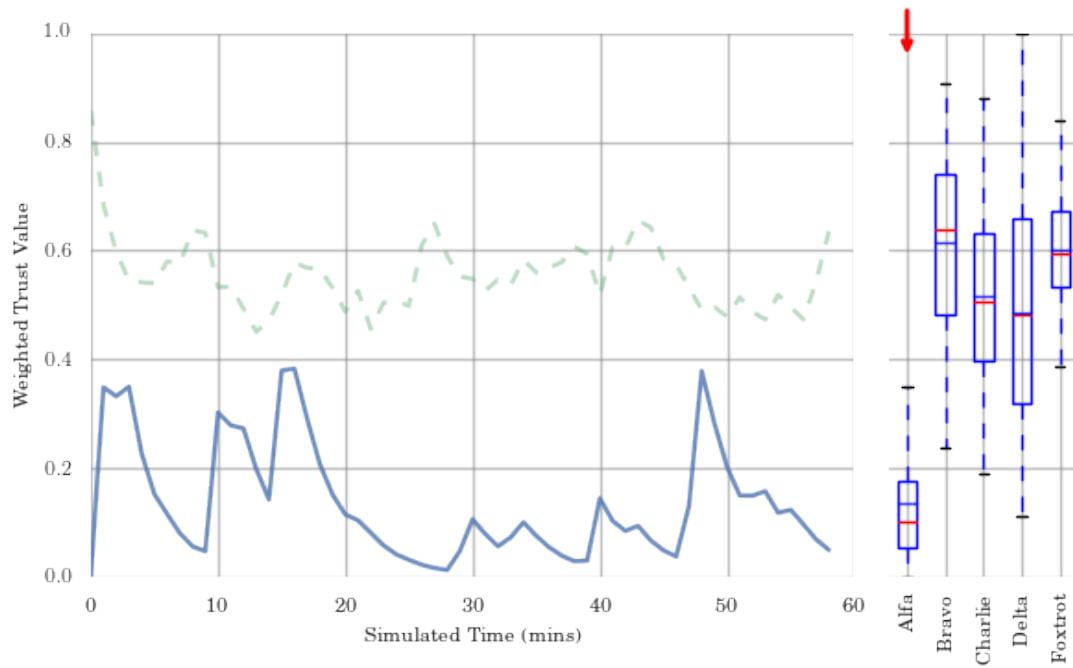
**Figure D.5:** STS Physical Metric Trust (showing mean of non-misbehaving nodes)



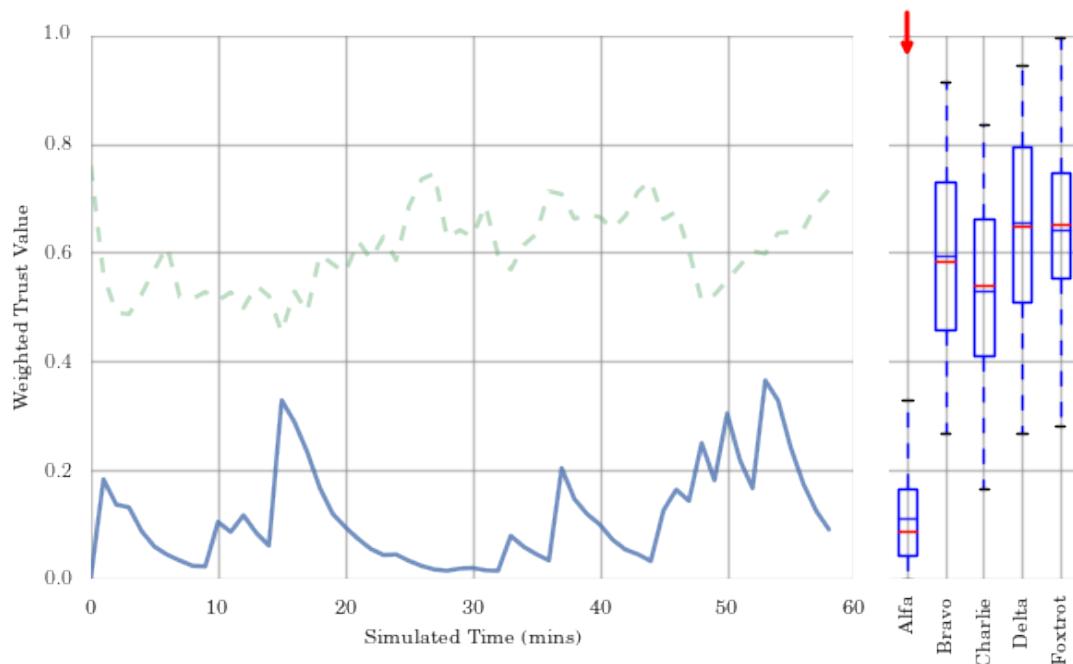
**Figure D.6:** STS Full Metric Trust (showing mean of non-misbehaving nodes)



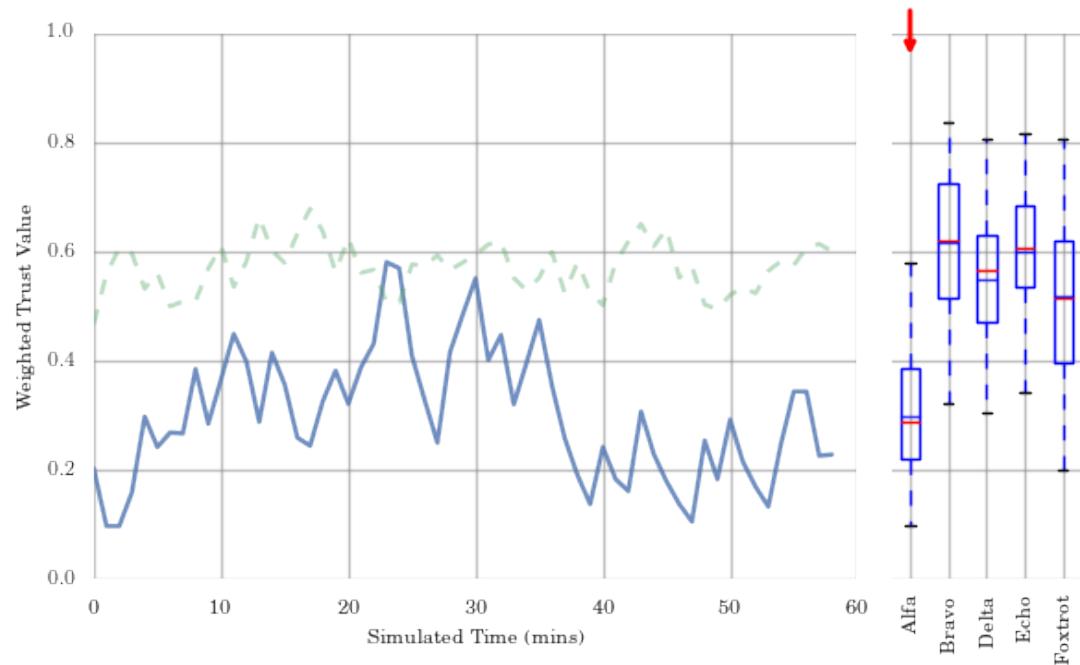
**Figure D.7:** Shadow Comms Metric Trust (showing mean of non-misbehaving nodes)



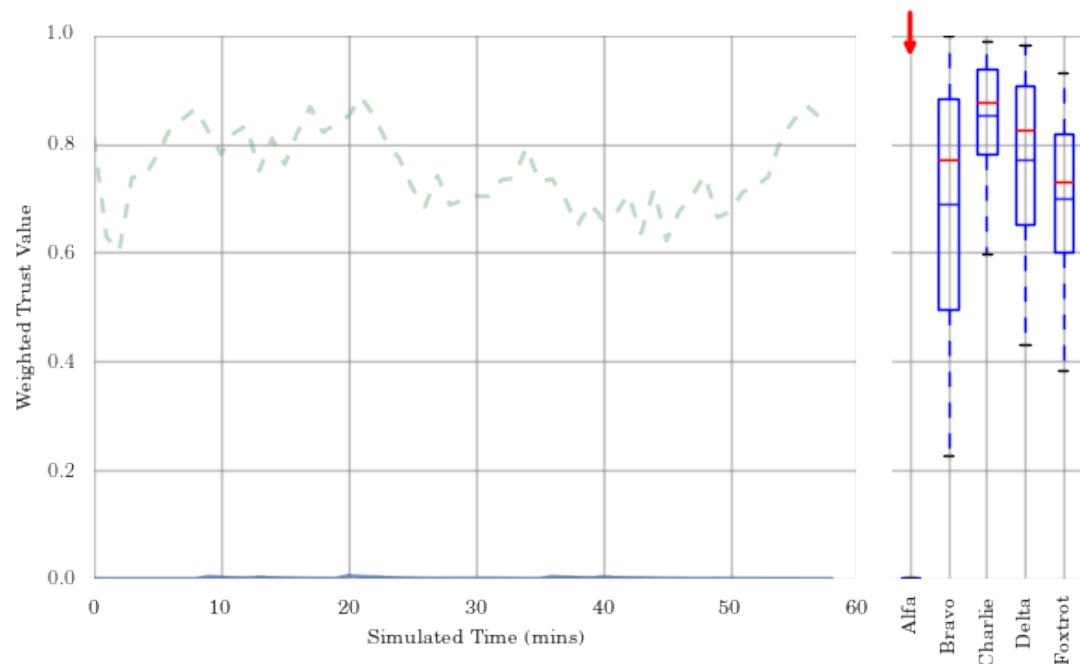
**Figure D.8:** Shadow Physical Metric Trust (showing mean of non-misbehaving nodes)



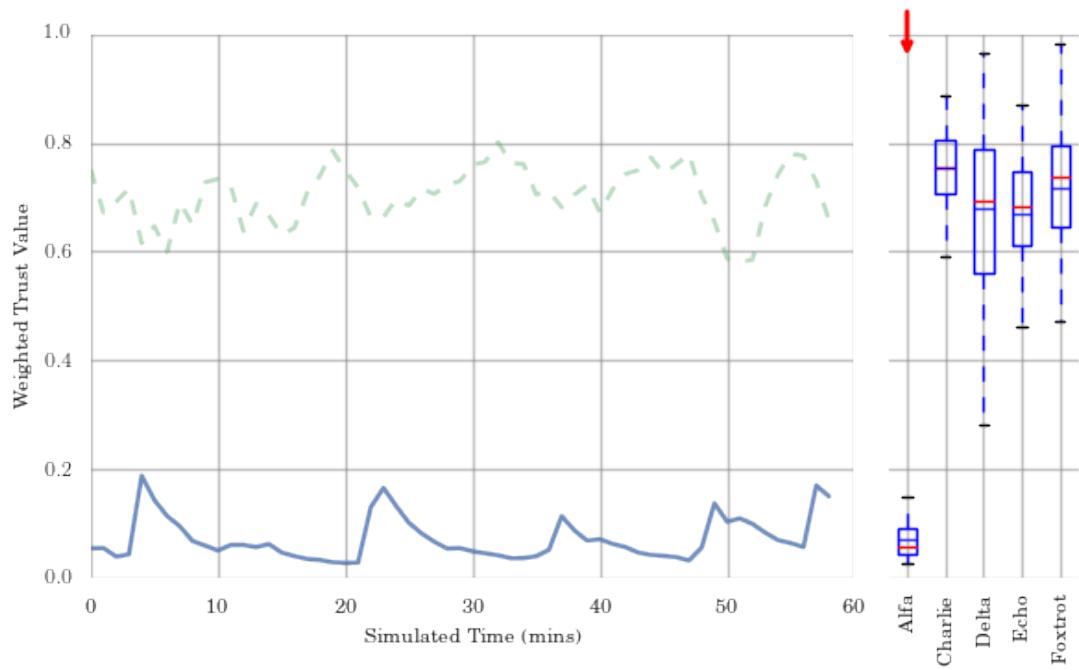
**Figure D.9:** Shadow Full Metric Trust (showing mean of non-misbehaving nodes)



**Figure D.10:** SlowCoach Comms Metric Trust (showing mean of non-misbehaving nodes)



**Figure D.11:** SlowCoach Physical Metric Trust (showing mean of non-misbehaving nodes)

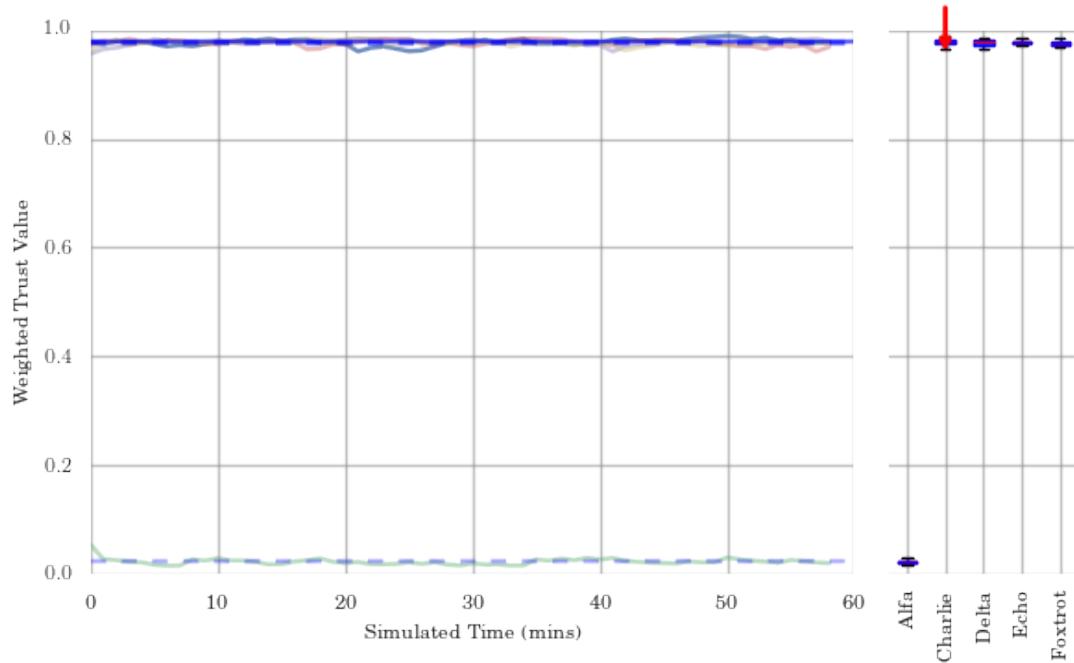


**Figure D.12:** SlowCoach Full Metric Trust (showing mean of non-misbehaving nodes)

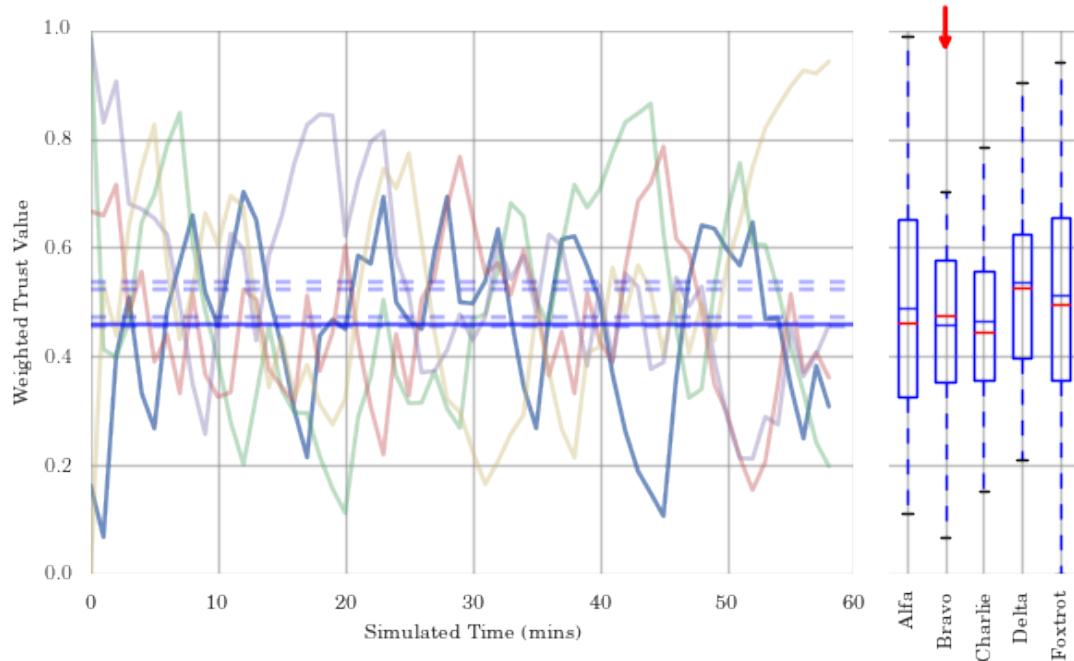
## D.2 From Subsection 7.2.6: Multi-Domain Trust Results

### Targeting Non-malicious nodes

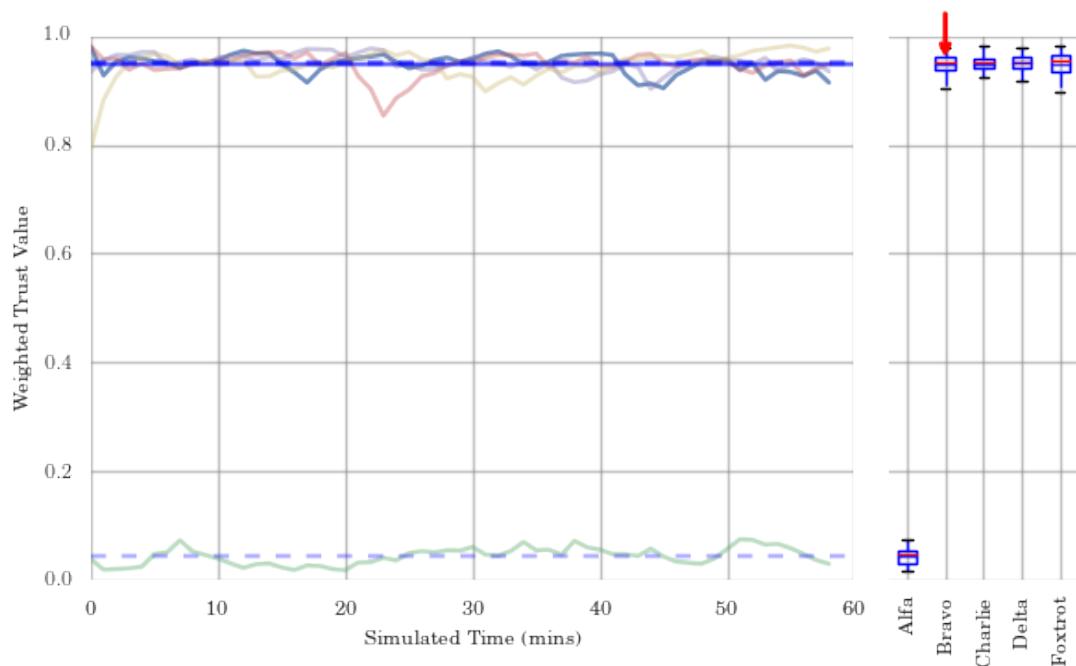
#### D.2.1 Per Node Breakdowns



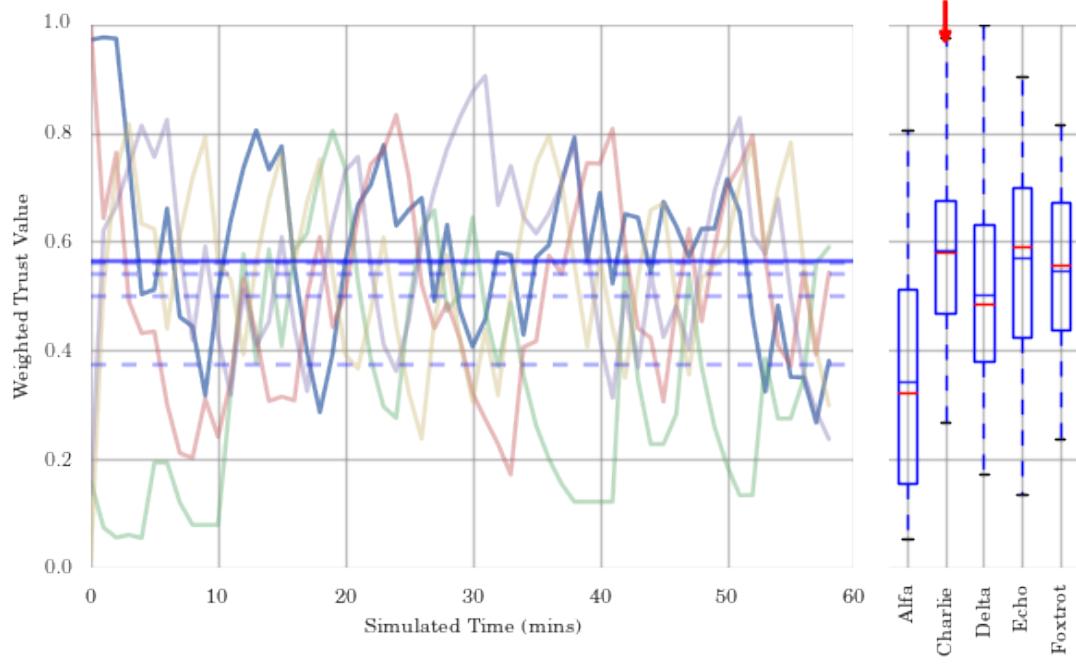
**Figure D.13:** MPC Comms Metric Trust (targeting non-malicious node)



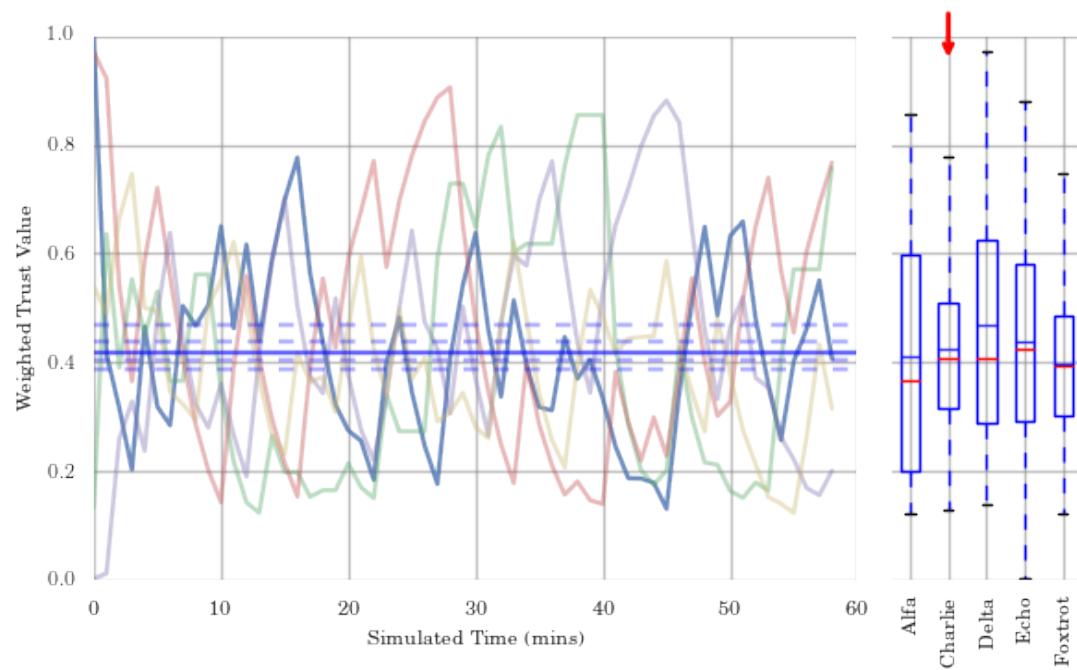
**Figure D.14:** MPC Physical Metric Trust (targeting non-malicious node)



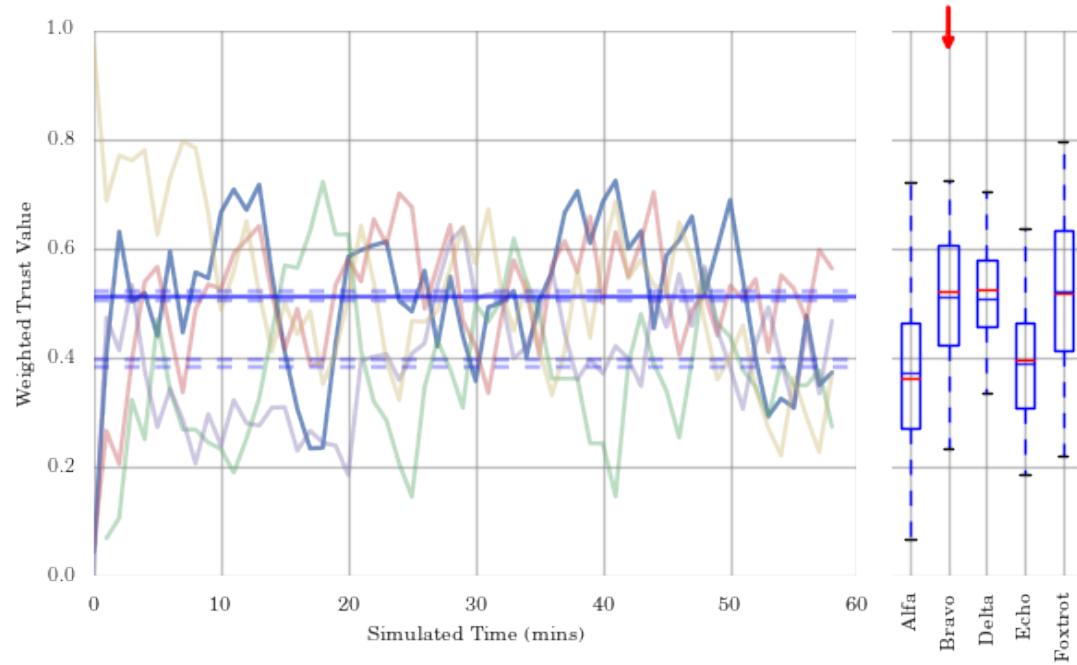
**Figure D.15:** MPC Full Metric Trust (targeting non-malicious node)



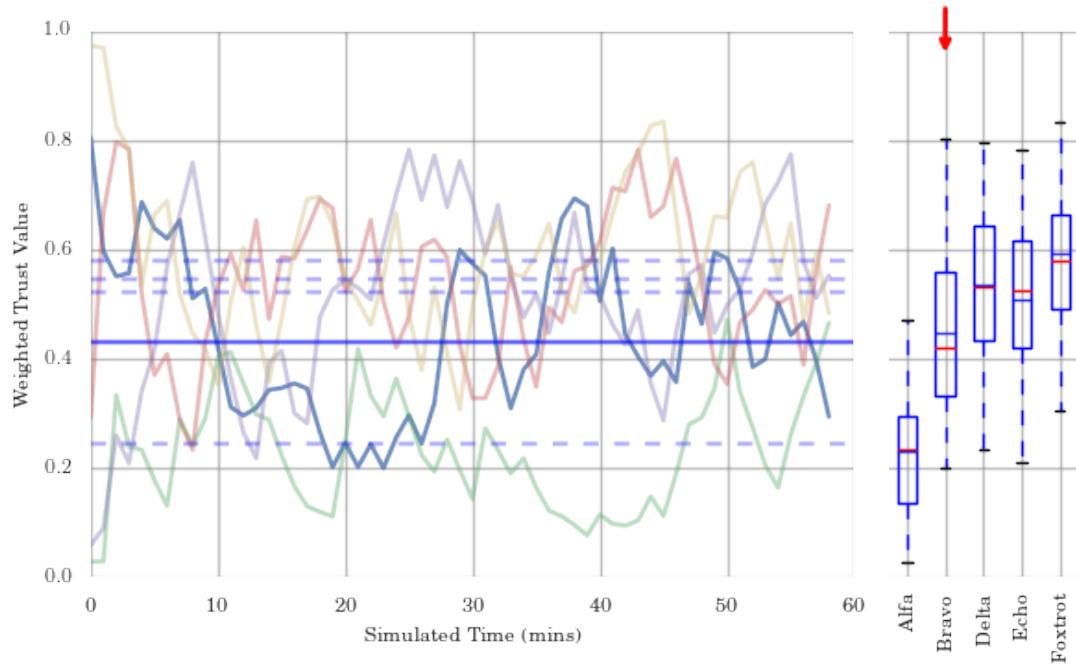
**Figure D.16:** STS Comms Metric Trust (targeting non-malicious node)



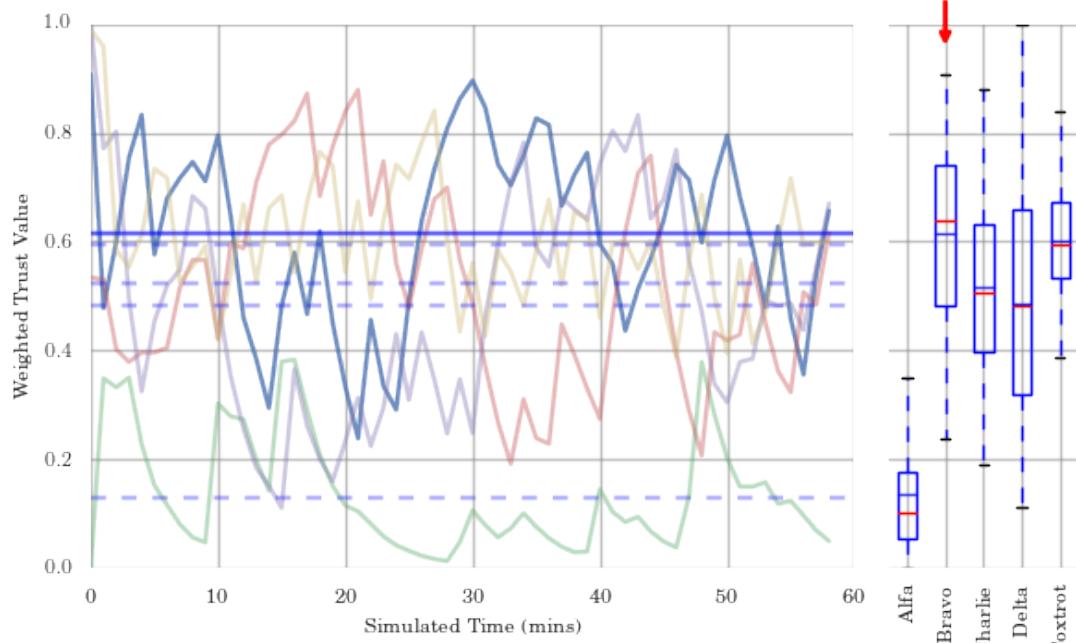
**Figure D.17:** STS Physical Metric Trust (targeting non-malicious node)



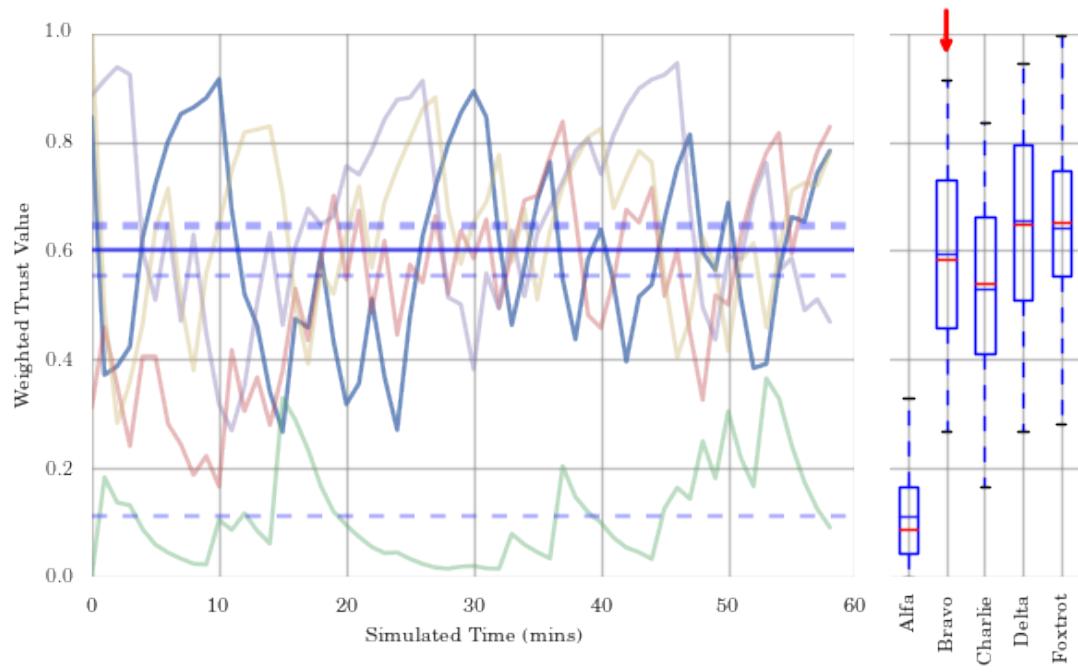
**Figure D.18:** STS Full Metric Trust (targeting non-malicious node)



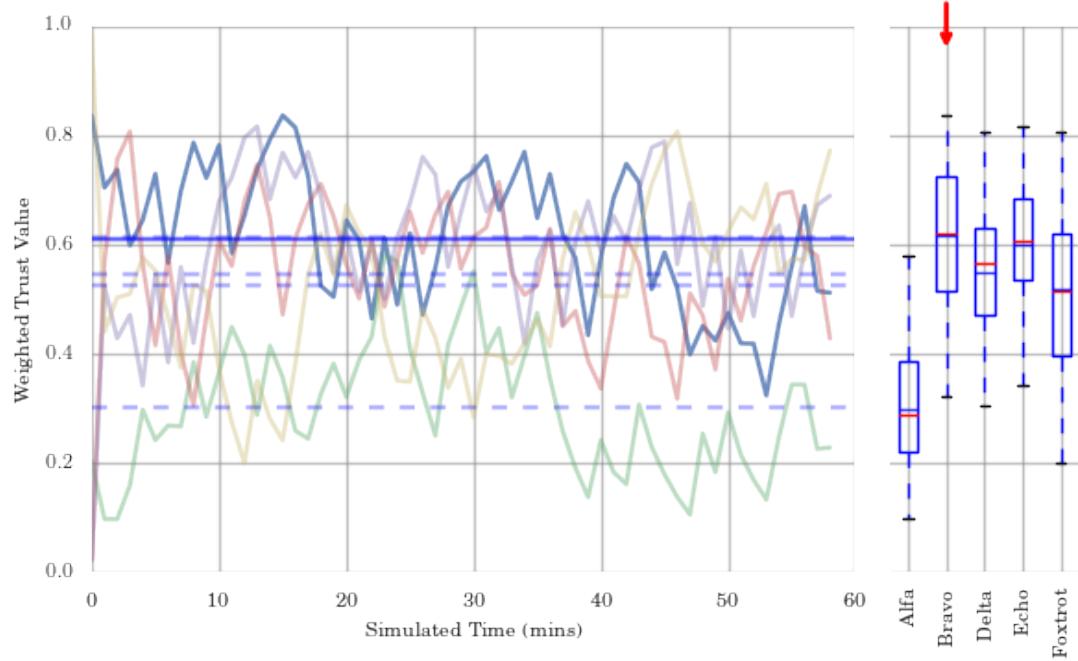
**Figure D.19:** Shadow Comms Metric Trust (targeting non-malicious node)



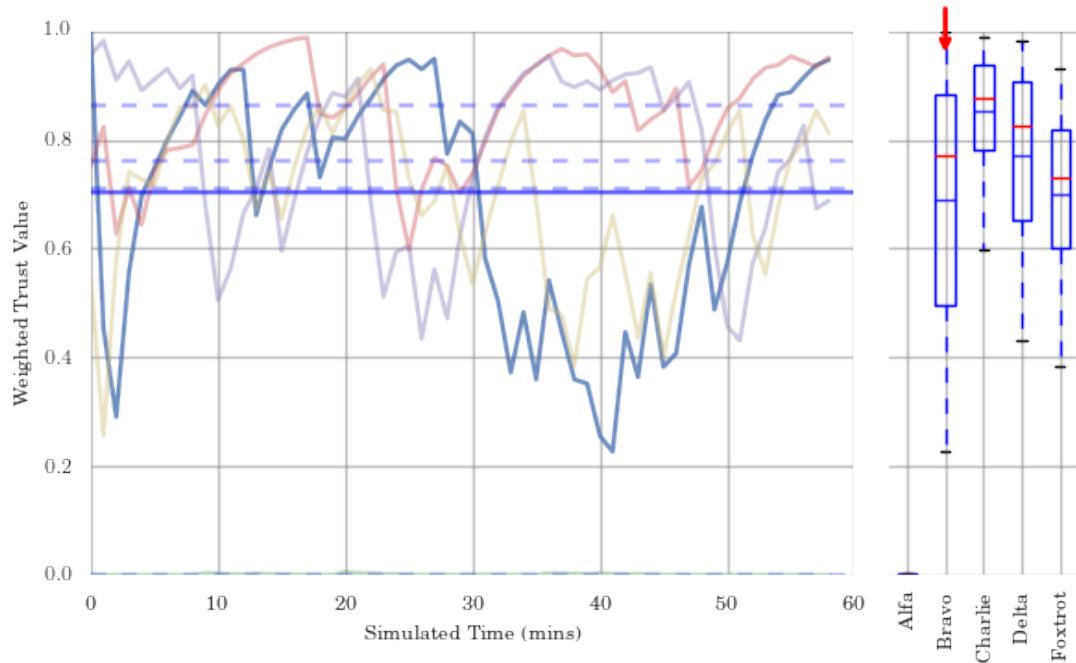
**Figure D.20:** Shadow Physical Metric Trust (targeting non-malicious node)



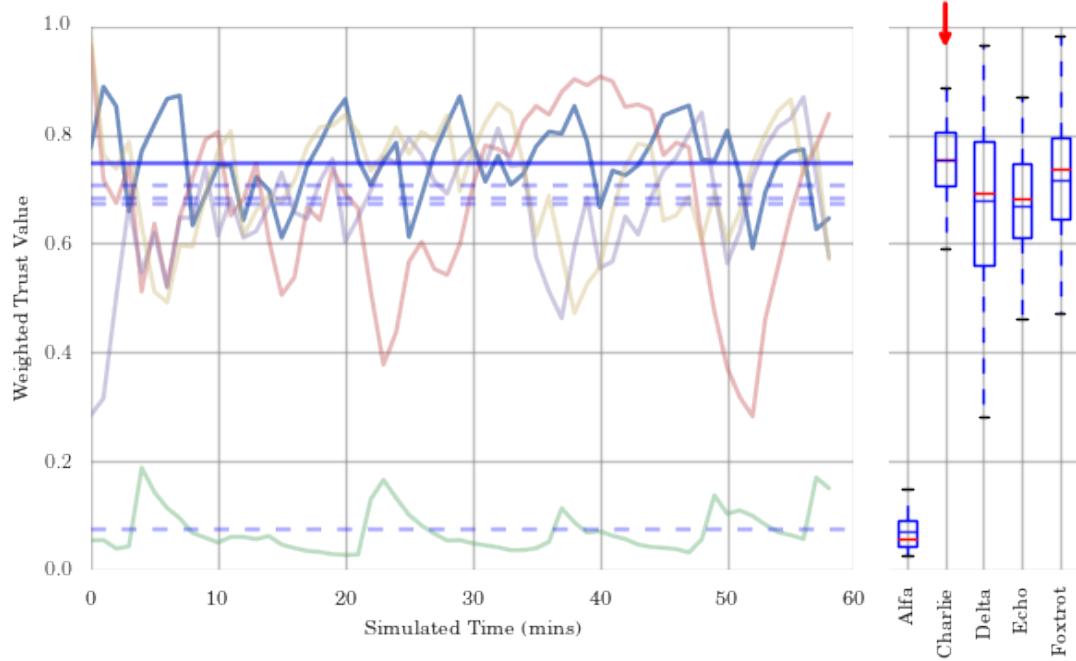
**Figure D.21:** Shadow Full Metric Trust (targeting non-malicious node)



**Figure D.22:** SlowCoach Comms Metric Trust (targeting non-malicious node)

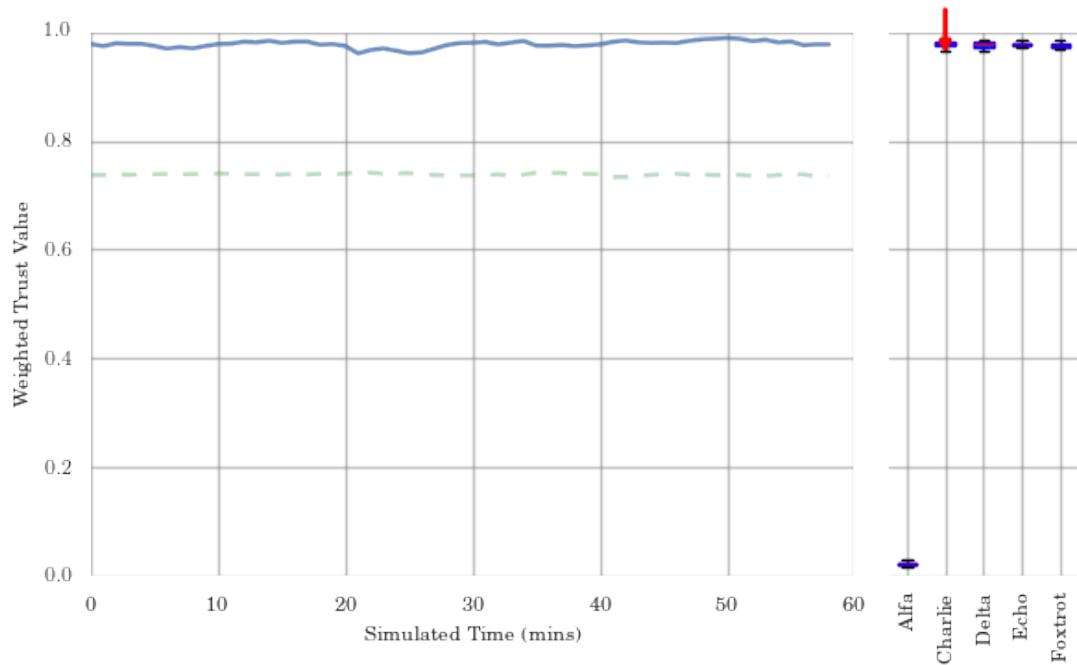


**Figure D.23:** SlowCoach Physical Metric Trust (targeting non-malicious node)

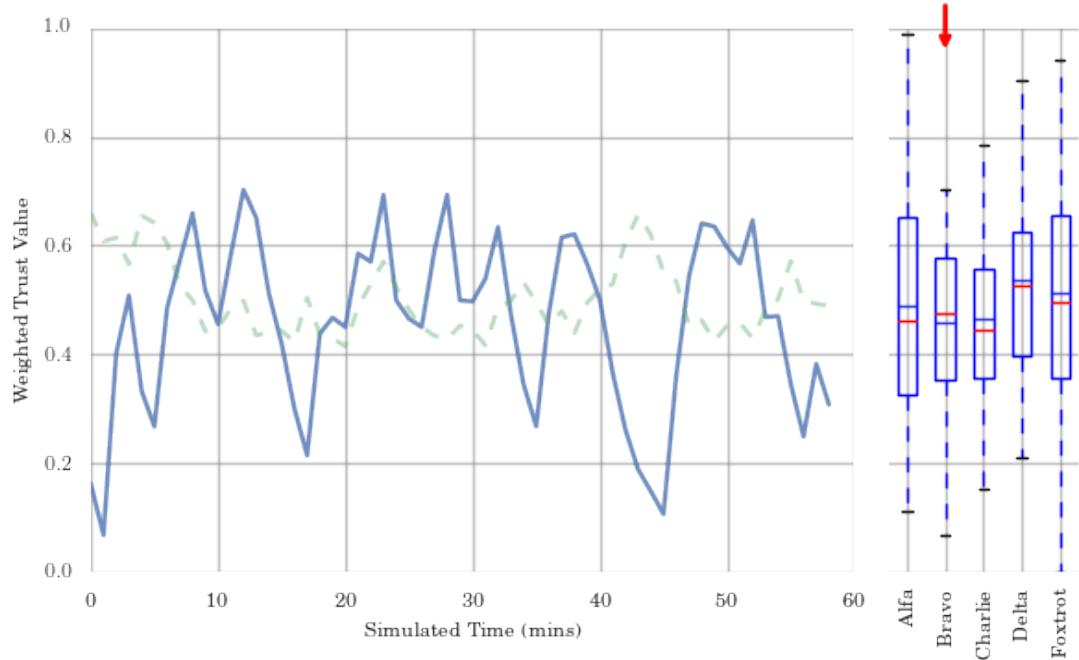


**Figure D.24:** SlowCoach Full Metric Trust (targeting non-malicious node)

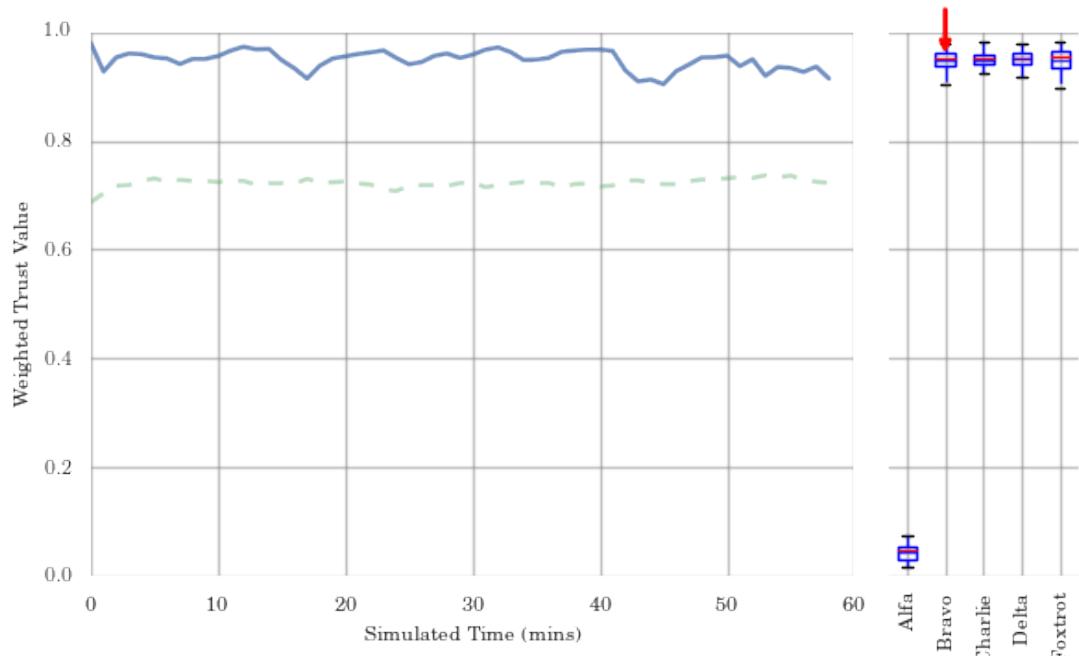
### D.2.2 Averaged across remaining cohort



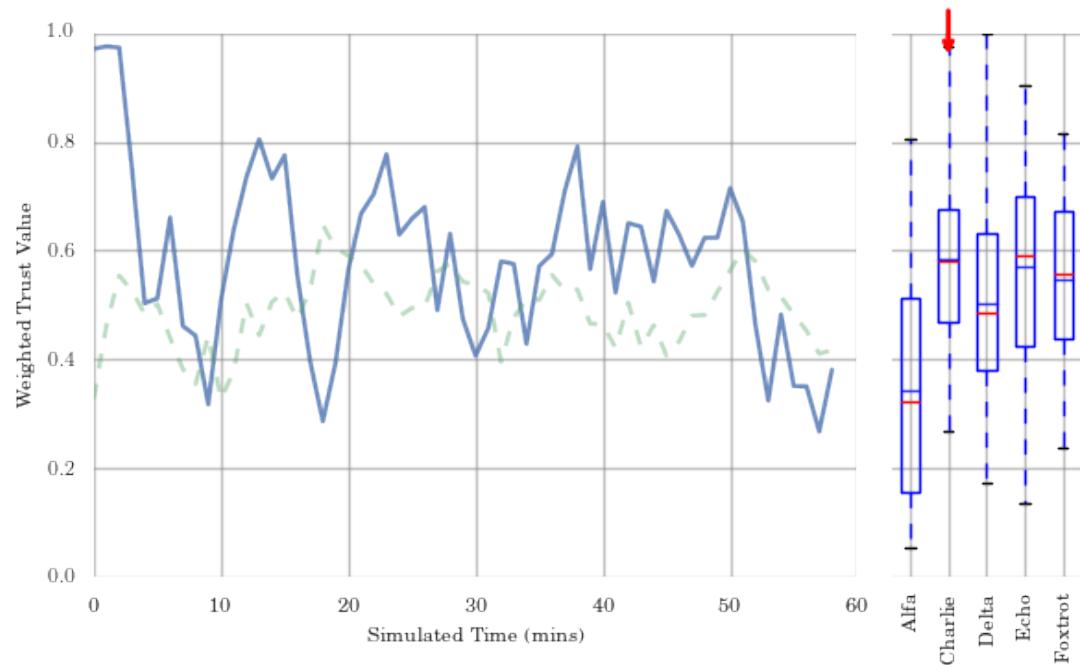
**Figure D.25:** MPC Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



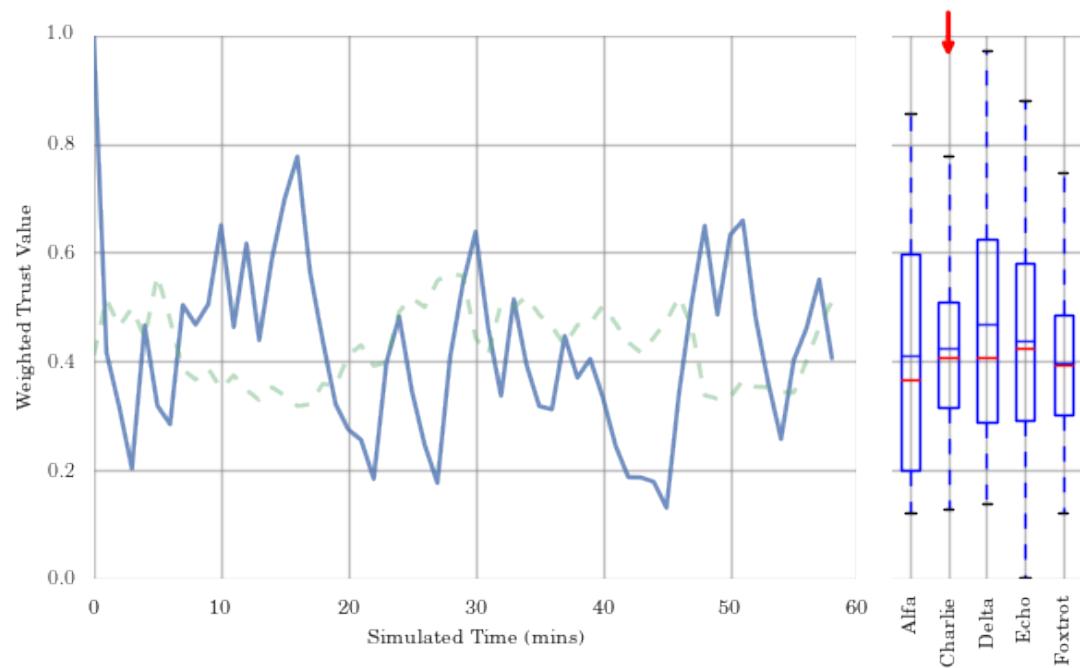
**Figure D.26:** MPC Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



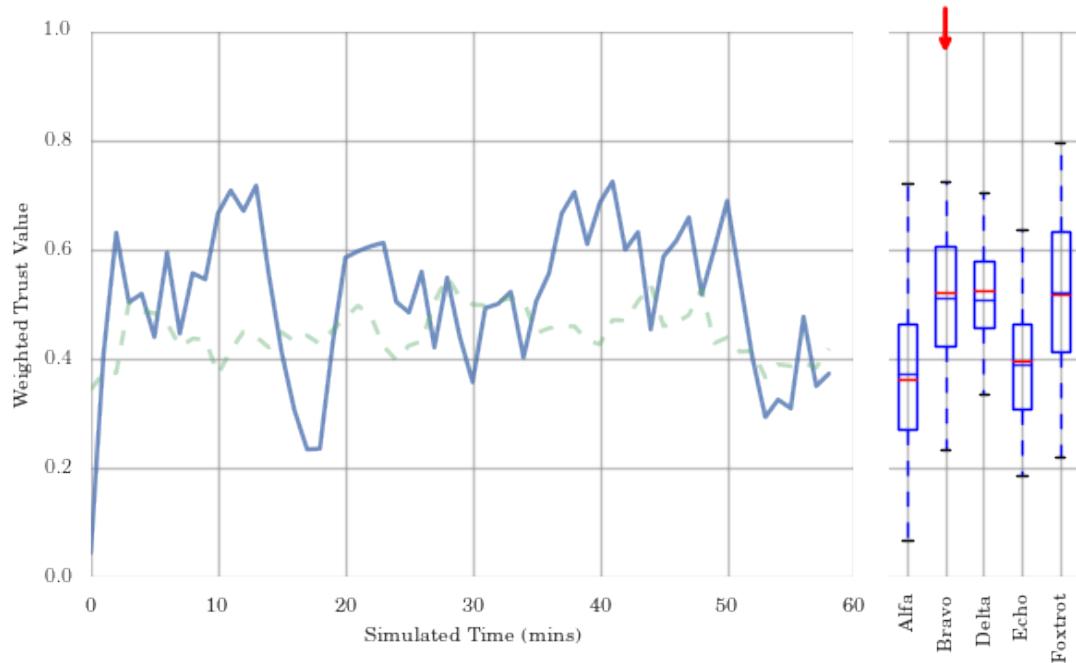
**Figure D.27:** MPC Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



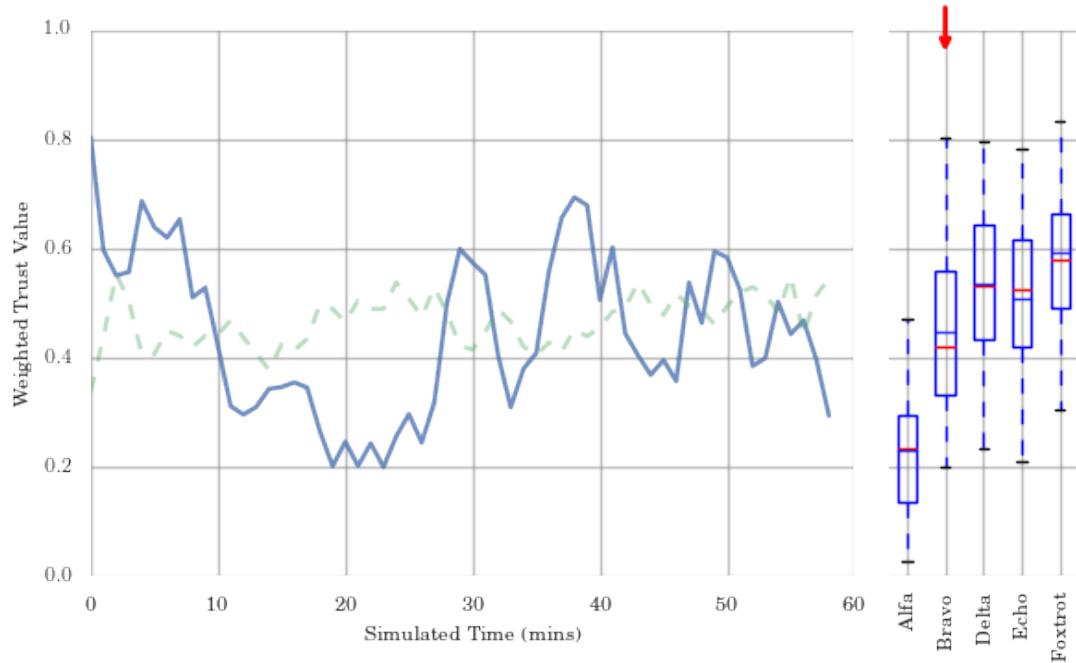
**Figure D.28:** STS Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



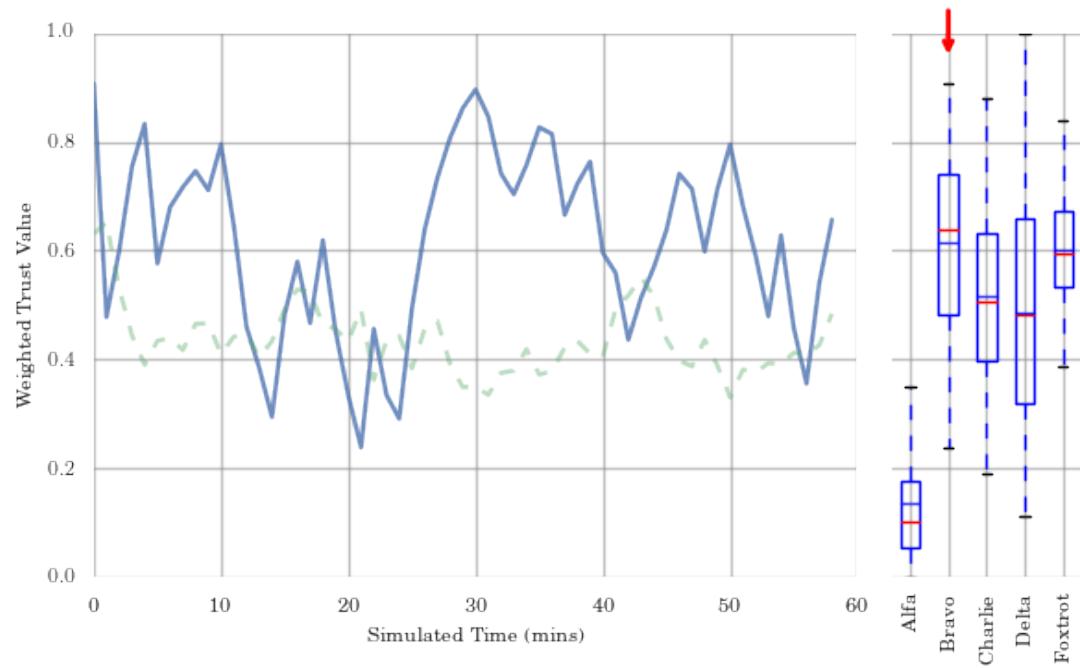
**Figure D.29:** STS Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



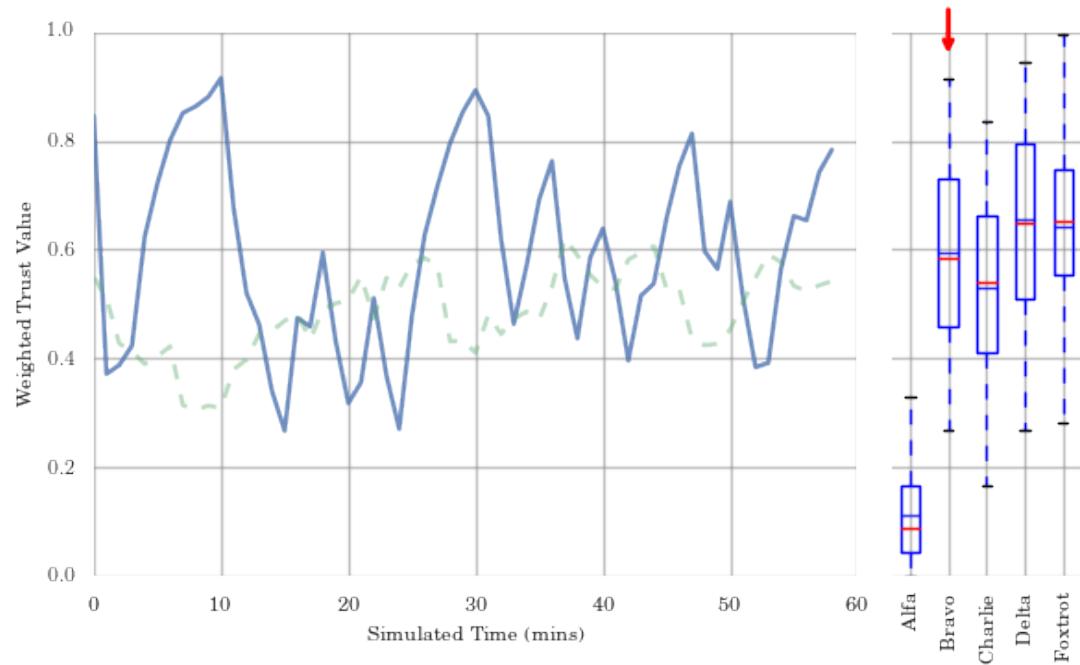
**Figure D.30:** STS Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



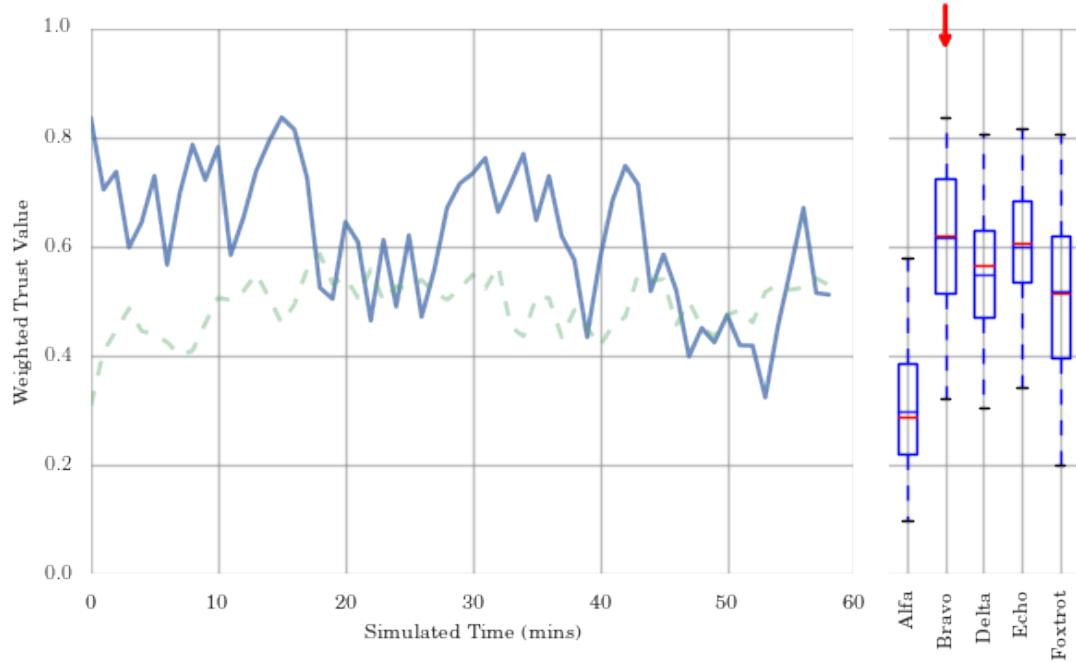
**Figure D.31:** Shadow Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



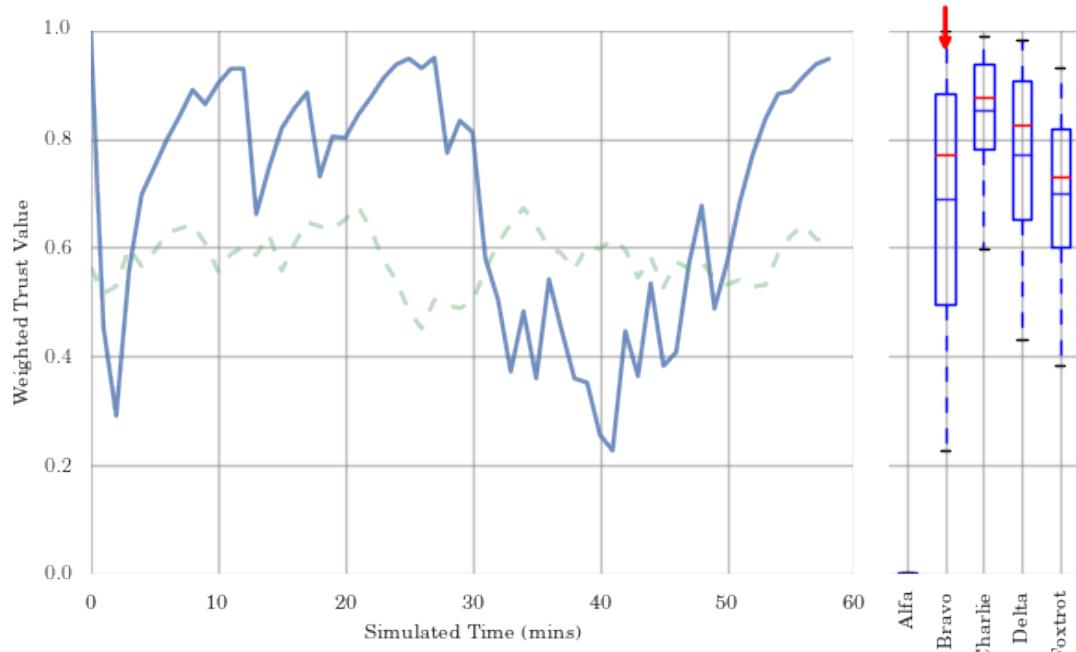
**Figure D.32:** Shadow Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



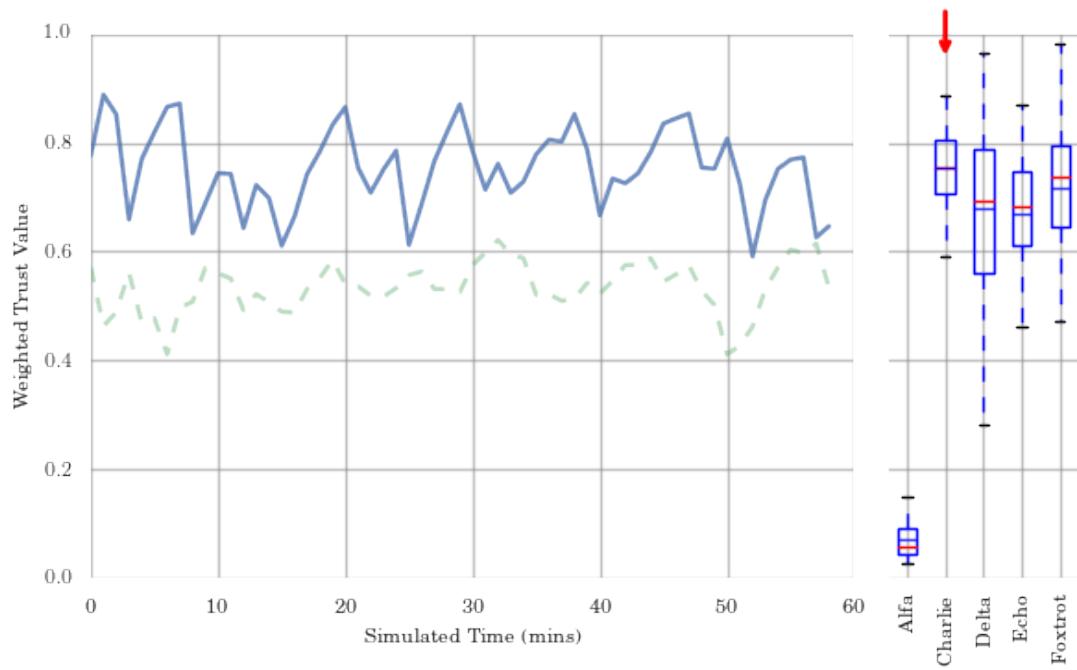
**Figure D.33:** Shadow Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



**Figure D.34:** SlowCoach Comms Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



**Figure D.35:** SlowCoach Physical Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)



**Figure D.36:** SlowCoach Full Metric Trust (targeting non-malicious node, showing mean of remaining cohort including malicious node)

# Todo list

|  |    |
|--|----|
| Expand Node Density discussion to include examples of sparse, dense, long/deep, fully connected networks . . . . .   | 2  |
| Finish Proactive Routing Protocols Table . . . . .   | 3  |
| Finish Reactive Routing Protocols Table . . . . .  | 4  |
| Write Hybrid Routing Protocols . . . . .   | 5  |
| Finish Hybrid Routing Protocols Table . . . . .  | 5  |
| more background on the operation of TTP/CA/PKI? . . . . .  | 5  |
| Trust as Assurance . . . . .   | 5  |
| Trust operation against capable attackers . . . . .  | 7  |
| Contributions . . . . .  | 7  |
| Conclusions including Layout . . . . .   | 7  |
| Fig. 2.1: Only reasonable delimitation of Trustee Operation that doesn't corrupt Mayers thesis is curring half way though Risk Taking, Outcomes and <i>maybe</i> perceived risk (as this is also affected by outcomes; may make more sense to have a separate augmented diagram) . . . . . | 10 |
| Liu and Wang do lots on this [21] as well as discussion regarding the entropic/probabilistic models of trust. This may be too much to throw in, might inject it later . . . . .  | 12 |
| Talk about trust vs untrust vs nontrust . . . . .  | 12 |
| Explore notations of transitivity and abstract trust synthesis . . . . .   | 12 |
| Redo relationship examples after Notation finished . . . . .   | 13 |
| Need to discuss how trust is established a) initially among a co-launched group, b) with a newcomer and c) with a returner (Li and Singhal [11], Liu [22], Theodorakopoulos and Baras [25]) . . . . .  | 14 |
| Expand introduction and plan the rest of the section . . . . .   | 15 |
| Possibly expand this discussion . . . . .  | 16 |
| Ref Table 2.5 there may be a case to discuss the breakdown of "Plan, Decide, Execute, Inform", possibly a nice onion-style graphic . . . . .   | 16 |
| Possibly worth looking at the Definition environment from amsthm to look after definitions like this . . . . .   | 18 |
| Need to provice a linking section to the next blocks about Design/Operational Trust  | 18 |

|   |    |
|---|----|
| No idea how to phrase this citation correctly; it's "my" work that was generated for DSTL and don't want to waste any more space backing it up; can I get away with just citing myself? . . . . . | 18 |
| Rethink using these questions at all; opens up to awkward questioning that isn't answered in the thesis . . . . .   | 18 |
| Need to check in with JP/JGF on status of JANUS. IIRC Janus dropped the whole idea of negotiating capabilities . . . . .  | 20 |
| Need to squeeze in something about Block 4 above is the focus of this work. Possibly could live in the conclusions . . . . .  | 20 |
| Needs references . . . . .  | 22 |
| Need to check security status of this source . . . . .  | 22 |
| Need to check security status of this source . . . . .  | 23 |
| Need to check security status of this source . . . . .  | 23 |
| ReDo this later . . . . .   | 24 |
| Could really do with a better / additional cite than this... . . . . .  | 24 |
| This isn't actually explained or justified in Kamvar so it may have been pulled out of his ass . . . . .  | 25 |
| Standard table . . . . .  | 25 |
| Emphasise Threat Surface discussion . . . . .   | 25 |
| Expand background detail on more frameworks . . . . .   | 27 |
| This $\rho$ bugs me; it should really be $p(O)$ based on Bayes Theorem . . . . .  | 27 |
| This makes absolutely no sense without a few diagrams . . . . .   | 27 |
| Want at least CONFIDANT and Fuzzy in here for contrast . . . . .  | 28 |
| Replace Fig. 2.4 with vector one . . . . .  | 28 |
| Fix equation links on this page after finishing grey stuff . . . . .  | 30 |
| Plot and explain the point of Whitenization (or move these back to the appendix) . . . . .  | 30 |
| Actual Conclusion of Trust Background . . . . .   | 30 |
| Introduction to Maritime . . . . .  | 31 |
| Best to discuss notation here . . . . .   | 32 |
| this might be better as a table . . . . .   | 32 |
| Need to discuss Speed of Sound Profiles . . . . .   | 33 |
| Possibly need to switch this with the Francois Garrison model which, depending on your source, is the refined version (or vice versa . . . . .  | 34 |
| check what $s$ and $w$ are in this . . . . .  | 35 |
| Vectorise and Label . . . . .   | 36 |
| expand this, justify AUVNetSim, reactive mobility, python compatibility, SimPy Etc. . . . .   | 36 |
| Summary of Akyildiz02/05 . . . . .  | 36 |
| Typical AUV missions, payloads, and available equipment . . . . .   | 36 |
| Future Applications of AUVs . . . . .   | 36 |
| it would be worth while going through this verification explicitly as an appendix .   | 41 |

|   |     |
|---|-----|
| Need to have a discussion about mission configurations at some point . . . . .  | 42  |
| redo these graphs with wider separations 1000m . . . . .  | 46  |
| Another interesting aspect is the behaviour of the Enqueued Packet lines and e2e<br>delay lines; They “Bump”; no idea why yet . . . . .   | 46  |
| This does NOT make for easy comparison between graphs as the scaling is different<br>for each mobility, but I need to think about how to fairly solve this . . . . .                  | 50  |
| Attempt to Formalise the relationship between separation, offered load, through-<br>put and delay . . . . .   | 50  |
| Double Check These Numbers Before Release . . . . .   | 50  |
| expand this section to include discussion and results of single mobility models . .   | 51  |
| this is a place holder for actual information . . . . .   | 51  |
| This should be moved back . . . . .   | 53  |
| Probably do away with this, repeated in a few other places . . . . .  | 54  |
| Look at redoing this with other mobilities (particularly distributed lawnmower) . .   | 56  |
| Explain the Minority Classifier . . . . .   | 67  |
| In the thesis, we’re concerned about a lot more than just the all mobile results . .  | 70  |
| Need to actually show physical only trust measurements . . . . .  | 80  |
| referencing the right equ in the wrong place . . . . .  | 81  |
| Duplicating C6 Metric Weighting Section . . . . .   | 82  |
| Come back to this and talk about redundancy . . . . .   | 82  |
| This isn’t right. DT doesn’t include it’s own value! . . . . .  | 85  |
| Could do with a conceptual graphic showing what these look like, although it’d<br>be messy as all hell . . . . .  | 85  |
| Could also do with a investigation into the deviation of T’s; so far most of this<br>analysis averages everything, which is almost certainly not the best approach<br>alone . . . . . | 85  |
| Answer what happens when you vary MPC power variation . . . . .   | 85  |
| Haven’t worked out a clever way of automatically generating both the basic domain<br>and alternate domain texes easily . . . . .  | 86  |
| Figure: trust bella single mobile selfish . . . . .   | 97  |
| Figure: trust bella allbut1 mobile selfish . . . . .  | 97  |
| Figure: beta trust bella static joint . . . . .   | 98  |
| Figure: beta trust bella single mobile joint . . . . .  | 98  |
| Figure: beta trust bella allbut1 mobile joint . . . . .   | 98  |
| Figure: beta trust bella all mobile joint . . . . .   | 98  |
| Figure: Indicitive Future MCM Scenario . . . . .  | 99  |
| Check Security . . . . .  | 102 |
| Check Security . . . . .  | 104 |
| don’t think classification is the right word here . . . . .   | 105 |
| eqs of sequence buffers and partial derivs . . . . .  | 106 |



# Bibliography

- [1] Jonny Milliken and David Linton. Prioritisation of citizen-centric information for disaster response. *Disasters*, pages n/a—n/a, 2015. ISSN 1467-7717. doi: 10.1111/dis.12168. URL <http://dx.doi.org/10.1111/dis.12168>.
- [2] J.W. Nicholson and A.J. Healey. Underwater Acoustic Communications and Networking: Recent Advances and Future Challenges. *Mar. Technol. Soc. J.*, 42(1):103–116, 2008. ISSN 00253324. doi: 10.4031/002533208786861263. URL [http://qub.library.ingentaconnect.com/content/mts/mts\\_j/2008/00000042/00000001/art00008](http://qub.library.ingentaconnect.com/content/mts/mts_j/2008/00000042/00000001/art00008).
- [3] S. Selvakennedy, S. Sinnappan, and Yi Shang. A biologically-inspired clustering protocol for wireless sensor networks. *Comput. Commun.*, 30(14-15):2786–2801, 2007. ISSN 01403664. doi: 10.1016/j.comcom.2007.05.010. URL <http://linkinghub.elsevier.com/retrieve/pii/S0140366407002083>.
- [4] J. Jubin and J.D. Tornow. The DARPA packet radio network protocols. *Proc. IEEE*, 75(1):21–32, 1987. ISSN 0018-9219. doi: 10.1109/PROC.1987.13702.
- [5] I. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci. Wireless sensor networks: a survey. *Comput. Networks*, 38(4):393–422, 2002. ISSN 13891286. doi: 10.1016/S1389-1286(01)00302-4. URL <http://linkinghub.elsevier.com/retrieve/pii/S1389128601003024>.
- [6] S Corson and J Macker. Mobile Ad hoc Networking (MANET): Routing Protocol Performance Issues and Evaluation Considerations. RFC 2501, RFC Editor, jan 1999.
- [7] Elizabeth M Royer. An Analysis of the Optimum Node Density for Ad hoc Mobile Networks. *IEEE Int. Conf. Commun.*, pages 857–861, 2001.
- [8] Charles E. Perkins, Pravin Bhagwat, Charles E. Perkins, and Pravin Bhagwat. Highly Dynamic Destination-Sequenced Distance-Vector Routing {(DSDV)} for Moblie Computers. *Proc. ACM Conf. Commun. Archit. Protoc. Appl.*, pages 234–244, 1994.
- [9] D. B Johnson and D. a Maltz. Dynamic source routing in ad hoc wireless networks. *Kluwer Int. Ser. Eng. Comput. Sci.*, pages 153–179, 1996.

- [10] Mehran Abolhasan, Eryk Dutkiewicz, and Tadeusz Wysocki. A review of routing protocols for mobile ad hoc networks. *Ad Hoc Networks*, 2(1):1–22, 2004. ISSN 15708705. doi: 10.1016/S1570-8705(03)00043-X. URL <http://linkinghub.elsevier.com/retrieve/pii/S157087050300043X>.
- [11] Huaizhi Li and Mukesh Singhal. Trust Management in Distributed Systems. *Computer (Long. Beach. Calif.)*, 40(2):45–53, 2007. ISSN 00189162. doi: 10.1109/MC.2007.76. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4085622>.
- [12] Sonja Buchegger and Jean-Yves Le Boudec. Performance analysis of the CONFIDANT protocol. In *Proc. 3rd ACM Int. Symp. Mob. ad hoc Netw. Comput. - MobiHoc '02*, pages 226–236. ACM Press, 2002. ISBN 1581135017. doi: 10.1145/513800.513828. URL <http://dl.acm.org/citation.cfm?id=513800.513828>.
- [13] Ji Guo, Alan Marshall, and Bosheng Zhou. A new trust management framework for detecting malicious and selfish behaviour for mobile ad hoc networks. *Proc. 10th IEEE Int. Conf. Trust. Secur. Priv. Comput. Commun. Trust. 2011, 8th IEEE Int. Conf. Embed. Softw. Syst. ICCESS 2011, 6th Int. Conf. FCST 2011*, pages 142–149, 2011. doi: 10.1109/TrustCom.2011.21. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6120813>.
- [14] Andrew Bolster and Alan Marshall. Single and Multi-metric Trust Management Frameworks for Use in Underwater Autonomous Networks. In *Trust. 2015 IEEE*, volume 1, pages 685–693, aug 2015. doi: 10.1109/Trustcom.2015.435. URL [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=7345343](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=7345343).
- [15] John D Lee and Katrina A See. Trust in automation: designing for appropriate reliance. *Hum. Factors*, 46(1):50–80, 2004. ISSN 0018-7208. doi: 10.1518/hfes.46.1.50.30392.
- [16] Tetsushi Okumura, Jeanne M. Brett, William W. Maddux, and Peter H. Kim. Cultural Differences in the Function and Meaning of Apologies. *Int. Negot.*, 16:405–425, 2011. ISSN 1382-340X. doi: 10.1163/157180611X592932.
- [17] Roger C Mayer, James H Davis, and F David Schoorman. An Integrative Model of Organizational Trust. *Acad. Manag. Rev.*, 20(3):709–734, jul 1995. ISSN 03637425. doi: 10.2307/258792. URL <http://www.jstor.org/stable/258792>.
- [18] Julian B Rotter. A new scale for the measurement of interpersonal trust1. *J. Pers.*, 35(4):651–665, 1967. ISSN 1467-6494. doi: 10.1111/j.1467-6494.1967.tb01454.x. URL <http://dx.doi.org/10.1111/j.1467-6494.1967.tb01454.x>.
- [19] Yan Lindsay Sun, Rhode Island, Z Han, and K J R Liu. Defense of trust management vulnerabilities in distributed networks. *IEEE Commun. Mag.*, 46(2):112–119, 2008. URL [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=4473092](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4473092).

- [20] D Harrison McKnight and Norman L Chervany. The meanings of trust. *Measurement*, 55455(612):86, 1996. ISSN 0277786X. doi: 10.1117/12.304574. URL <http://citeserx.ist.psu.edu/viewdoc/summary?doi=10.1.1.155.1213>.
- [21] K. J. Ray Liu and Beibei Wang. *Cognitive Radio Networking and Security: A Game-Theoretic View*. 2010. ISBN 9780521762311. doi: 10.1017/CBO9780511778773. URL [#}reader\\_{\\_}0521762316.](http://www.amazon.com/Cognitive-Radio-Networking-Security-Game-Theoretic/dp/0521762316/ref=sr_1_10?s=books&ie=UTF8&qid=1413413370&sr=1-10&keywords=cognitive+radio)
- [22] K J R Liu. Information theoretic framework of trust modeling and evaluation for ad hoc networks. *IEEE J. Sel. Areas Commun.*, 24(2):305–317, 2006. ISSN 07338716. doi: 10.1109/JSAC.2005.861389. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1589110>.
- [23] Surya Pavan, Kumar Gudla, and N Preeti. An Overview of Reputation and Trust in Multi Agent System in Disparate Environments. 5(3):498–504, 2015.
- [24] Jerome Robbins and Robert Wise. West Side Story, 1961.
- [25] George Theodorakopoulos and John S Baras. Trust evaluation in ad-hoc networks. *Proc. 2004 ACM Work. Wirel. Secur. WiSe 04*, (October):1, 2004. doi: 10.1145/1023646.1023648. URL <http://portal.acm.org/citation.cfm?doid=1023646.1023648>.
- [26] Nomy Arpaly. *Unprincipled virtue : an inquiry into moral agency*. Oxford University Press, Oxford; New York, 2003. ISBN 0195152042 9780195152043.
- [27] Allen Hunter, New Social, Movements Author, Allen Hunter Review, Allen Hunter Source, and Springer Stable Url. Post-Marxism and the New Social Movements. *Theory Soc.*, 17(6):885–900, 2016. ISSN 03042421, 15737853. URL <http://www.jstor.org/stable/657793>.
- [28] Daniel Halberstam and Roderick M Hills. State Autonomy in Germany and the United States. *Ann. Am. Acad. Pol. Sci.*, 574:173–184, 2001. ISSN 00027162. URL <http://www.jstor.org/stable/1049063>.
- [29] Wolfram Richter. Is Europe ready to give up national autonomy for the sake of the euro?, 2012. URL <http://www.theguardian.com/commentisfree/2012/jul/15/europe-economists-letters-national-autonomy>.
- [30] R. Alami, R. Chatila, S. Fleury, M. Ghallab, and F. Ingrand. An Architecture for Autonomy. *Int. J. Rob. Res.*, 17:315–337, 1998. ISSN 0278-3649. doi: 10.1177/027836499801700402.

- [31] George A Bekey. *Autonomous robots : from biological inspiration to implementation and control.* 2005. ISBN 0262025787. URL <http://books.google.com/books?hl=en&lr=&id=3xwbia2DpmoC&oi=fnd&pg=PR13&dq=Autonomous+Robots+From+Biological+Inspiration+to+Implementation+and+Control&ots=WxngXPbihr&sig=7G8VA4GRaU0wcOsAFbfPi5uAK18>.
- [32] Stan Franklin and Art Graesser. Is it an Agent, or just a Program?: A Taxonomy for Autonomous Agents. In *Intell. agents III agent Theor. Archit. Lang.*, pages 21–35. Springer, 1997.
- [33] H. M. Huang. Autonomy Levels for Unmanned Systems ( ALFUS ) Framework Volume I : Terminology Unmanned Systems Working Group Participants 1 National Institute of Standards and Technology. *Framework*, I(September):29, 2004.
- [34] R.R. Murphy. *Introduction to AI robotics*, volume 108. 2000. ISBN 0262133830. doi: 10.1111/j.1464-410X.2011.10513.x. URL <http://www.ncbi.nlm.nih.gov/pubmed/21917105>.
- [35] Stuart Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach, 3rd edition.* 2009. ISBN 0136042597. doi: 10.1017/S0269888900007724. URL [http://portal.acm.org/citation.cfm?id=1671238&coll=DL&dl=GUIDE&CFID=190864501&CFTOKEN=29051579\\$&delimter" 026E30F\\$npapers2://publication/uuid/4B787E16-89F6-4FF7-A5E5-E59F3CFEFE88](http://portal.acm.org/citation.cfm?id=1671238&coll=DL&dl=GUIDE&CFID=190864501&CFTOKEN=29051579$&delimter).
- [36] Sebastian Thrun. Toward a Framework for Human-Robot Interaction. *Human-Computer Interact.*, 19:9–24, 2004. ISSN 0737-0024. doi: 10.1207/s15327051hci1901&2\_2.
- [37] Michael Wooldridge and Nicholas R. Jennings. Intelligent agents: theory and practice, 1995. ISSN 0269-8889.
- [38] Thomas B. Sheridan and William L. Verplank. Human and Computer Control of Undersea Teleoperators. *ManMachine Syst. Lab Dep. Mech. Eng. MIT Grant N0001477C0256*, page 343, 1978.
- [39] M R Endsley and D B Kaber. Level of automation effects on performance, situation awareness and workload in a dynamic control task. *Ergonomics*, 42(3):462–492, 1999. ISSN 0014-0139. doi: 10.1080/001401399185595.
- [40] NATO Standardization Office. STANAG 4586 STANDARD INTERFACES OF UAV CONTROL SYSTEM (UCS) FOR NATO UAV INTEROPERABILITY Ed: 3. Technical report, NATO, Brussels, Belgium, 2012. URL <http://nso.nato.int/nso/zPublic/stanags/current/4586eed03.pdf>.
- [41] Mary L. Cummings, Sylvain Bruni, and Paul J. Mitchell. Chapter 2<BR> Human Supervisory Control Challenges in Network-Centric Operations, 2010. ISSN 1557234X.

- [42] Neya Systems LLC. The JAUS Toolset. URL <http://jaustoolset.org/>.
- [43] American Society of Testing and Materials. ASTM F2500 - 07 Standard Practice for Unmanned Aircraft System (UAS) Visual Range Flight Operations. Technical report, 2007. URL <http://www.astm.org/Standards/F2500.htm>.
- [44] American Society of Testing and Materials. ASTM F2541-06 Standard Guide for Unmanned Undersea Vehicles (UUV) Autonomy and Control. Technical report, 2006. URL <http://www.astm.org/Standards/F2541.htm>.
- [45] Jessie Y. C. Chen, Michael J. Barnes, and Michelle Harper-Sciarini. Supervisory Control of Multiple Robots: Human-Performance Issues and User-Interface Design. *IEEE Trans. Syst. Man, Cybern. Part C (Applications Rev.)*, 41(4):435–454, 2011. ISSN 1094-6977.
- [46] Aaron Mehta. Political, Financial Threads Underscore German Euro Hawk Saga. *Def. News*, jun 2013.
- [47] Nick Johnson, Pedro Patron, and David Lane. The importance of trust between operator and AUV: Crossing the human/computer language barrier. *Ocean. 2007 - Eur.*, pages 1–6, jun 2007. doi: 10.1109/OCEANSE.2007.4302408. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4302408>.
- [48] Nicholas A. R. Johnson and David M. Lane. Narrative monologue as a first step towards advanced mission debrief for AUV operator situational awareness. *2011 15th Int. Conf. Adv. Robot.*, pages 241–246, 2011. doi: 10.1109/ICAR.2011.6088618.
- [49] Sepandar D. Kamvar, Mario T. Schlosser, and Hector Garcia-Molina. The Eigen-trust algorithm for reputation management in P2P networks. *12th Int. Conf. World Wide Web (WWW)*, page 640, 2003. ISSN 1581136803. doi: 10.1145/775240.775242. URL <http://portal.acm.org/citation.cfm?doid=775152.775242>.
- [50] Charikleia Zouridaki, Brian L Mark, Marek Hejmo, and Roshan K Thomas. A quantitative trust establishment framework for reliable data packet delivery in MANETs. *Proc. 3rd ACM Work. Secur. ad hoc Sens. networks*, pages 1–10, 2005. ISSN 0926227X. doi: 10.1145/1102219.1102222.
- [51] Jie Li, Ruidong Li, Jien Kato, Jie Li, Peng Liu, and Hsiao-Hwa Chen. Future Trust Management Framework for Mobile Ad Hoc Networks. *IEEE Commun. Mag.*, 46(4):108–114, apr 2007. ISSN 01636804. doi: 10.1109/MCOM.2008.4481349. URL [http://ieeexplore.ieee.org/xpls/abs{\\\_}all.jsp?arnumber=4212452http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4481349](http://ieeexplore.ieee.org/xpls/abs{\_}all.jsp?arnumber=4212452http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4481349).
- [52] Jin-hee Cho, Ananthram Swami, and Ing-ray Chen. A survey on trust management for mobile ad hoc networks. *Commun. Surv. & Tutorials*, 13(4):562–583, 2011. URL [http://ieeexplore.ieee.org/xpls/abs{\\\_}all.jsp?arnumber=5604602](http://ieeexplore.ieee.org/xpls/abs{\_}all.jsp?arnumber=5604602).

- [53] MEG E G Moe, BE E Helvik, and SJ J Knapskog. TSR: Trust-based secure MANET routing using HMMs. ... *symposium QoS Secur.* ..., pages 83–90, 2008. URL <http://dl.acm.org/citation.cfm?id=1454602>.
- [54] Junhai Luo, Xue Liu, Yi Zhang, Danxia Ye, and Zhong Xu. Fuzzy trust recommendation based on collaborative filtering for mobile ad-hoc networks. *2008 33rd IEEE Conf. Local Comput. Networks*, pages 305–311, 2008. doi: 10.1109/LCN.2008.4664184. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4664184>.
- [55] Liang Hong Liang Hong, Wu Chen Wu Chen, Li Gao Li Gao, Guoqing Zhang Guoqing Zhang, and Cai Fu Cai Fu. Grey theory based reputation system for secure neighbor discovery in wireless ad hoc networks. *Futur. Comput. Commun. (ICFCC), 2010 2nd Int. Conf.,* 2, 2010. doi: 10.1109/ICFCC.2010.5497609.
- [56] Jim Partan, Jim Kurose, and Brian Neil Levine. A survey of practical issues in underwater networks. *Proc. 1st ACM Int. Work. Underw. networks WUWNet 06*, 11(4):17, 2006. ISSN 15591662. doi: 10.1145/1161039.1161045. URL <http://portal.acm.org/citation.cfm?doid=1161039.1161045>.
- [57] Milica Stojanovic. On the relationship between capacity and distance in an underwater acoustic communication channel, 2007. ISSN 15591662. URL <http://www.mit.edu/~millitsa/resources/pdfs/bwdx.pdf>.
- [58] Xavier Lurton. *An introduction to underwater acoustics: principles and applications*. Springer Praxis Books. Springer Berlin Heidelberg, 2002. ISBN 9783540784807. URL <https://books.google.fr/books?id=PFXgLQAACAAJ>.
- [59] Kenneth V. Mackenzie. Nineterm equation for sound speed in the oceans. *J. Acoust. Soc. Am.*, 70(3):807, 1981. ISSN 00014966. doi: 10.1121/1.386920.
- [60] R F W Coates. *Underwater acoustic systems*. A Halstead Press book. John Wiley & Sons Canada, Limited, 1989. ISBN 9780470215449. URL <https://books.google.co.uk/books?id=0qUeAQAAIAAJ>.
- [61] Chiara Petrioli and Roberto Petroccia. SUNSET: Simulation, emulation and real-life testing of underwater wireless sensor networks. *Proc. IEEE UComms 2012*, 2012. URL [http://reti.dsi.uniroma1.it/UWSN{\\_}Group/publications/pdf/2012/sunset.pdf](http://reti.dsi.uniroma1.it/UWSN{_}Group/publications/pdf/2012/sunset.pdf).
- [62] Anuj Sehgal, Iyad Tumar, and Jürgen Schönwälder. AquaTools: An Underwater Acoustic Networking Simulation Toolkit. *IEEE, Ocean. Sydney*, ..., 2010. URL [http://www.researchgate.net/publication/233835116{\\_}AquaTools{\\_}{\\_}An{\\_}Underwater{\\_}Acoustic{\\_}Networking{\\_}Simulation{\\_}Toolkit\\_file/32bfe510a858c844b7.pdf](http://www.researchgate.net/publication/233835116{_}AquaTools{_}{_}An{_}Underwater{_}Acoustic{_}Networking{_}Simulation{_}Toolkit_file/32bfe510a858c844b7.pdf).

- [63] Xinjie Chang. Network Simulations with OPNET. In *Proc. 31st Conf. Winter Simul. Simulation—a Bridg. to Futur. - Vol. 1*, WSC '99, pages 307–314, New York, NY, USA, 1999. ACM. ISBN 0-7803-5780-9. doi: 10.1145/324138.324232. URL <http://doi.acm.org/10.1145/324138.324232>.
- [64] Andrea Caiti. Cooperative distributed behaviours of an AUV network for asset protection with communication constraints. *Ocean. 2011 IEEE-Spain*, 2011. URL [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=6003463](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6003463).
- [65] Klaus Müller and Tony Vignaux. SimPy: Simulating Systems in Python. *ON-Lamp.com Python DevCenter*, feb 2003. URL <http://www.onlamp.com/pub/a/python/2003/02/27/simpy.html?page=2>.
- [66] Josep Miquel and Jornet Montana. AUVNetSim: A Simulator for Underwater Acoustic Networks. *Program*, pages 1–13, 2008. URL <http://users.ece.gatech.edu/jmj3/publications/aувnetsim.pdf>.
- [67] Andrej Stefanov and Milica Stojanovic. Design and performance analysis of underwater acoustic networks. *IEEE J. Sel. Areas Commun.*, 29(10):2012–2021, 2011. ISSN 07338716. doi: 10.1109/JSAC.2011.111211.
- [68] Kaixin Xu, Mario Gerla, Sang Bae, and Hoc Networks. Effectiveness of RTS / CTS Handshake in IEEE. . . , 2002. *Globecom'02. Ieee*, 56:1–14, 2002. ISSN 15708705. doi: 10.1049/el. URL [http://ieeexplore.ieee.org/xpls/abs\\_all.jsp?arnumber=1188044](http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1188044).
- [69] Jared Cordasco and Susanne Wetzel. Cryptographic Versus Trust-based Methods for MANET Routing Security. *Electron. Notes Theor. Comput. Sci.*, 197(2):131–140, 2008. ISSN 15710661. doi: 10.1016/j.entcs.2007.12.022.
- [70] Son-Cheol Yu and Tamaki Ura. A System of Multi-AUV Interlinked With a Smart Cable For Autonomous Inspection of Underwater Structures. *Int. J. Offshore Polar Eng.*, 14(04), 2004.
- [71] K Asakawa, J Kojima, Y Kato, S Matsumoto, N Kato, T Asai, and T Iso. Design concept and experimental results of the autonomous underwater vehicle {AQUA EXPLORER} 2 for the inspection of underwater cables. *Adv. Robot.*, 16(1):27–42, jan 2002. ISSN 01691864. doi: 10.1163/156855302317413727. URL <http://dx.doi.org/10.1163/156855302317413727>.
- [72] Brian Morr. All Quiet on the AUV Front. *Underw. Mag.*, (February):1–5, 2003.
- [73] a. Matos, N. Cruz, a. Martins, and F. Lobo Pereira. Development and implementation of a low-cost LBL navigation system\nfor an AUV. *Ocean. '99. MTS/IEEE. Rid. Crest into 21st Century. Conf. Exhib. Conf. Proc. (IEEE Cat. No.99CH37008)*, 2:774–779, 1999. ISSN 01977385. doi: 10.1109/OCEANS.1999.804906.

- [74] Jeff Snyder. Doppler Velocity Log (DVL) navigation for observation-class ROVs. *MTS/IEEE Seattle, Ocean. 2010, (Dvl)*:1–9, 2010. ISSN 00933651. doi: 10.1109/OCEANS.2010.5664561.
- [75] Bjorn Jalving, Kenneth Gade, Ove Kent Hagen, and Karstein Vestgård. A toolbox of aiding techniques for the HUGIN AUV integrated inertial navigation system. *Model. Identif. Control*, 25(3):173–190, 2004. ISSN 03327353. doi: 10.1109/OCEANS.2003.178505.
- [76] Xixiang Liu, Xiaosu Xu, Yiting Liu, and Lihui Wang. Kalman filter for cross-noise in the integration of SINS and DVL. *Math. Probl. Eng.*, 2014(Dvl), 2014. ISSN 15635147. doi: 10.1155/2014/260209.
- [77] Stefan B Williams, Paul Newman, Gamini Dissanayake, and Hugh Durrant-Whyte. Autonomous underwater simultaneous localisation and map building. *Robot. Autom. 2000. Proceedings. ICRA '00. IEEE Int. Conf.*, 2:1793–1798, 2000. ISSN 1050-4729. doi: 10.1109/ROBOT.2000.844855.
- [78] Ji Guo. Trust and Misbehaviour Detection Strategies for Mobile Ad hoc Networks. 2012.
- [79] Rob McEwen and Knut Streitlien. Modeling and control of a variable-length auv. *Proc 12th UUST*, pages 1–42, 2006. URL <http://www.mbari.org/staff/rob/uustrep.pdf>.
- [80] Andrew Bolster. Analysis of Trust Interfaces in Autonomous and Semi-Autonomous Collaborative MHPC Operations. Technical report, The Technical Cooperation Program, 2014.
- [81] R. B. Dean and W. J. Dixon. Simplified Statistics for Small Numbers of Observations. *Anal. Chem.*, 23(4):636–638, 1951. ISSN 0003-2700. doi: 10.1021/ac60052a025. URL <http://pubs.acs.org/doi/abs/10.1021/ac60052a025>.
- [82] Andrew Bolster and Alan Marshall. A Multi-Vector Trust Framework for Autonomous Systems. In *2014 AAAI Spring Symp. Ser.*, pages 17–19, Stanford, CA, 2014. ISBN 9781577356448. URL <http://www.aaai.org/ocs/index.php/SSS/SSS14/paper/viewFile/7697/7724>.
- [83] L Breiman. Random forests. *Mach. Learn.*, pages 5–32, 2001. ISSN 0885-6125. doi: 10.1023/A:1010933404324. URL <http://link.springer.com/article/10.1023/A:1010933404324>.
- [84] Sifeng Liu and Yi Lin. *Grey System Theory and Application*. Number 1. Springer-Verlag Berlin Heidelberg, 2011. ISBN 978-1-61284-490-9. doi: 10.1109/GSIS.2011.6044018. URL [http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6044018\\$&delimiter](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6044018$&delimiter) "026E30F\$nhttp://www.springer.com/physics/complexity/book/978-3-642-16157-5.

- [85] Fengchao Zuo. Determining Method for Grey Relational Distinguished Coefficient. *SIGICE Bull.*, 20(3):22–28, jan 1995. ISSN 0893-2875. doi: 10.1145/202081.202086. URL <http://doi.acm.org/10.1145/202081.202086>.