

**“Trust in Autonomy is to DoD what safety is to FAA”**



Air  
Land  
Sea  
Space  
Cyberspace

Innovation. In all domains.

# Creating Trust in Autonomous Systems – the Trust V

Dr. Dan Zwillinger  
Gari Palmer  
Anne Selwyn

Safe & Secure Systems &  
Software Symposium

12 June 2014

# Trust In Autonomous Systems – Overview

## What it is

- Develop and prototype methods for creating both Systems Trust and Operational Trust via the “Trust V”; enhancing system acquisition and use.
- Develop and prototype methods for trust measurement and assessment

## Why it’s important

- As Raytheon’s systems move from automated to autonomous, successful deployment and adoption will require both Systems Trust and Operational Trust.

## Key Technologies

- Raytheon developed Trust techniques
  - “Chatter” and “What Agent”
  - “Future Risk Prediction”
- Trust Calibration

## Key 2014 Deliverables

- Design, develop and add trust to 2+ Raytheon programs
- MOEs for “relative” trust calculation

### Notational User Interfaces

#### Chatter

#### Chatter Settings

| Information by Time or Event        |  | Unusual Values  | Save Data  |
|-------------------------------------|--|---|--|
| <b>Time Summary</b><br>Not Reported | Threat change<br>Down 2 levels<br><b>Frequency</b><br>Every instance | <input type="checkbox"/> Not shown<br><input checked="" type="checkbox"/> 3 $\sigma$ deviations<br>(after 20 samples) | <input type="checkbox"/> No<br><input checked="" type="checkbox"/> Yes<br>Data.txt |
|                                     |  | Reset data  |  |

#### What Agent

#### What Agent

Off

**Past – What did happen?** (replay of previously collected data)

| Information by Time or Event        |  | Data from          |
|-------------------------------------|--|--------------------|
| <b>Time Summary</b><br>Not Reported | Threat change<br>Down 2 levels<br><b>Frequency</b><br>Every instance | Data.txt<br>Browse |

**Future – What might happen?**

| Information by Event           | Save Data  |
|--------------------------------|--|
| Threat change<br>Down 2 levels | <input type="checkbox"/> No<br><input checked="" type="checkbox"/> Yes<br>Future.txt |

#### Future Risk Prediction

#### Future Risk Prediction

Risk unchanged for next 5 minutes

Assessment confidence 89%

Show assessment details

Need an affordable way to create Trust for autonomous systems

# Trust – What / Why

## What is Trust?

“The attitude that an agent will help achieve an individual’s goals in a situation characterized by uncertainty and vulnerability.”

## Why do we care?

“For commanders and operators in particular, these challenges can collectively be characterized as a lack of trust that the autonomous functions of a given system will operate as intended in all situations.”

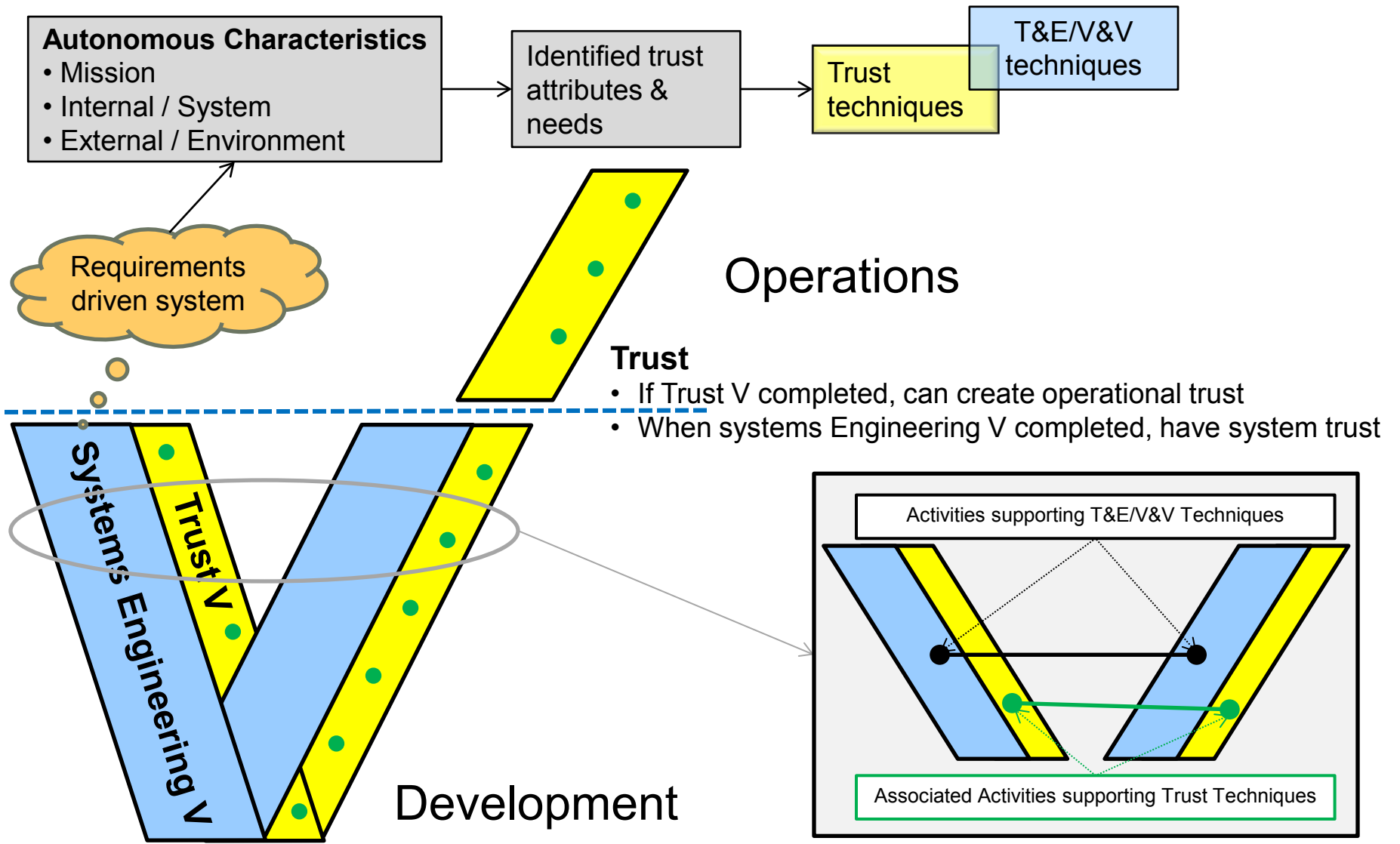
[Defense Science Board – The Role of Autonomy in DoD Systems](#) – emphasis in the original

## Trust goals

- *must* deliver value added capabilities
- *must* address the needs for different autonomy types
  - ... does system have a learning capability?
  - ... does system work in unstructured environments?
  - ...
- *should* work within existing development process
- *should* support creation of modular components

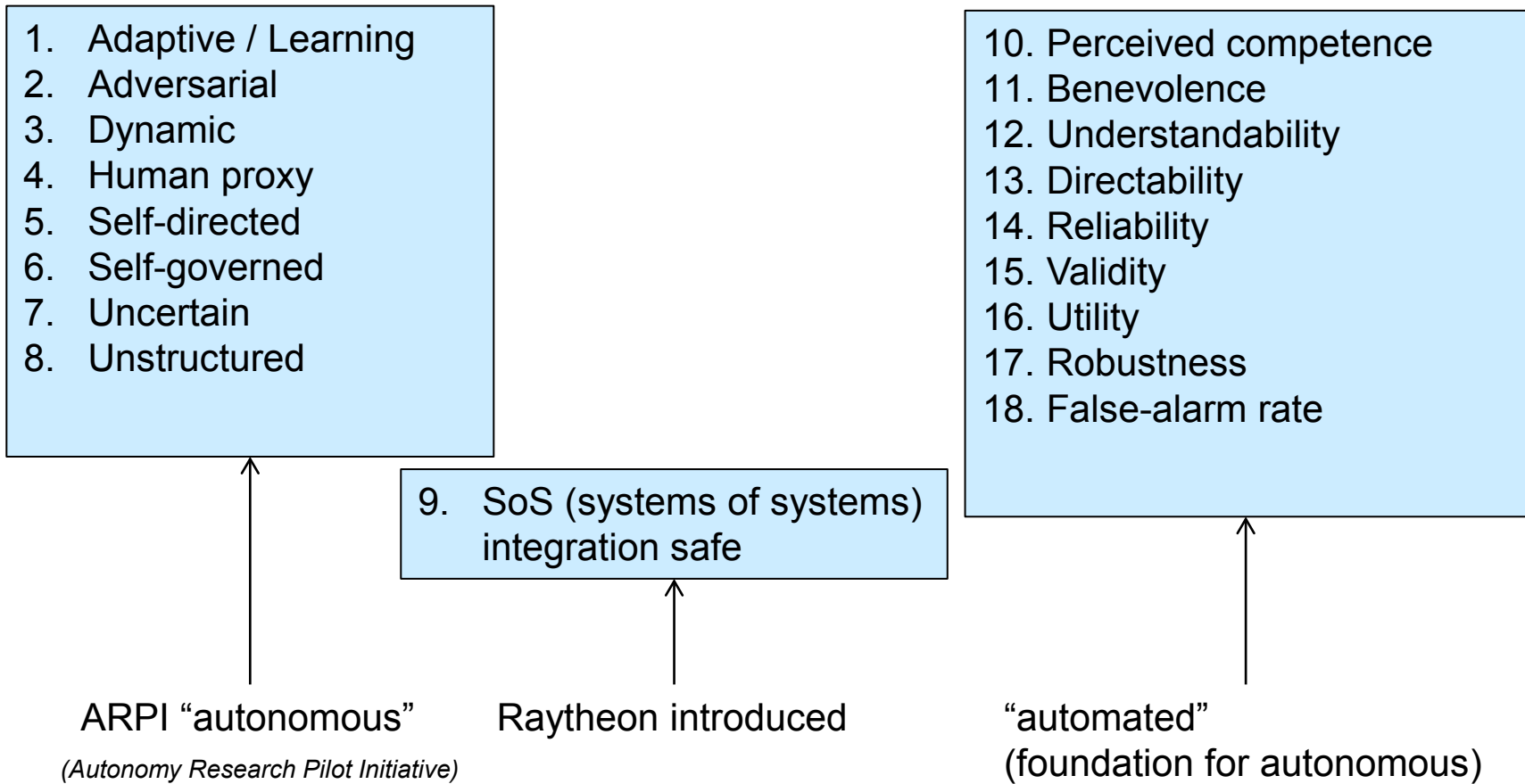
**While “Trust” adds belief in a result, it does not change the result**

# Systems Engineering V and "Trust V"



**Key insight: treat "Trustability" as a non-functional requirement**

# 18 Dimensions of Trust



**Need to create Trust in each dimension**

# Autonomy Trust V techniques

|                               | Analysis preceding Requirements             | 2 Requirements Analysis Process  | 3 Architectural Design Process  | 4 Implementation Process  | 6 Verification Process  | 8 Validation Process  | 9 Operation Process  |
|-------------------------------|---|--|---|---|---|---|--|
| broadly applicable methods    | For most of these, can use: (*) Simulation; | For most of these, can use: (*) Dynamic Fault Trees<br>For most of these must have (*) Incorporation of MOEs |   |   | For most of these, can use: (*) Semantic Q&A Capability; (*) Future scenario prediction   | For most of these, can use: (*) Semantic Q&A Capability; (*) Future scenario prediction       | For most of these, can use: (*) Semantic Q&A Capability; (*) Future scenario prediction<br>(*) Automated capture, analysis and feedback loop for Measures of trust |
| Autonomy characteristics      | Adaptive / Learning                         |  |   |   |   |   |  |
|                               | Adversarial                                 | Information Assurance; Cyber Protection / Security; IFF Capability   | Information Assurance; Cyber Protection / Security; IFF Capability                                    | Cyber Protection / Security; IFF Capability   | Information Assurance; Anti-tamper  | Negative testing  |  |
|                               | Dynamic                                     |  | Boundary definitions; Temporal considerations; Goal driven based on Beliefs, Desires, and Intention   |   | Graceful degradation; Boundary testing  | Testing to failure  |  |
|                               | Human interaction                           |  | Flexible autonomy; Human profiling; Require human adjustable interaction                              |   | Unified mission-based language for human / machine interaction  |   | Performance confirmation   |
|                               | Self-directed                               | Create positive statements of prevention of negative   | Build in Directability; Include positive statements;  | build in Directability  | Run time assurance  | Run time assurance (purpose); Formal methods based on Beliefs,                                | Test against CONOPS  |
|                               | Self-governed                               |  |   |   |   |   |  |
|                               | Uncertain                                   | Prediction of Mission Threats  | Measure of uncertainty & confidence; Reasoning of Situational Awareness for dynamic threat prediction | Sensor reliability analysis   |   |   | "Turing" Test<br>Future Risk Prediction method   |
| Unstructured                  |   | Environmental abstraction and characterization of mission threats; Measure of environment modeling           |   | Inference engine to categorize Mission Threats; Boundary testing; Safety mechanisms | Extreme testing to failure; Hierarchical tested across abstraction levels; State transitions  |   |  |
| Automation - trust attributes | SoS integration safe                        | TLYF (test like you fly); ConOps Analysis  |   |   |   | TLYF  | SoS certification; Thread based system testing   |
|                               | Perceived competence                        |  | Executable Models; Calibrated trust   |   | *Model CONOPS -- needed for "Semantic Q&A Capability"; *Identified Human trust needs -- needed for "ATCI"   |   | ATCI (Calibrated Trust)  |
|                               | Benevolence                                 | Hazard analysis  | SMRI (SW Mishap Risk Index); Adaptable safeguards   | Domestically manufactured components  | Adaptable safeguards  | Negative testing  |  |
|                               | Understandability                           |  | Restrictive language (e.g., Behavior Trees); Model Based Engineering; DoDAF artifacts                 | DoDAF artifacts   | Traceability; *Model DoDAF artifacts -- needed for "Semantic Q&A Capability"; Machine transparent; "dead code" removal; traceability  |   | Semantic Q&A   |
|                               | Directability                               |  | Requirements specification  |   |   | Demonstration of override   | Demonstration of override  |
|                               | Reliability                                 | MTBF analysis; MTTFF analysis; FMECA   | Independent concurrent SW/HW efforts; use of "taint checking" use of "trademarking"                   | Redundancy; Fault tolerant processing   | Run time assurance; Reuse; Mean time to Recovery; Graceful Degradation; Model Based Engineering; FMECA; autonomous footprint as small as possible; Automated Unit Test; Continuous Integration; Automatic code generation | Hybrid systems verification; Accelerated aging; Test Optimization; Run time assurance; DO178C |  |
|                               | Validity                                    | External SME Validation  | Formal Methods Acceptance; Model Based Engineering  |   |   | Formal Methods Acceptance (purpose)   | External SME Validation;   |
|                               | Utility                                     |  | Statement of performance  |   | Run time assurance  | Proof of performance; Run time assurance  | Test against CONOPS  |
|                               | Robustness                                  | Monte Carlo simulations; Fault Tree Analysis   |   | Redundancy; Fault tolerant processing   | Reuse; DFMA; Graceful Degradation   | Longevity Testing; Fault injection; Black Box Testing;  |  |

Many table entries – Trust Dimension versus Development Step

# Autonomy Trust V techniques – top left corner

|                             | Analysis preceding Requirements  | 2 Requirements Analysis Process  | 3 Architectural Design Process                 | 4 Implementation Process   |
|-----------------------------|--|--|--|--|
| broadly applicable methods  | For most of these, can use:<br>(* Simulation;                            | For most of these, can use:<br>(* Dynamic Fault Trees<br>For most of these must have<br>(* Incorporation of MOEs |  |  |
| <b>Adaptive / Learning</b>  |  |  |  |  |
| <b>Adversarial</b>          | Information Assurance;<br>Cyber Protection / Security;<br>IFF Capability | Information Assurance;<br>Cyber Protection / Security;<br>IFF Capability   | Cyber Protection / Security;<br>IFF Capability | Information Assurance;<br>Anti-tamper  |
| <b>Dynamic</b>              |  | Boundary definitions;<br>Temporal considerations;<br>Goal driven based on Beliefs, Desires,<br>and Intention     |  | Graceful degradation;<br>Boundary testing  |
| <b>Human interaction</b>    |  | Flexible autonomy;<br>Human profiling  |  | Unified mission-based language for<br>human / machine interaction                                      |
| <b>Self-directed</b>        | Create positive statements<br>of prevention of negative                  | Build in Directability;  | build in Directability                         | Run time assurance   |
| <b>Self-governed</b>        |  | Include positive statements;   |  |  |
| <b>Uncertain</b>            | Prediction of Mission<br>Threats   | Measure of uncertainty & confidence;<br>Reasoning of Situational Awareness for<br>dynamic threat prediction      | Sensor reliability analysis                    |  |
| <b>Unstructured</b>         |  | Environmental abstraction and<br>characterization of mission threats;<br>Measure of environment modeling         |  | Inference engine to categorize Mission<br>Threats;<br>Boundary testing;<br>Safety mechanisms           |
| <b>SoS integration safe</b> | TLYF (test like you fly);<br>ConOps Analysis                             |  |  |  |
| <b>Perceived competence</b> |  | Executable Models;<br>Calibrated trust   |  | Model CONOPS -- needed for "Semantic<br>Q&A Capability";<br>Identified Human trust needs -- needed for |
| <b>Benefit</b>              |  | Adaptable safeguards   | components                                     |  |

**Many table entries – Trust Dimension versus Development Step**

# Autonomy Trust V MOEs

| type   |                      | Operational Trust Metrics   |   |  |
|--|----------------------|---|---|--|
|  |                      | passive   | active  | historical   |
| Autonomy characteristics                                   | Adaptive / Learning  |   |   | Determine if system response to similar inputs has changed over time   |
|  | Adversarial          |   | Evaluate results of system operation applied to operator/SME supplied what-if examples                                  | Number of adversaries identified, and percentage of those addressed  |
|  | Dynamic              | If operator must accept tainted results before they can be used then number of tainted results accepted by operator |   | <i>(if has occurred)</i> Assessment of quality of graceful degradation   |
|  | Human interaction    |   | SUS score   |  |
|  | Self-directed        |   |   | Percentage of times that system "took appropriate action" without operator intervention given an   |
|  | Self-governed        |   |   | Operator/SME rating of historical examples of system response to actually observed   |
|  | Uncertain            |   |   |  |
|  | Unstructured         |   |   |  |
| SoS integration safe                                       |                      |   |   |  |
| Trust  | Perceived competence | Percentage of time operator chooses to not override system (decrease);  | Evaluate results of system operation applied to operator/SME supplied what-if examples                                  | Number of "functions" executed per unit time by computer (vs humans) (desire increased);<br>Ratio for completion of "Mission-Critical Objectives" vs. "Secondary Objectives" (increase);<br>Time to respond to critical events (decrease);<br>Time to respond to non-critical events (decrease). |
|  | Benevolence          |   | "Lead it into temptation and see if it delivers evil" ==> Create "tempting" scenarios and assess if malevolence results | Percentage of time that system operates outside of "safeguarded" regions   |
|  | Understandability    | Number operator queries per unit time;<br>Percentage of queries which are followed immediately by another query     | Evaluate results of query usefulness by operator/SME;<br>SUS score  | Time required for operator training (desire decrease)  |
|  |                      | Response time when operator overrides system;   |   |  |
| <b>Many table entries– Trust Dimension versus MOE type</b> |                      |   |   |  |
|  | back to system       |   |   |  |



# “Chatter” and “Why/What Agent”

## 1. Operator trust in the system is enhanced when

- The **user understands** why the system performs as it does
- Communications are in the **language of the user**
- The user is pleasantly **surprised** by meaningful yet unexpected information
  - Think “automated performance monitoring”

## 2. The speed of trust is enhanced when

- The operator controls the **frequency of communication**
- The operator controls the **information content level**
  - Layered control in many systems
  - HW: assemblies to sub-assemblies to components
  - SW: libraries to modules to functions
- The operator controls the **information details**

## 3. Communications between system and user

| Attribute                   | Tool | Chatter          | What Agent      | Why Agent       |
|-----------------------------|------|------------------|-----------------|-----------------|
| Query structure             |      | pre-defined      | scripted        | free-form       |
| Content level               |      | operator selects | script driven   | operator driven |
| Information details         |      | operator selects | script driven   | operator driven |
| Frequency of communications |      | operator selects | operator driven | operator driven |
| Driven by                   |      | system           | operator        | operator        |

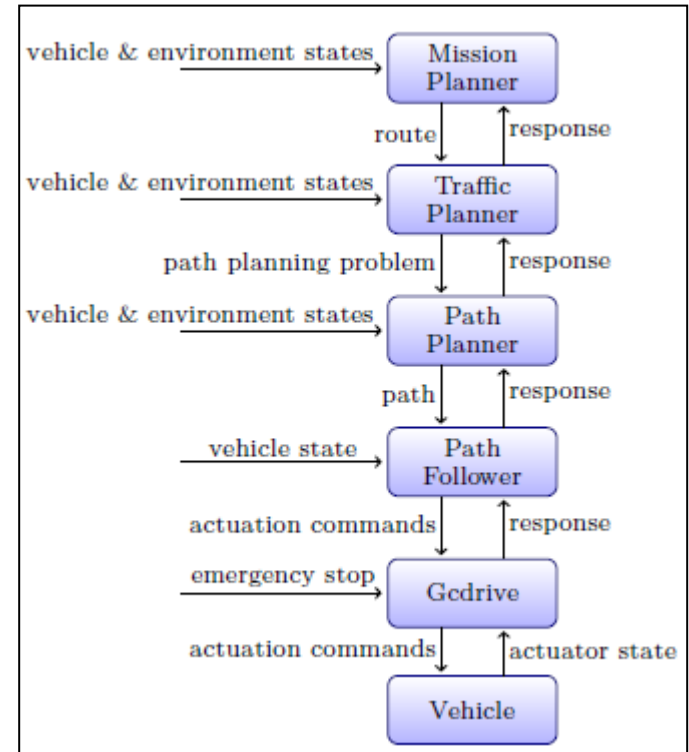
# Many systems use a layered control model

## BASIS - Fast Attack Craft / Fast Inshore Attack Craft (FAC/FIAC)

1. Threat Change
2. Behavioral Change
3. Last Trigger Change
4. Input Change

## Alice – Caltech’s entry in DARPA’s 2007 Urban Challenge – “navigation protocol stack” is

1. Mission Planner
2. Traffic Planner
3. Path Planner
4. Path Follower
5. Gcdrive
6. Vehicle



Want to know – at level (1-4) or (1-6)

- What changed?
- Why did it change? (recursively)

“Synthesis of Control Protocols for Autonomous Systems”, Wongpiromsarn, Topcu, & Murray, December 2012, Figure 1, [https://www.google.com/url?q=http://www.cds.caltech.edu/~murray/preprint s/wtm12-us\\_s.pdf&sa=U&ei=k9ZGU66qPlamsAS154CoBw&ved=0CBsQFjAA&usg=AFQjCNEb4KRpsDhTx9cqrthDzCZ7\\_zJnzQ](https://www.google.com/url?q=http://www.cds.caltech.edu/~murray/preprint%2Fs/wtm12-us_s.pdf&sa=U&ei=k9ZGU66qPlamsAS154CoBw&ved=0CBsQFjAA&usg=AFQjCNEb4KRpsDhTx9cqrthDzCZ7_zJnzQ)

**Chatter naturally supports hierarchical control systems**

# “Why/What Agent” and “Chatter” graphical representation

**Why/What Agent** → Operator initiated queries

- Either
  - What → Template driven → “What is the reason for \_\_\_\_\_?” *easier*
  - Why → Free form input → “Why did the system do \_\_\_\_\_?” *harder*
- Operator can drill down to desired hierarchy level as needed
- Operator controls “level” of response

## Why/What

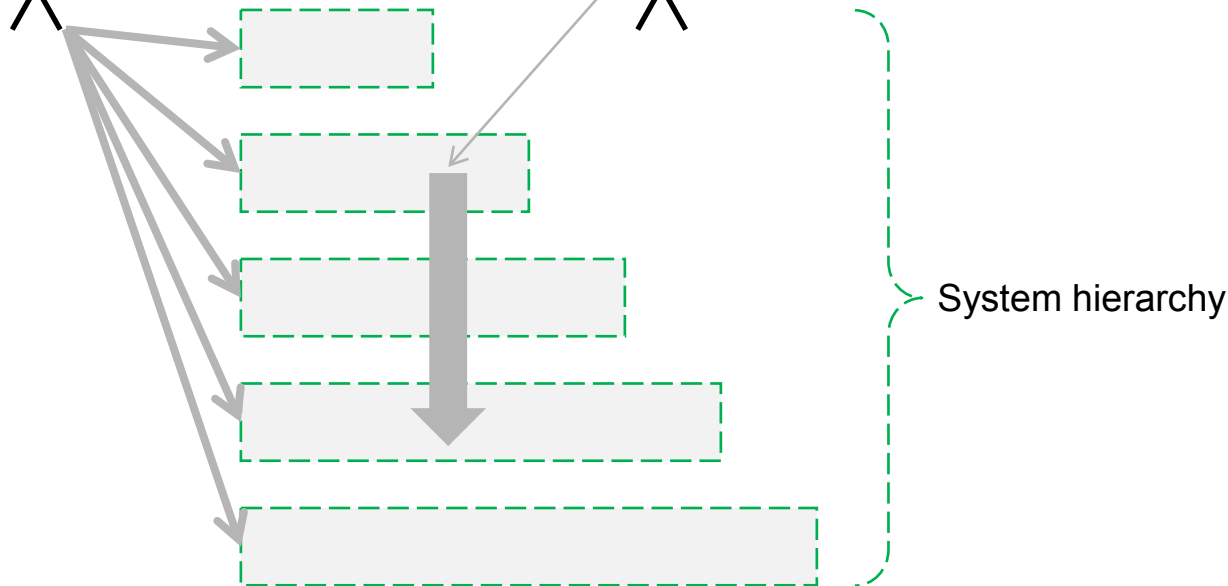
- Responses from anywhere



## Chatter

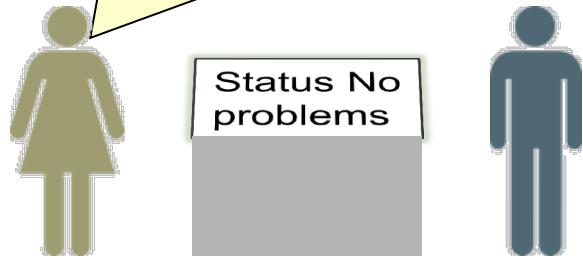
→ System initiated communications

- Frequency of communication
- Information content level
- Information details

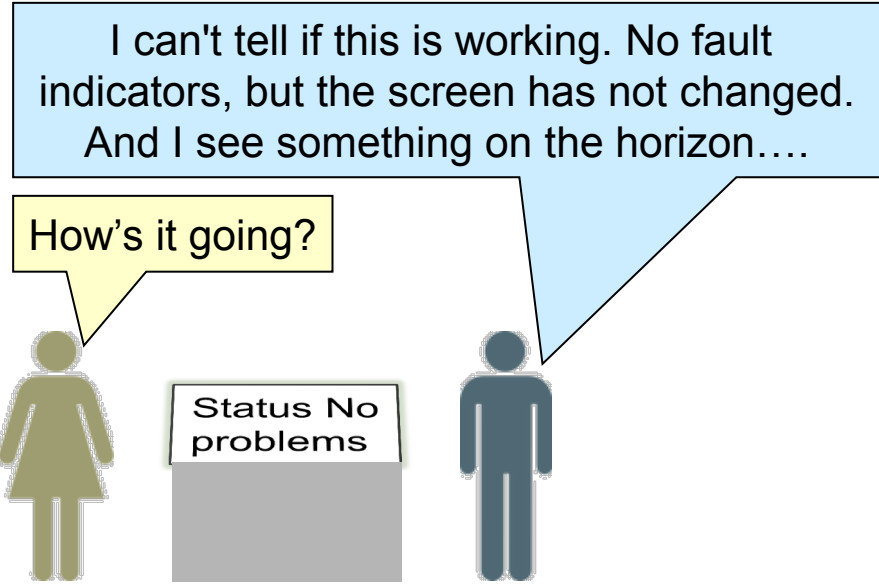


# Example of Chatter use

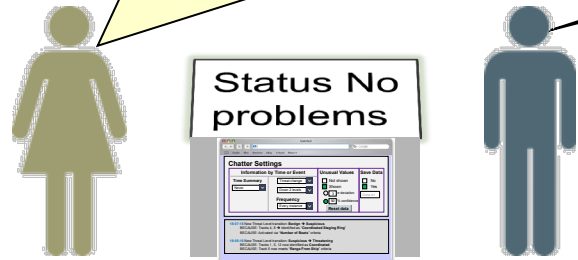
Here's the system to monitor.  
Let me know when it says "Alert"



2 days go by



Look at the Chatter. It's processing many measurements and its seen the objects on the horizon. They're not a threat.



# “Chatter” – notional interface



Untitled

Apple .Mac Amazon eBay Yahoo! News ▾

## Chatter Settings

| Information by Time or Event                     | Unusual Values  | Save Data  |
|--|---|--|
| <p><b>Time Summary</b></p> <p>Not Reported ▾</p> | <p>Threat change ▾</p> <p>Down 2 levels ▾</p> <p><b>Frequency</b></p> <p>Every instance ▾</p> | <p><input type="checkbox"/> Not shown</p> <p><input checked="" type="checkbox"/> 3 <math>\sigma</math> deviations<br/>(after 20 samples)</p> <p><input type="checkbox"/> No</p> <p><input checked="" type="checkbox"/> Yes</p> <p><i>Data.txt</i></p> <p><b>Reset data</b></p> |

**18:07:18** New Threat Level transition: **Benign** → **Suspicious**  
 BECAUSE: Tracks 4, 8 → identified as “**Criteria ABC**”  
 BECAUSE: Activated via “**Boat Criteria 3**” criteria

**18:05:16** New Threat Level transition: **Suspicious** → **Threatening**  
 BECAUSE: Tracks 1, 5, 12 now identified as **Coordinated**  
 BECAUSE: Track 5 now meets “**Criteria DEF**” criteria

Applied to Fast Attack Craft / Fast Inshore Attack Craft (FAC/FIAC)

# “Chatter” – Operator controls

- Threat change
- Behavioral level
- Internal Change
- Input Change

### Chatter Settings for BASIS

| Information by Time or Event |                                    | Unusual Values  | Save Data  |
|------------------------------|------------------------------------|---|--|
| <b>Time Summary</b><br>Never | Threat change<br>Down 2 levels     | <input type="checkbox"/> Not shown<br><input checked="" type="checkbox"/> Shown<br>3 $\sigma$ deviation<br>90 % confidence<br><input type="button" value="Reset data"/> | <input type="checkbox"/> No<br><input checked="" type="checkbox"/> Yes<br>Data.txt |
|                              | <b>Frequency</b><br>Every instance |   |  |

- Myself Only
- Down 1 level
- Down 2 levels
- Down 3 levels

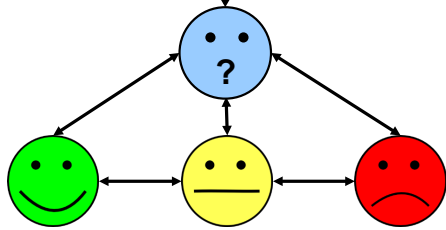
- Never
- Every 1 minute
- Every 5 minutes
- Every 15 minutes
- Every hour

- Never
- Every Instance
- Every 10th instance
- Every 100th instance
- Every 1000th instance


Operator controls the level/details/frequency of communications

# Future Risk Prediction – notional interface


System startup




**Future Risk Prediction**

 Risk unchanged for next Assessment confidence unknown unknown


**Future Risk Prediction**

 Risk unchanged for next Assessment confidence 5 minutes 89%


**Future Risk Prediction**

 Risk unchanged for next Assessment confidence 4 minutes 80%

**Future Risk Prediction**

 Risk unchanged for next Assessment confidence 2 minutes 92%  
Risk driven by

**Future Risk Prediction**



Risk unchanged for next Assessment confidence 5 minutes 89%

|                          |    |
|--------------------------|----|
| Number of similar states | 23 |
| 1 had negative results   | 4% |

Current Pursuit

|                           |     |
|---------------------------|-----|
| Number of states required | 15  |
| State match threshold     | 0.1 |
| Time increment (min)      | 1   |

Explanation of high risk assessment

Output shows independent risk assessment of system capability

# BACKUP



# Definitions from ARPI *(Autonomy Research Pilot Initiative)*

## *Automation*

The system functions with no/little human operator involvement. However, the system performance is limited to the specific actions it has been designed to do. Typically these are well-defined tasks that have predetermined responses, i.e. rule-based responses in reasonably well-known and structured environments.

## *Autonomy*

Systems which have a set of intelligence-based capabilities that allow it to respond within a bounded domain to situations that were not pre-programmed or anticipated in the design (i.e., decision-based responses) for operations in unstructured, dynamic, uncertain, and adversarial environments. Autonomous systems have a degree of self-governance and self-directed behavior and must be adaptive to and/or learn from an ever-changing environment (with the human's proxy for decisions).

# Dimensions Describing the Basis of Trust

TABLE 1: Summary of the Dimensions That Describe the Basis of Trust

| Study                          | Basis of Trust   | Summary Dimension  |
|--------------------------------|--|--|
| Barber (1983)                  | Competence<br>Persistence<br>Fiduciary responsibility  | Performance<br>Process<br>Purpose  |
| Butler & Cantrell (1984)       | Competence<br>Integrity<br>Consistency<br>Loyalty<br>Openness  | Performance<br>Process<br>Process<br>Purpose<br>Process  |
| Cook & Wall (1980)             | Ability<br>Intentions  | Performance<br>Purpose   |
| Deutsch (1960)                 | Ability<br>Intentions  | Performance<br>Purpose   |
| Gabarro (1978)                 | Integrity<br>Motives<br>Openness<br>Discreetness<br>Functional/specific competence<br>Interpersonal competence<br>Business sense<br>Judgment | Process<br>Purpose<br>Process<br>Process<br>Performance<br>Performance<br>Performance<br>Performance |
| Hovland, Janis, & Kelly (1953) | Expertise<br>Motivation to lie   | Performance<br>Purpose   |
| Jennings (1967)                | Loyalty<br>Predictability<br>Accessibility<br>Availability   | Purpose<br>Process<br>Process<br>Process   |
| Kee & Knox (1970)              | Competence<br>Motives  | Performance<br>Purpose   |
| Mayer et al. (1995)            | Ability<br>Integrity<br>Benevolence  | Performance<br>Process<br>Purpose  |
| Mishra (1996)                  | Competency<br>Reliability<br>Openness<br>Concern   | Performance<br>Performance<br>Process<br>Purpose   |
| Moorman et al. (1993)          | Integrity<br>Willingness to reduce uncertainty<br>Confidentiality<br>Expertise<br>Tactfulness<br>Sincerity<br>Congeniality<br>Timeliness     | Process<br>Process<br>Process<br>Performance<br>Process<br>Process<br>Performance<br>Performance     |
| Rempel et al. (1985)           | Reliability<br>Dependability<br>Faith  | Performance<br>Process<br>Purpose  |
| Sitkin & Roth (1993)           | Context-specific reliability<br>Generalized value congruence   | Performance<br>Purpose   |
| Zuboff (1988)                  | Trial and error experience<br>Understanding<br>Leap of faith   | Performance<br>Process<br>Purpose  |

## Review of trust dimensions

- *Trust in Automation: Designing for Appropriate Reliance* by Lee & See
- <http://www.engineering.uiowa.edu/~csl/publications/pdf/leesee04.pdf>

## List of trust dimensions

- *benevolence*
- *directability*
- *false-alarm rate*
- *perceived competence*
- *reliability*
- *robustness*
- *understandability*
- *utility*
- *validity*
- *Trust in Automation*, Hoffman, Johnson, Bradshaw, Underbrink
- <http://www.jeffreybradshaw.net/publications/50.%20Trust%20in%20Automation.pdf>