Calibrated Trust vs Blind Faith

Testing AI/ML Powered Emerging Technologies for Aviation



Col Dan "Animal" Javorsek, PhD

Commander, AFOTEC Det 6 Director, F-35 US Operational Test Team

21 September 2022













The Augmented Age is the next great era

Enabled by software defined systems augmenting human decisions



Distribution A: Approved for Public Release, Distribution Unlimited

7



History of Al





Assumptions when applying AI to air combat

- Air combat problems consist of two time constants
 - immediate tactic selection
 - long-term strategic impact
- The air combat problem is manifold bounded









Neural nets are just spreadsheets on steroids

 Each layer stretches and squashes the data space until the manifolds are cleanly separated









Evolution of the AI playground for strategy and decision making

 $\overline{\triangleleft}$



comment

Deep Blue beats 12-time world champion and grandmaster Garry Kasparov in 3 games to zero (with 3 draws)



C AFP





AlphaGo beats 18-time world champion Lee Sedol in 4 games to one



Lee Jin-man / Associated Press



STAR

AlphaStar

AlphaStar beats TLO and

MaNa (two of the best pro

gamers in the world) 10-0

© Google Deepmind



AlphaDogfight

AlphaMosaic

Mosaic ACE TA3, ACK and others

Once useful models are developed, AI will be applied to determine the best ways to beat the game



Ernest et al.

11

202













Mishap Report: F-16 Controlled Flight Into Terrain (CFIT)



https://youtu.be/WkZGL7RQBVw



Mishap Report: F-16 Controlled Flight Into Terrain (CFIT)



ttps://youtu.be/WkZGL7RQBVw



Source: Scaled Composites

N62270



Autonomy was born in the air, now behind, why?

Commercial vehicle autonomy closing in on full autonomy [5]... ...while aerial combat autonomy struggles to reach conditional autonomy [3]





miles

What is the vehicle autonomy approach to get to level 5? A matter of quantity of quality data \rightarrow Optimized for time to Market

- 2 challenges must be solved to get to market as fast as possible:
 - Performance: is it capable of handling what's required
 - Trust: expect it to deploy capability as required

by the way: these are the same characteristics we build/assess while training human pilots

Insight in their metric

What is measured is managed -- so what do they measure? Assessment rolled up into a single metric: Miles per Disengagement

Also from RAND report:

Miles alone not enough

Autonomous vehicles would have to be driven hundreds

of millions of miles and sometimes hundreds of billions

of miles to demonstrate their reliability in terms of fatali-

fleets would take tens and sometimes hundreds of years

aim is to demonstrate their performance prior to releas-

to drive these miles—an impossible proposition if the

• Therefore, at least for fatalities and injuries, test-driving

alone cannot provide sufficient evidence for demonstrat-

ing them on the roads for consumer use.

ing autonomous vehicle safety.

Under even aggressive testing assumptions, existing

Key findings

ties and injuries.



- Waymo announced in January 2020 that its autonomous cars have driven tens of billions of miles through computer simulations and 20 million miles on public roads in 25 cities"
- ...and they are just 1 of 36 companies testing in California alone!



agement

0

disen

e L

Õ

Need operator engagement to sufficiently build trust





All the significant players are using a human on-the-loop configuration Miles per disengagement verifies performance, builds trust, enables mileage volume

These Companies Are Testing Self-Driving Cars in California

Number of autonomous vehicles registered to be tested on public roads in California*



100s of autonomous vehicles testing on public roads in California alone...

















...what does every single vehicle have in common? human-on-the-loop



Another byproduct of approach is incremental feature rollout

Deploy minimum viable assistance autonomy features to build trust incrementally

LEVELS OF AUTONOMY



Driver is still the captain

relinquish certain controls at certain times to driver assist technologies.

of the ship, but can

No Automation

The full-time performance by the human driver of all aspects of the dynamic driving task, even when enhanced by warning or intervention systems *(i.e. Cars from 20 years ago)*

Driver Assistance

A driver assist system of *either* steering or acceleration/ deceleration using information about the driving environment and with the expectation that the human driver perform all remaining aspects the driving task.

Partial Automation

Driver assist systems are able to control both steering and acceleration/deceleration by using information about the driving environment while the human driver performs all remaining aspects of the driving task (*l.e. The system executes the steering/accerlation and deceleration, while the human monitors and acts as the fallback "system"*).

Conditional Automation

automated driving systems perform all aspects of the driving task, with the expectation that the human driver will respond and intervene when needed.

High Automation

Performance by an automated driving system of all aspects of the driving task, even if a human driver does not respond appropriately to a system request or warning for human intervention. (*delivers the capability in a defined area/geofenced area and weather constraints*)

Full Automation

The full performance of driving by an automated driving system under all roadway and environmental conditions that can also be managed by a human driver, but human intervention is not needed. *(unrestricted area and weather)*

Front Assist Blind Spot Monitor Rear Traffic Alert Lane Assist Park Assist Adaptive Cruise Park Pilot Control



Incremental feature roll-out of increasing levels of automation: verifies performance while building trust



Super-human effectiveness of "pilot-assist features" demonstrated Death Claw and AlphaDogfight have both shown super-human performance in specialized tasks

ALIAS's Death Claw

- Automated final portion of gun targeting to keep gun pipper on targets—a difficult & high workload task
- Flight tests using NF-16 VISTA test aircraft (USAF TPS) in November 2017
- Substantial performance improvements with unanimous, positive pilot comments

ACE's AlphaDogfight

- Eight contractor teams competed for top AI agent before engagement with Weapons Instructor Course grad
- Al agent defeated Banger 5-0 in high aspect BFM due to superior aiming and decision making





Demonstrated effectiveness of automated aiming in air-to-ground and air-to-air employment



Air Force is currently betting on a different approach





The spectrum of AI approaches have already been explored





The tradeoff between novelty and explainability



AOG: And-Or-Graphs MLN: Markov Logic Network CRF: Conditional Random Fields SRL: Statistical Relational Learning HBN: Hierarchical Bayesian Networks SVM: Support Vector Machine



Hierarchical/Hybridized solutions result in best of both worlds

- A Hierarchical/Hybridized AI solution blends policies from:
 - Reinforcement learning artificial neural networks (referred to as multi-layered perceptrons)
 - Markov decision processes for the incorporation of uncertainties
 - Hierarchy and rules-based policies for explainability





Hierarchical/Hybridized solutions combine two techniques



- Performance increase via additional states
- Exponential design and test burden
- Straightforward verification; unbounded validation



- Performance via simultaneous state processing
- Exponential training burden (simulation)
- Validation via training; unbounded verification

Must allow for, and even anticipate, simulation/reality mismatch



Reconsider Validation, Verification, and Accreditation (VV&A)

• One example of a trust methodology currently being explored





Divide the pilot's attention...

...measure workload distribution

AlphaDogfight Trials

UIRTUAL FINALS 8.18-20.2020

DARPA







Autonomy applications for Air Combat go beyond dogfighting



Distribution A: Approved for Public Release, Distribution Unlimited

Human-On-The-Loop Manned Unmanned Teaming

Fleet Learning 33



Our current systems are not as deterministic as we think

Consider performing envelope expansion testing on the following two aircraft



- Federated subsystems
- Gaussian => Predictable
- Fewer test points
- Predict, Test, Validate Works
- Build-Up Approach Works
- Complicated



- Interconnected subsystems
- Non-Gaussian => Unpredictable
- More test points
- Predict, Test, Validate Tough
- Build-Up Approach Tricky
- Complex



Structure of a system: Simple





Structure of a system: Complicated





Structure of a system: Complex







Distribution A: Approved for Public Release, Distribution Unlimited



M&S, Digital Twins, and a continuous TEVVA Process



• Entry Criteria

• Steps

- 1. Define the M&S/digital twin intended use
- 2. Evaluate risk
- 3. Define accreditation criteria
- 4. Select the candidate M&S/digital twin
- 5. Develop TEVV plans
- 6. Perform TEVV
- 7. Prepare accreditation support packages
- 8. Accreditation assessment
- 9. Accreditation decision
- Exit Criteria



Must convert all manned aircraft to autonomy test assets

Air Force does not have test assets to obtain flight hours required

275 million miles required to prove a self driving vehicle is at least as trustworthy as a human – how many hours does the Air Force have?



To go from level 2 to level 5:



Aircraft: Full Squadrons

□ Flight hours: 1000s

Metric: Hours per disengagement







Greatly improved autonomy is coming to the cockpit





- The Augmented age of the 21st Century is upon us
 - Complete with Al
 - Hardware-enabled, software-defined systems
- Winning in autonomy is about volume: quantity of quality training data
 - Time to market: Must have a human-on-the-loop to adequately build performance and trust quickly
 - US military is taking a different approach effectively skipping human-on-theloop but this is changing
- We must dispel the myth that we cannot write a requirement for adaptability or perform continuous V&V
 - In the modern era full of complex systems we must modernize our Test, Evaluation, Validation, Verification, and Accreditation (TEVVA) processes
 - By choosing the right metrics, performing flight test suitability checks, and flipping the current role of M&S/digital twins we have a chance to tackle 21st Century challenges





