# AFRL

# Autonomy Capability for Next Generation Air and Space Dominance

Dr. Kevin "Nagel" Schmidt

AFRL Autonomy Capability Team (ACT3)

# Agenda

- **What is AFRL's ACT3**

- **Applied State-of-the-art AI to USAF and USSF challenges**
  - **Autonomous Air Combat Operations (AACO)**
  - **Space Autonomy**
  - **Air/Space Safe, Ethical Autonomy**

- **Next Generation AI**
  - **Disruptive Capabilities = Cutting Edge Research & Development**
  - **QuEST = Qualia Exploitation of Sensing Technology**

Approved for public release: distribution unlimited. Case Number AFRL-2023-1547

# What is AFRL's Autonomy Capability Team (ACT3)?

- Applying industry best practices to **scale AI solutions** across the DAF (USAF/USSF)
- Organic AI **expertise from across AFRL** enterprise, growing group since 2017
- Out of the box solutions from our people **outside the gate**
- Dev/Ops and Research in **custom IT environment**
- Integrated Civ, Ctr, Mil Teams with focus, **technical depth**



Autonomous Horizons v2: https://www.airuniversity.af.edu/AUPress/Display/Article/1787830/autonomous-horizons-the-way-forward

# Air/Space Force Cognitive Engine

A NEW business model for invention / development / fielding of AI

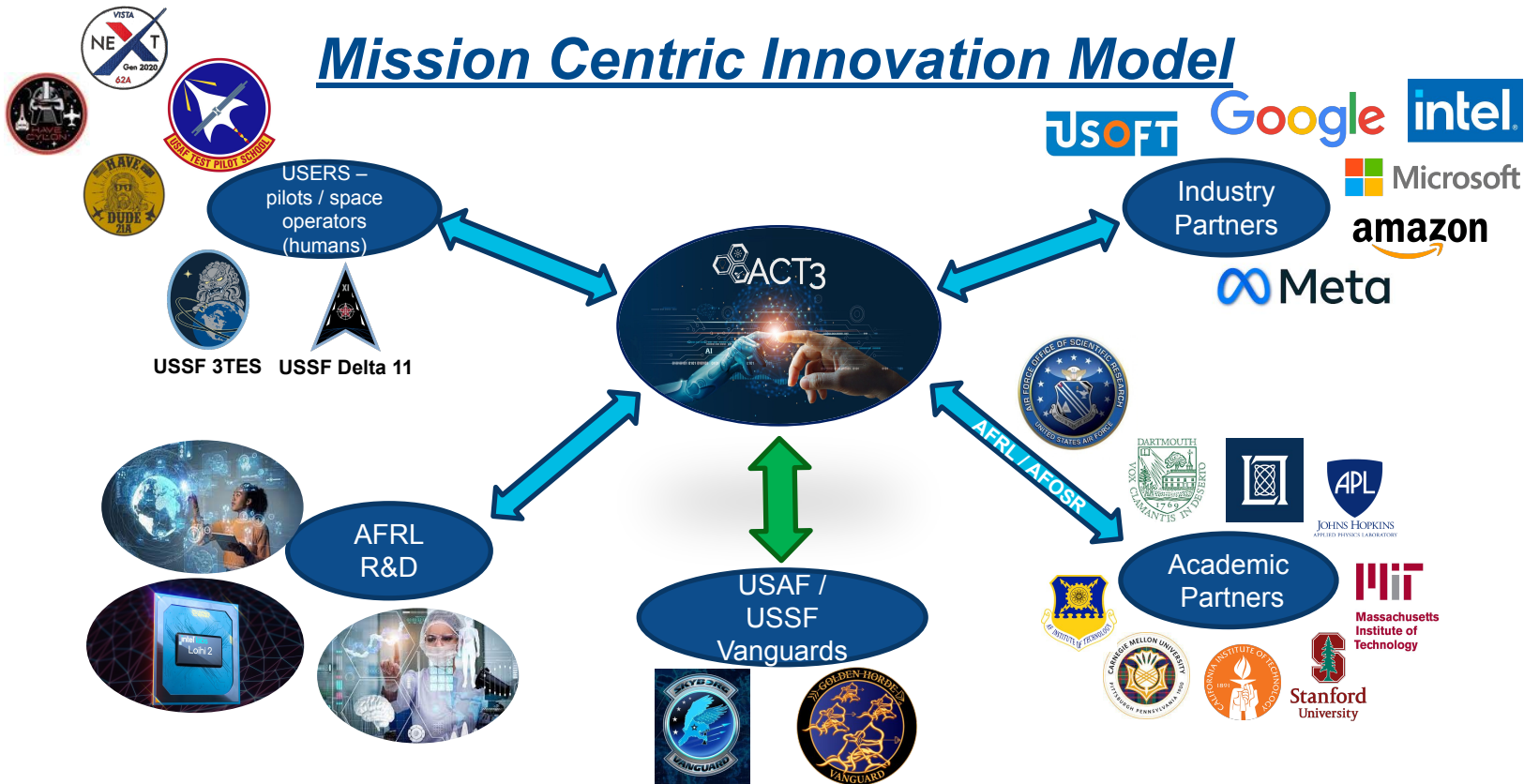| | |
|---|---|
| Execute AI pilot projects to gain momentum | Develop an AI / Analytics Strategy |
| Infrastructure: build an in-house AI / Analytics Team | Develop internal and external communications |
| Provide broad AI / Analytics Training | **Applying industry best practices to scale AI solutions across the AF/USSF** |

Mission Centric Innovation Model

User-Producer Innovation with the Art of the Possible

# Applying the State-of-the-Art AI to USAF/USSF Challenges

# Significant Achievements of Reinforcement Learning

RL has already shown superhuman performance in several big "human versus machine" competitions - high dimensional state spaces, partial observability, complex strategy
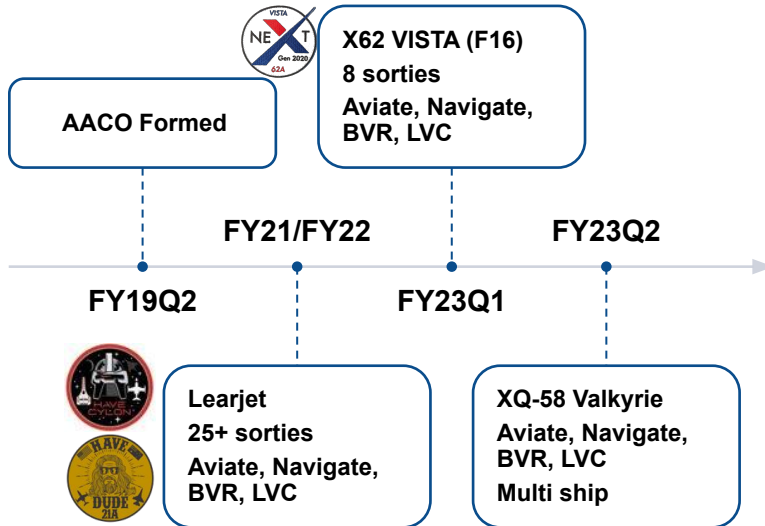
- AlphaGo - Deepmind trained agents to play GO (2016)
- OpenAI 5 - OpenAI trained agents to play DOTA 2 (2019)
- AlphaStar - Deepmind trained agents to play StarCraft 2 (2019)
- Alpha Dog Fight Trials - Heron Systems trained agents to aerial dogfight (2020)
- GPT-4 passes Bar Exam, US Medical Licensing Exam, college course exams (2023)



Photo Courtesy of DARPA

# Autonomous Air Combat Operations

***Vision Statement:*** AFRL ACT3 Autonomous Air Combat Operation (AACO) was formed to Demonstrate AI Tactical Autopilot engaging in multi-ship / multirole Beyond Visual Range (BVR) & ISR combat operations with tactical proficiency



**AACO Formed**

**X62 VISTA (F16)**
**8 sorties**
**Aviate, Navigate, BVR, LVC**

FY21/FY22

FY23Q2

FY19Q2

FY23Q1

**Learjet**
**25+ sorties**
**Aviate, Navigate, BVR, LVC**

**XQ-58 Valkyrie**
**Aviate, Navigate, BVR, LVC**
**Multi ship**



**Impacts**
Generating unprecedented data sets to support future AF needs
*AACO 100% government owned tools, agents, and capabilities are the base of DARPA AIR*

# Current Core Focus: Cooperating Space

Quickly react, plan, and decide on appropriate courses of action for inspection tasking in proximity operations in support of In-space Servicing, Assembly, and Manufacturing (ISAM).

For the following design reference missions (DRMs):

**DRM1**. Chief spacecraft provides sensing and state estimation for multiple deputy spacecraft.

**DRM2**. Multiple deputy spacecraft collaborate to inspect a chief spacecraft.



With the following high-level General Technical Objectives (GTOs):

**GTO1**. Develop flexible **human-autonomy interfaces** that accept operator-based mission preferences and provide course of action visualization, comparison, and selection mechanisms

**GTO2**. Develop **reinforcement learning-based neural network multi-agent controllers** that incorporate mission- and task-specific operator preferences in identification of courses of action

**GTO3**. Develop **run time assurance** approaches that mitigate hazards and allow the autonomy to stay on mission for longer in the face of unexplored system parameters, scenarios, and/or poorly modeled aspects of the system.

# Space Trusted Autonomy Levels

- Active Partner in the NRO/NASA/USSF Space Trusted Autonomy group under the Space Science and Technology Partnership forum

- Defined *Space Trusted Autonomy Readiness (STAR) Levels*
  - Defines **capability** and **trust** levels
  - Includes background on trust in automation research
  - **Consistent with other forms of readiness levels**: original TRLs, algorithm RLs, manufacturing RL, commercialization RL, data RL, machine learning RL, etc.

Kerianne L. Hobbs, Joseph B. Lyons, Martin S. Feather, Benjamin P Bycroft, Sean Phillips, Michelle Simon, Mark Harter, Kenneth Costello, Philip C. Slingerland, Yuri Gawdiak, Stephen Paine, "Space Trusted Autonomy Readiness Levels," IEEE Aerospace, Big Sky, MT, March 4-11, 2023. Preprint: https://arxiv.org/pdf/2210.09059

THE AIR FORCE RESEARCH LABORATORY

# Ethics, Safety and Trust in AI

- There is no such thing as Ethical AI, only Ethical use of AI
  - ***Trust*** is defined as a willingness to accept vulnerability in situations characterized by uncertainty.
  - ***Safety*** is freedom from harm during operations.
  - ***Ethics*** are rules created by societies and cultures governing moral and just usage of a technology.
- Ethical considerations begin at design time and continue through operational use
- Potential to use Model Cards or Datasheets for Datasets as communication of ethics from design to use

[Draft] Joseph B. Lyons, Kerianne Hobbs, Steve "Cap" Rogers, Scott H. Clouse, "Responsible (Use of) AI," in Understanding the Role of Humanity in Responsible Deployment of Intelligent Technologies in Socio-technical Ecosystems, Special Issue of Frontiers in Neuroergonomics, 2023.
[In Draft] Kerianne Hobbs, Bernard Li, "Safety, Trust, and Ethics Considerations for Human-AI Teaming in Aerospace Control," AIAA SciTech, 8–12 January 2024, Orlando, FL.

# Verifying and Guaranteeing Safety

- Run Time Assurance
- Safe RL, Rigorous Evaluation
- Neural Network Verification
- Hazard Analysis
- Assurance Cases

**Assured Controller**

Primary Controller $\xrightarrow{u_{des}(x)}$ RTA Mechanism $\xrightarrow{u_{act}(x, u_{des})}$ Plant

$x$

Selected Publications

1. Kerianne L. Hobbs, Mark L. Mote, Matthew C.L. Abate, Samuel D. Coogan, and Eric M. Feron "**Run Time Assurance for Safety Critical Systems:** An Introduction to Safety Filtering Approaches for Complex Control Systems," IEEE Control Systems Magazine. April 2023. Preprint: https://arxiv.org/pdf/2110.03506.pdf

2. [Submitted] Jonathan Rowanhill, Ashlie B. Hocking, Aditya Zutshi, Kerianne L. Hobbs "Conformance of Run-Time Assurance to MIL-HDBK-516C through Verified System **Arguments**" Digital Avionics Systems Conference 2023.

3. Hobbs, K., Heiner, B., Busse, L., Rowanhill, J., Hocking, A. B., Zutshi, A., "Systems Theoretic Process Analysis of a Run Time Assured Neural Network Control System," AIAA SciTech, National Harbor, MD, Jan. 2023. Preprint: https://arxiv.org/pdf/2209.00552

4. Diego Manzanas Lopez, Hoang-Dung Tran, Taylor T. Johnson, Stanley Bak, Xin Chen, Kerianne L. Hobbs, "Evaluation of **Neural Network Verification Methods** for Air to Air Collision Avoidance" AIAA Journal of Air Transportation Systems, 2022.

IEEE **Control systems**

APRIL 2023  VOLUME 43  NUMBER 2

**Runtime Assurance for Safety-Critical Systems**

AN INTRODUCTION TO SAFETY FILTERING APPROACHES FOR COMPLEX CONTROL SYSTEMS

KERIANNE L. HOBBS, MARK L. MOTE, MATTHEW C.L. ABATE, SAMUEL D. COOGAN, and ERIC M. FERON

**Safety Critical Control Systems**

IEEE Control Systems Society    IEEE

Approved for public release: distribution unlimited. Case Number AFRL-2023-1547
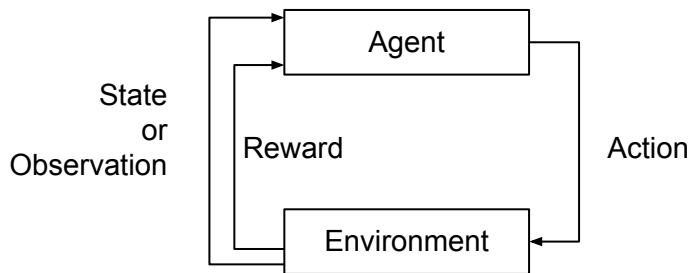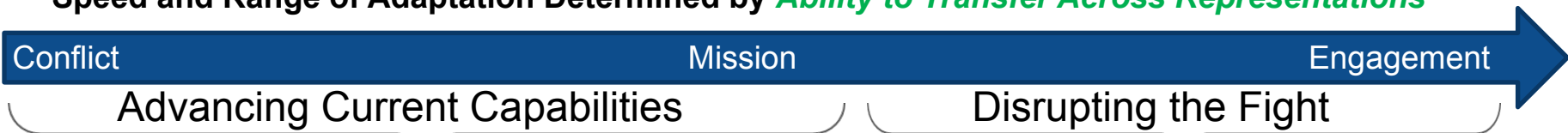
12

Creating the
Next Generation of AI

# Disruptive Capability Upgrades – Collaborative Combat Aircraft

**Speed and Range of Adaptation Determined by *Ability to Transfer Across Representations***

| Conflict | Mission | Engagement |
|---|---|---|

## Advancing Current Capabilities



State or Observation — Agent — Reward — Environment — Action

## Disrupting the Fight



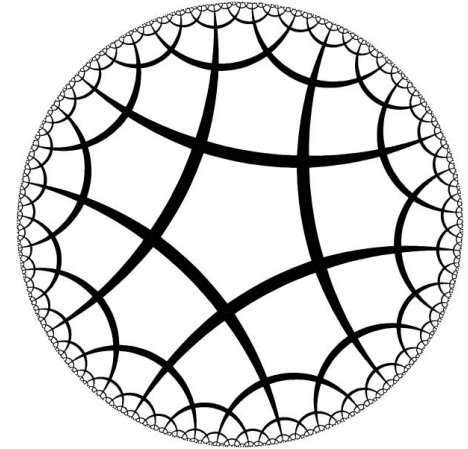- **Primary Problem #1: Efficient Agent Training**
- Catastrophic Interference
    - Effects of rewards and training curricula
    - Retraining components

**Need:**
- Real-time decentralized & distributed learning
- Within-mission model updates
- Robustness to disrupted, disconnected, intermittent, and low-bandwidth (DDIL) environments

**User-Producer Innovation with the <u>Art of the Possible</u>**

# Why Consciousness?



**Understand Human Cognition**



State
or
Observation

Agent

Reward

Action

Environment

**Build More Flexible Machines**



**Study Interesting Structures**

## Representations of Computational Processes

# Memory Consolidation



- Complementary Learning Systems

  - Memories first depend on the hippocampus

  - Hippocampus supports reinstatement of recent memories in the neocortex
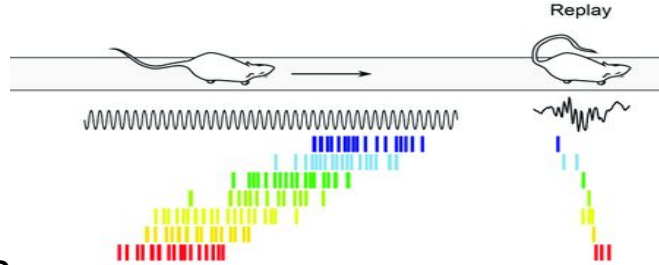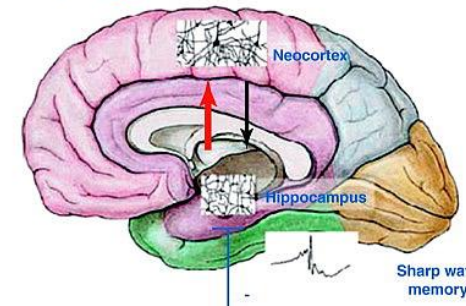
  - Neocortical synapses change a little on each reinstatement

- ***Interleaved learning and catastrophic interference***

McClelland, J.L., McNaughton, B.L., & O'Reilly, R.C. (1995). Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, *102*(3), 419–457. https://doi.org/10.1037/0033-295X.102.3.419

THE AIR FORCE RESEARCH LABORATORY

Approved for public release: distribution unlimited. Case Number AFRL-2023-1547

16

# Targeted Memory Reactivation



Schmidt, K., Larue, O., Kulhanek, R., Flaute, D., Veliche, R., Manasseh, C., Dellis, N., Culbertson, J., Clouse, H. & Rogers, S. (2023). Representational Tenets for Memory Athletics. arXiv preprint arXiv:2303.11944.

THE AIR FORCE RESEARCH LABORATORY

Approved for public release: distribution unlimited. Case Number AFRL-2023-1547

17

# Brain-Inspired AI

## System 1

- **Reinforcement-based mechanisms**

- Value of stimuli and actions are learned incrementally and through repeated experience

- Extracts statistical co-occurrences among features

- Neocortex

  Slow acquisition of structure

  Parametric

  Efficient representations for generalization

## System 2

- **Conscious memory**

- Instance based mechanisms

- Allow experiences to be encoded rapidly (in ''one shot'')

- Hippocampus

  - Rapid storage: individual experiences

  - Non-parametric instance-based system

  - Sparse non-overlapping representations (poor generalization)

Schmidt, K., Culbertson, J., Cox, C., Clouse, H., Larue, O., Molineaux, M., Rogers, S. (2021). What is it Like to Be a Bot: Simulated, Situated, Structurally Coherent Qualia (S3Q) Theory of Consciousness. https://arxiv.org/abs/2103.12638
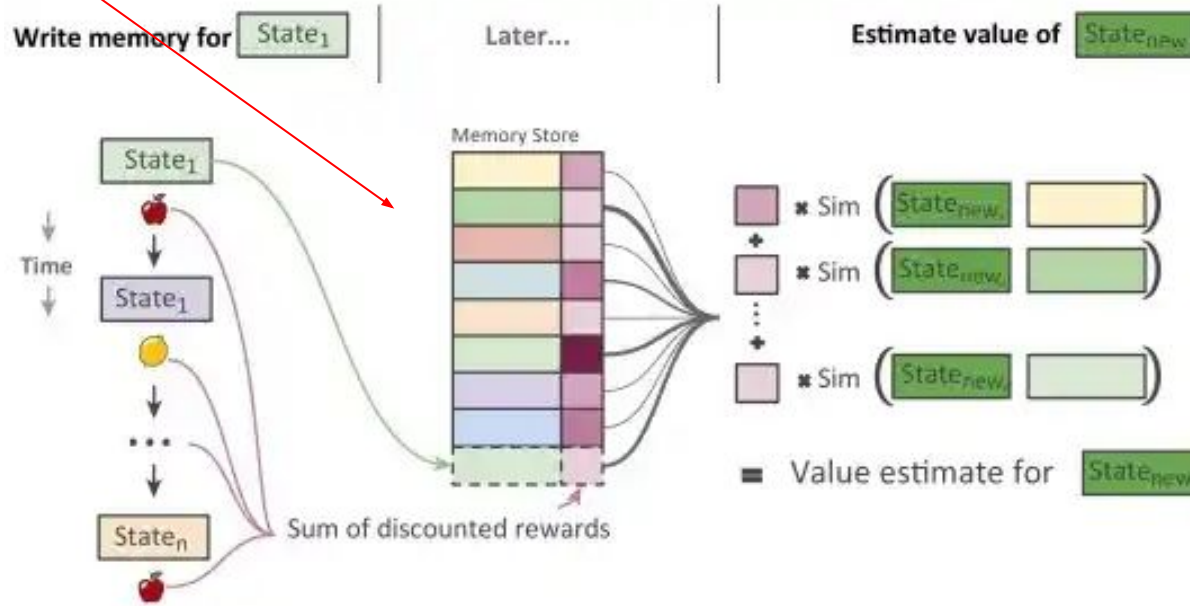
Han, C., Schmidt, K., Grandoit, E., Shu, P., McRobert, C., Reber, P. (2022). 'Cognitive Neuroscience of Implicit Learning: Implications for Complex Learning and Expertise, in Arthur S. Reber, and Rhianon Allen (eds), *The Cognitive Unconscious: The First Half Century*, New York, online edn, Oxford Academic.

THE AIR FORCE RESEARCH LABORATORY

Approved for public release: distribution unlimited. Case Number AFRL-2023-1547

18

**We need a suitable representation language**
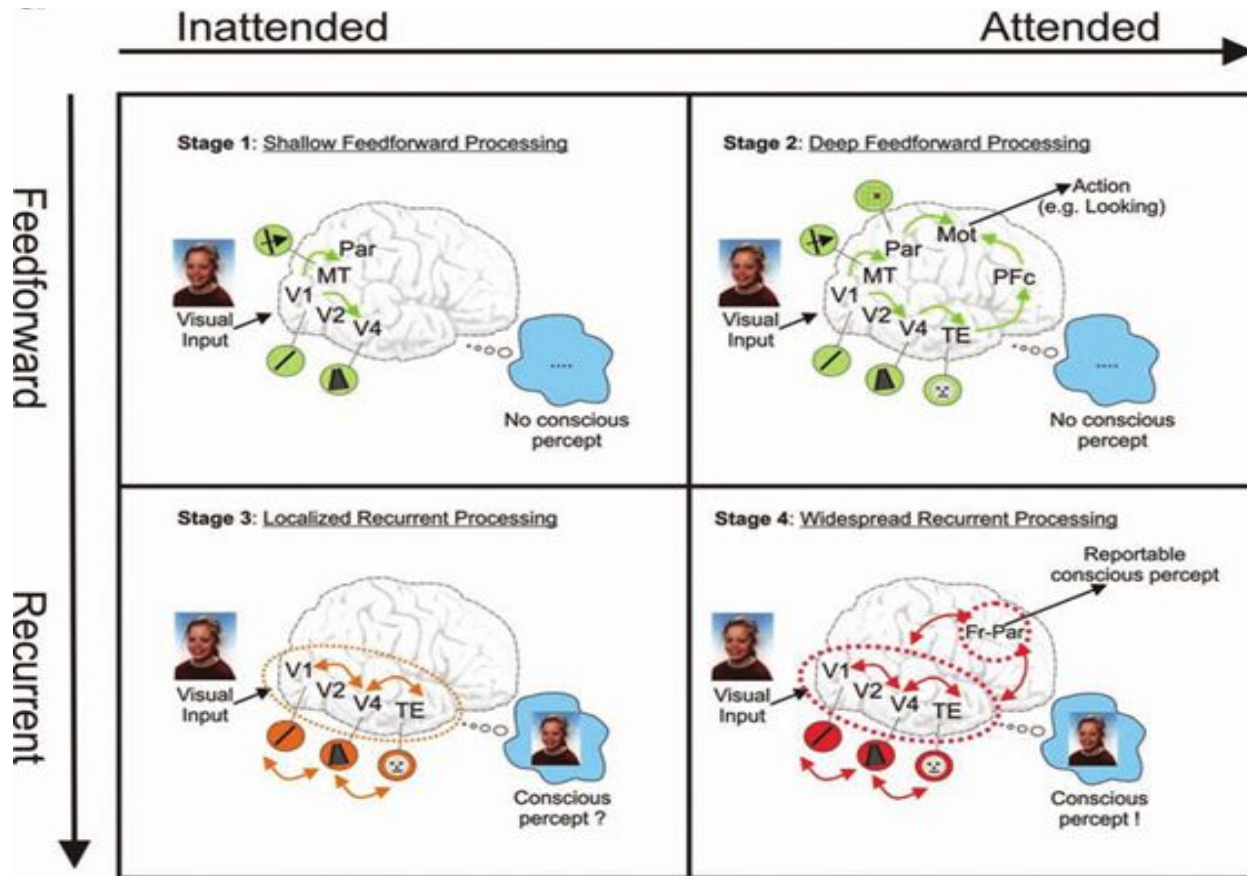
Review

# Reinforcement Learning, Fast and Slow

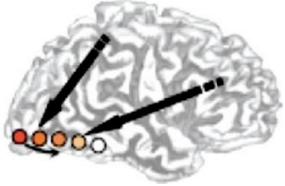Matthew Botvinick,[1,2,*] Sam Ritter,[1,3] Jane X. Wang,[1] Zeb Kurth-Nelson,[1,2] Charles Blundell,[1] and Demis Hassabis[1,2]

Feed forward universal function approximator cannot compute consciousness
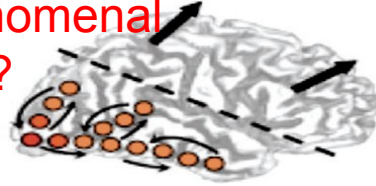


Attention is NOT all you need

Artificial Neural Network

Hierarchical LSTM graphical NN? Better suited for phenomenal content?

Transformer, LLM? Well suited for language?

- Subliminal Attended or Unattended

- Little Local Activity

- Feedforward

- No reportability

- **Conscious??**

- Recurrent Activation confined to sensorimotor processing

- No reportability

- Conscious

- Global Ignition / Gamma Synchrony / P3

- Durable Activation

- Reportable

# Questions?

## Join us @
# QuEST

Fridays at 12p eastern!!

Qualia Exploitation of Sensing Technology

Public Seminar for the World's Best on the Topic

# Kevin Schmidt, PhD
## ACT3 Senior Neuroscientist

Autonomy Capability Team (ACT3) | AFRL/RYZA
Sensors Directorate | Air Force Research Laboratory
Wright-Patterson AFB, OH

kevin.schmidt.15@us.af.mil

meet.google.com/dui-wwjj-fzr