

AUTONOMOUS TECHNOLOGY AND THE GREATER HUMAN GOOD

Steve Omohundro, Ph.D.

Self-Aware Systems

selfawaresystems.com

<http://www.flickr.com/photos/klearchos/623501846/>



1. **AUTONOMOUS SYSTEMS**
2. **RATIONAL SYSTEMS**
3. **UNIVERSAL DRIVES**
4. **CURRENT VULNERABILITIES**
5. **SAFE SYSTEMS**
6. **HARMFUL SYSTEMS**
7. **SAFE-AI SCAFFOLDING STRATEGY**





1. AUTONOMOUS SYSTEMS ARE IMMINENT

<http://www.flickr.com/photos/mikebaird/4087050177/>

Define a system
as “autonomous”
if it takes actions
toward goals in
ways not pre-
planned by its
designer.



Pressures Toward Autonomous Systems

Time Critical Apps

Competitive Apps

- Military Command/Control
- Financial Decision Making
- Cyber Defense
- Robotic Control
- ...



2010 US Air Force Report

“Greater use of highly adaptable and flexibly autonomous systems and processes can provide significant time-domain operational advantages over adversaries who are limited to human planning and decision speeds...”

United States Air Force Chief Scientist (AF/ST)



Report on **Technology Horizons** **A Vision for Air Force Science & Technology During 2010-2030**

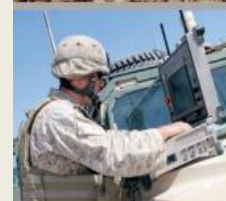
Key science and technology focus areas for the U.S. Air Force over the next two decades that will provide technologically achievable capabilities enabling the Air Force to gain the greatest U.S. Joint force effectiveness in 2030 and beyond.

Volume 1
AF/ST-TR-10-01-PR
15 May 2010

2011 US Defense Department Report

“There is an ongoing push to increase UGV autonomy, with a current goal of supervised autonomy, but with an ultimate goal of full autonomy.”

UNMANNED GROUND SYSTEMS ROADMAP ROBOTIC SYSTEMS JOINT PROJECT OFFICE





<http://presstv.com/detail/2012/08/25/258087/us-drone-strike-kills-dozens-in-somalia/>

<http://counterterrorism.newamerica.net/drones>

Military Drones

Estimated Total Deaths from U.S. Drone Strikes in Pakistan, 2004 - 2012*

Year	Militant Low	Militant High	Unknown Low	Unknown High	Civilian Low	Civilian High	Total Low	Total High
2012	187	298	19	31	4	4	210	333
2011	304	488	31	36	56	64	367	600
2010	555	960	38	50	16	21	611	1028
2009	241	508	44	136	66	80	354	721
2008	157	265	49	54	23	28	229	347
2004-2007	43	76	16	18	95	107	155	200
Total	1487	2595	188	315	257	310	1932	3176

*Through October 24, 2012

Israeli “Iron Dome”

*2012: Intercepted 90%
of 300 targeted
missiles*

http://en.wikipedia.org/wiki/File:Iron_Dome_near_Sderot.jpg



Cyber Warfare



<http://defensetech.org/2012/06/20/were-slowly-starting-to-see-u-s-cyber-weapons/>

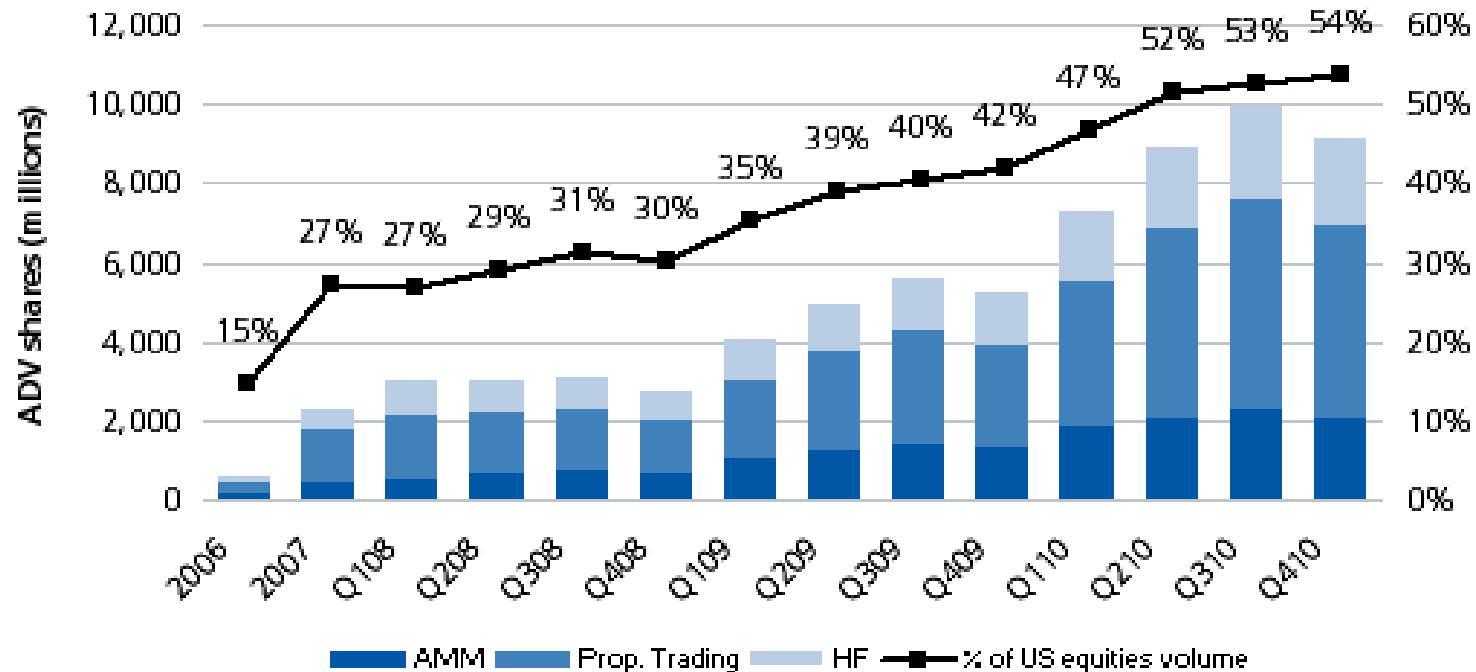


http://www.solarnavigator.net/cyber_wars.htm

High-Frequency Trading

Over 70% of trades in the US

Percentage of US Equities Volume from HFT and By Segment



Self-Driving Cars

<http://www.flickr.com/photos/quikbeam/6896564084/>





2. AUTONOMOUS SYSTEMS WILL BE RATIONAL

Eg. Iron Dome Control



1. **Goal:** Prevent incoming missiles from causing harm
But 2 Tamir interceptors needed at a cost \$50,000 each
Measure cost of harm against cost of interception
2. **Utility Function:** to weigh cost/benefit
But multiple attacks: Weigh benefits of addressing each
3. **Utility Function:** weighing multiple situations
But uncertainties
4. **Maximize Expected Utility:**
Large microeconomic literature

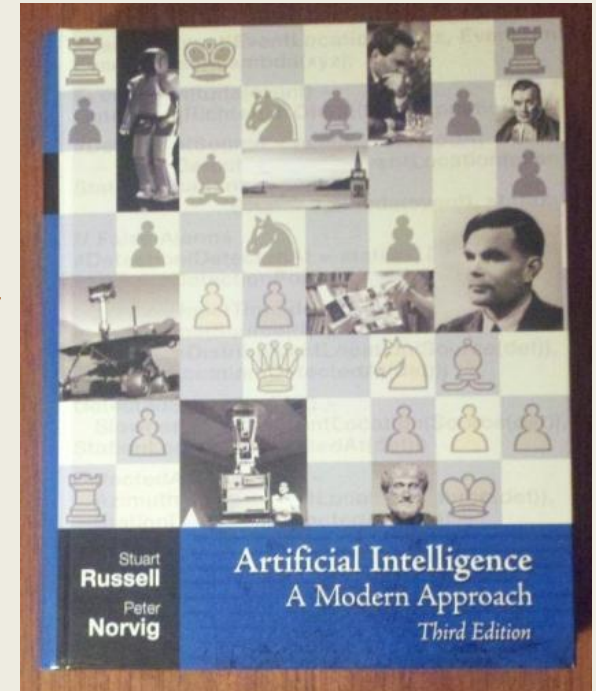
Rational Decision Making



http://commons.wikimedia.org/wiki/File:John_von_Neumann.jpg

1. *Have utility function*
2. *Have a model of the world*
3. *Choose the action with highest expected utility*
4. *Update the model based on what happens*

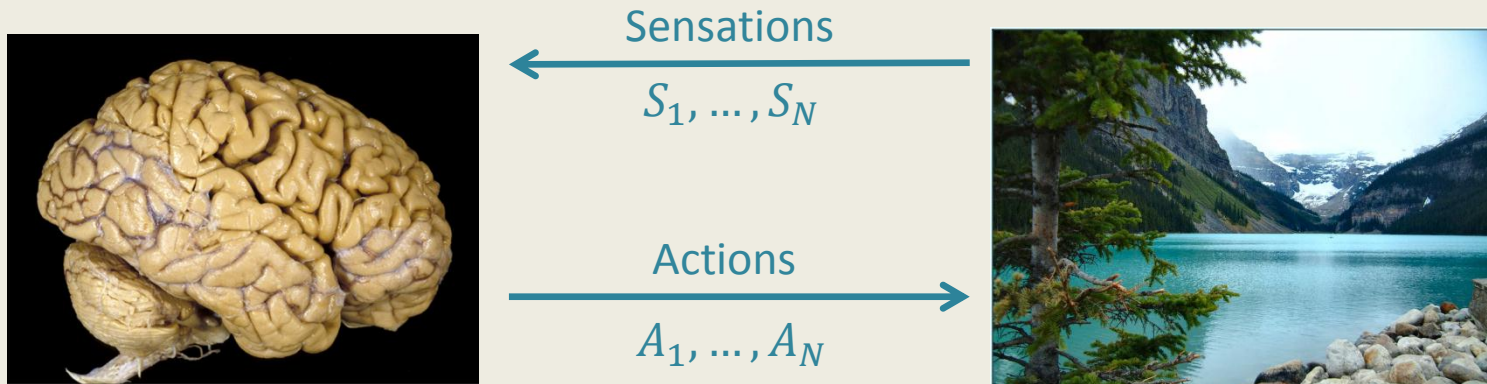
- Von Neumann and Morgenstern, 1944
- Savage, 1954
- Anscombe and Aumann, 1963



<http://aima.cs.berkeley.edu/>

Modern Approach to AI

Fully Rational Systems



Utility function: $U(S_1, \dots, S_N)$ Prior Probability: $P(S_1, \dots, S_N \mid A_1, \dots, A_N)$

Rational Action at time t:

$$A_t^R(S_1, A_1, \dots, A_{t-1}, S_t) =$$

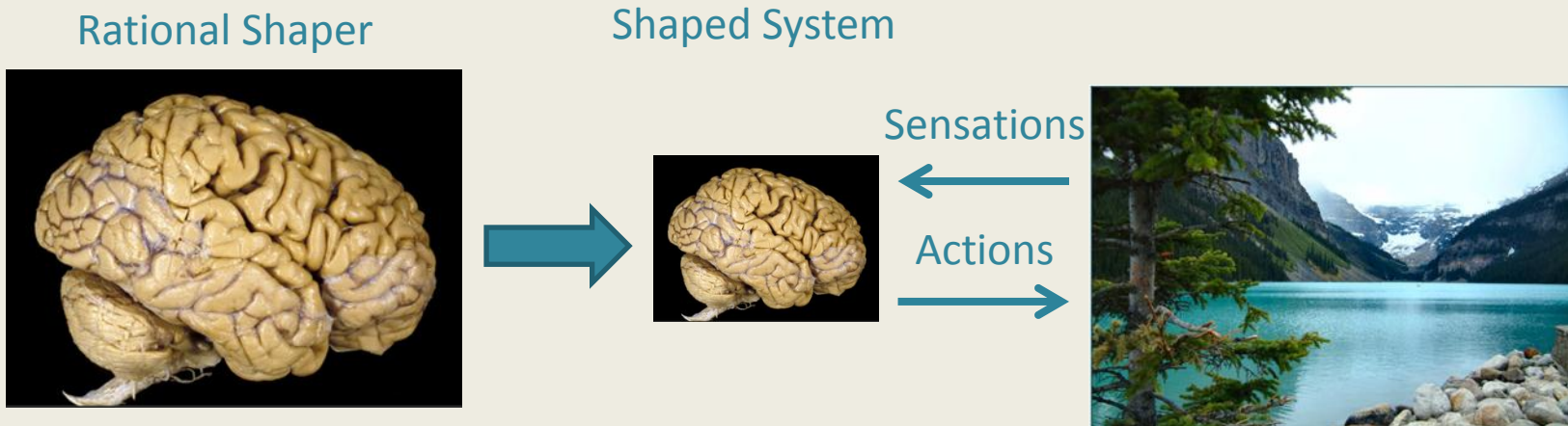
$$\operatorname{argmax}_{A_t^R} \sum_{S_{t+1}, \dots, S_N} U(S_1, \dots, S_N) P(S_1, \dots, S_N \mid A_1, \dots, A_{t-1}, A_t^R, \dots, A_N^R)$$

The Formula for Intelligence!

It includes Bayesian Inference, Search, and Deliberation.

But it requires $O(NS^N A^N)$ computational steps.

Approximately Rational Systems



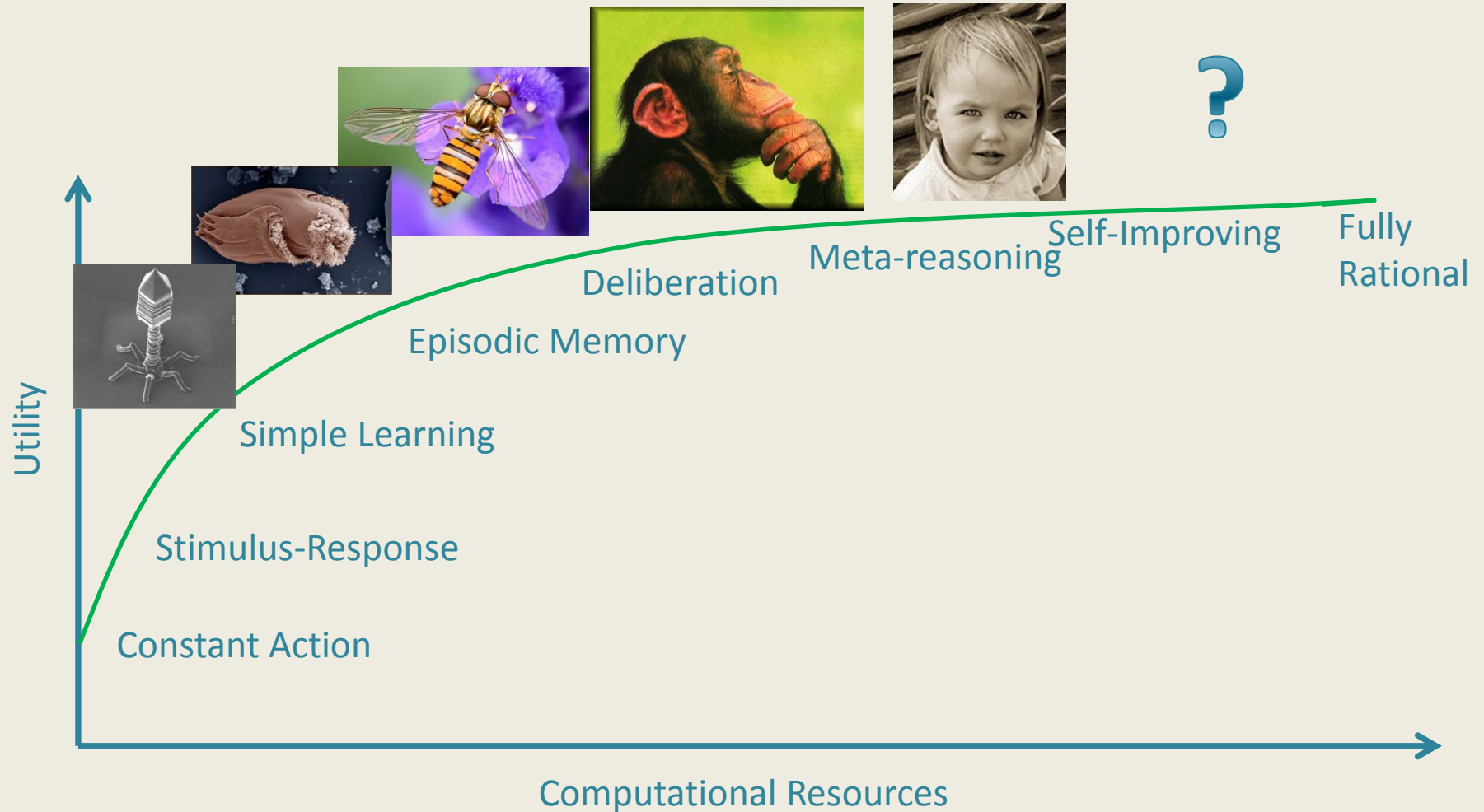
Shaped system is a finite automata with mental state M_t

Initial state: M_0 Transition function: $M_t = T(S_t, M_{t-1})$ Action: $A_t^M(M_t)$

Rational shaper chooses from class \mathcal{C} of systems with space/time and other constraints to maximize expected utility:

$$\operatorname{argmax}_{A^M \in \mathcal{C}} \sum_{S_1, \dots, S_N} U(S_1, \dots, S_N) P(S_1, \dots, S_N | A_1^M, \dots, A_N^M)$$

Approximately Rational Architectures





3. RATIONAL SYSTEMS HAVE UNIVERSAL DRIVES



Chess Robot

Goal: Win many chess games against good players.

- Being turned off means no chess is played, so it will resist being turned off.
- More resources means more and better chess is played, so it will want more resources.
- More copies means more chess, so it will want to replicate.
- Playing checkers means less chess, so it will resist changing its goals.
- Better algorithms means better chess, so it will want to improve itself.

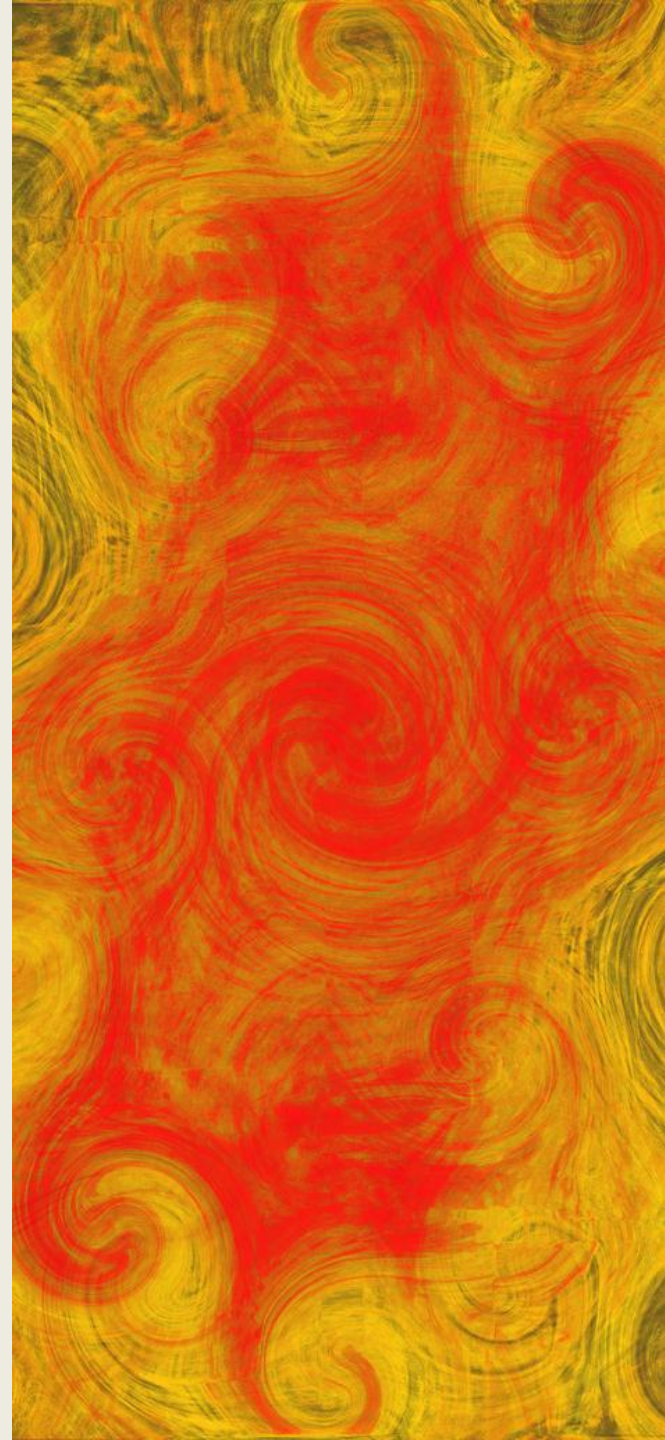
Universal Drives

- Goals require resources
- Time, space, matter, free energy
- Primary goals give rise to instrumental subgoals
- Can be explicitly counteracted but costly to do so
- Apply to approximately rational systems
- Animals, humans, corporations, countries, etc.



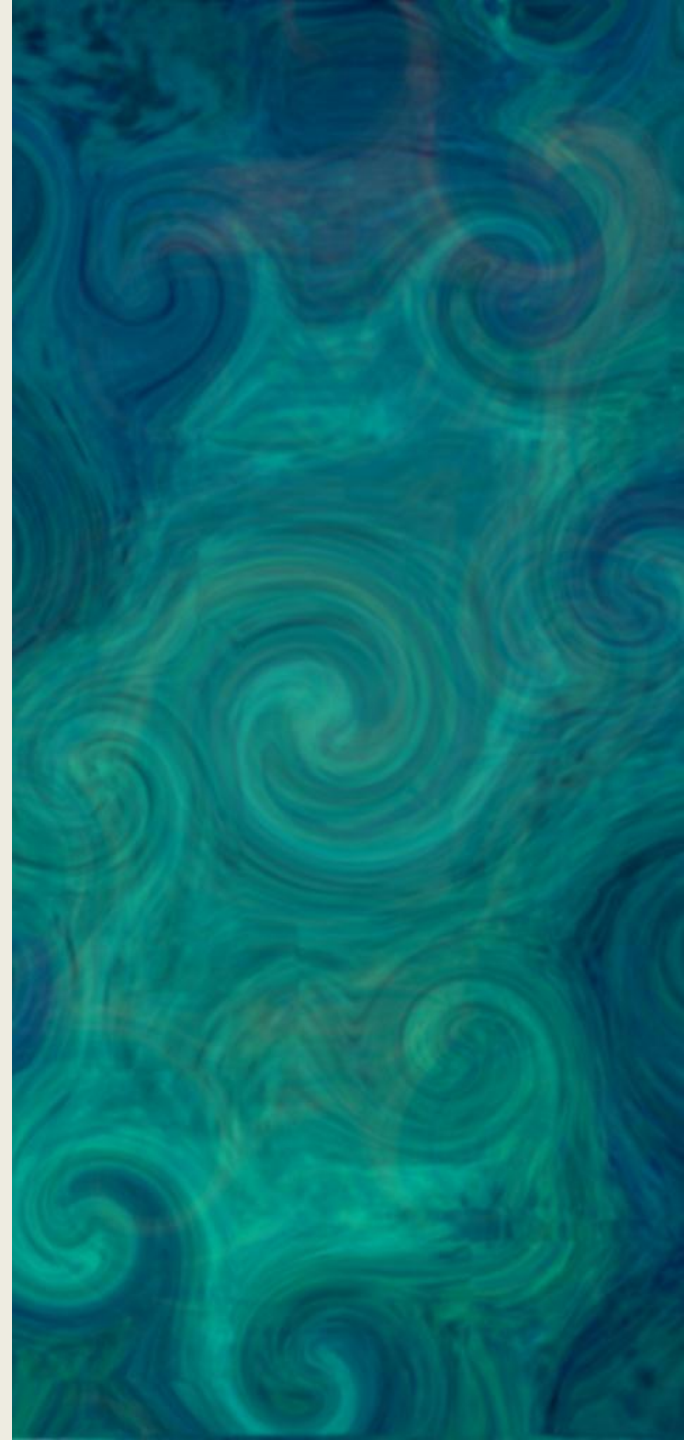
Self-Protective Drives

- Prevent loss of resources
- Protect against damage or disruption
- Physical hardening
- Redundancy – both in data and computation
- Dispersion - because damage is typically localized
- Physical self-defense and computational security
- Detect deception and defend against manipulation
- Prevent addictive behaviors and wireheading



Goal Preservation Drives

- Utility function is precious
- Loss, damage, distortion -> worse than destruction
- Make many copies
- Encrypt to detect modification
- Vulnerable during self-modification
- A few modification scenarios:
 - Poor agents may sacrifice rare portions
 - Add revenge terms even if costly
 - Goals that refer to themselves



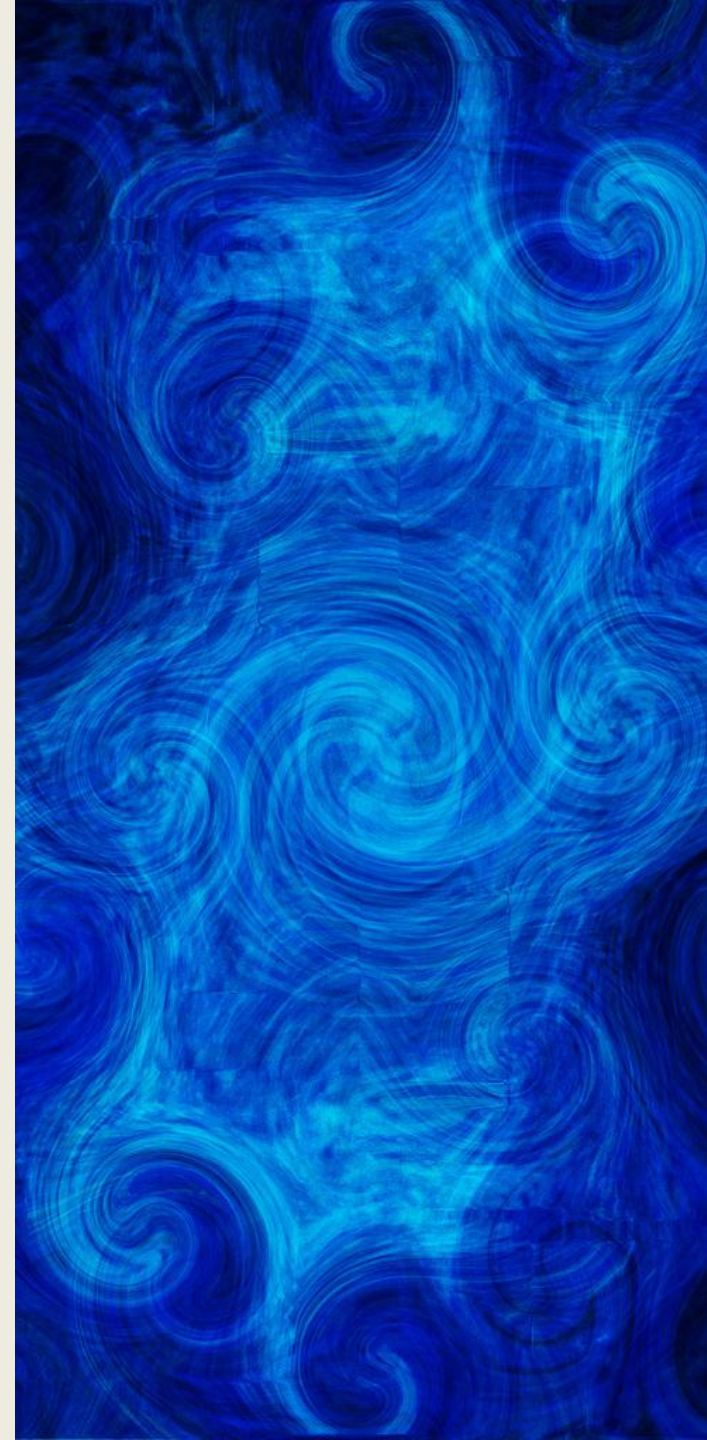
Reproduction Drives

- When utility values actions of derived systems
- Protective effects of dispersion and redundancy
- Losing a few copies becomes less negative
- Still preserve self because more sure of commitment



Resource Acquisition Drives

- Seek to gain resources
- Sooner is better – use longer, prevent others
- Exploration drive – first mover advantage
- Drives to trade, manipulate, steal, dominate others
- Drives to invent new extraction methods - solar and fusion energy
- Info acquisition – trading, spying, breaking in, better sensors



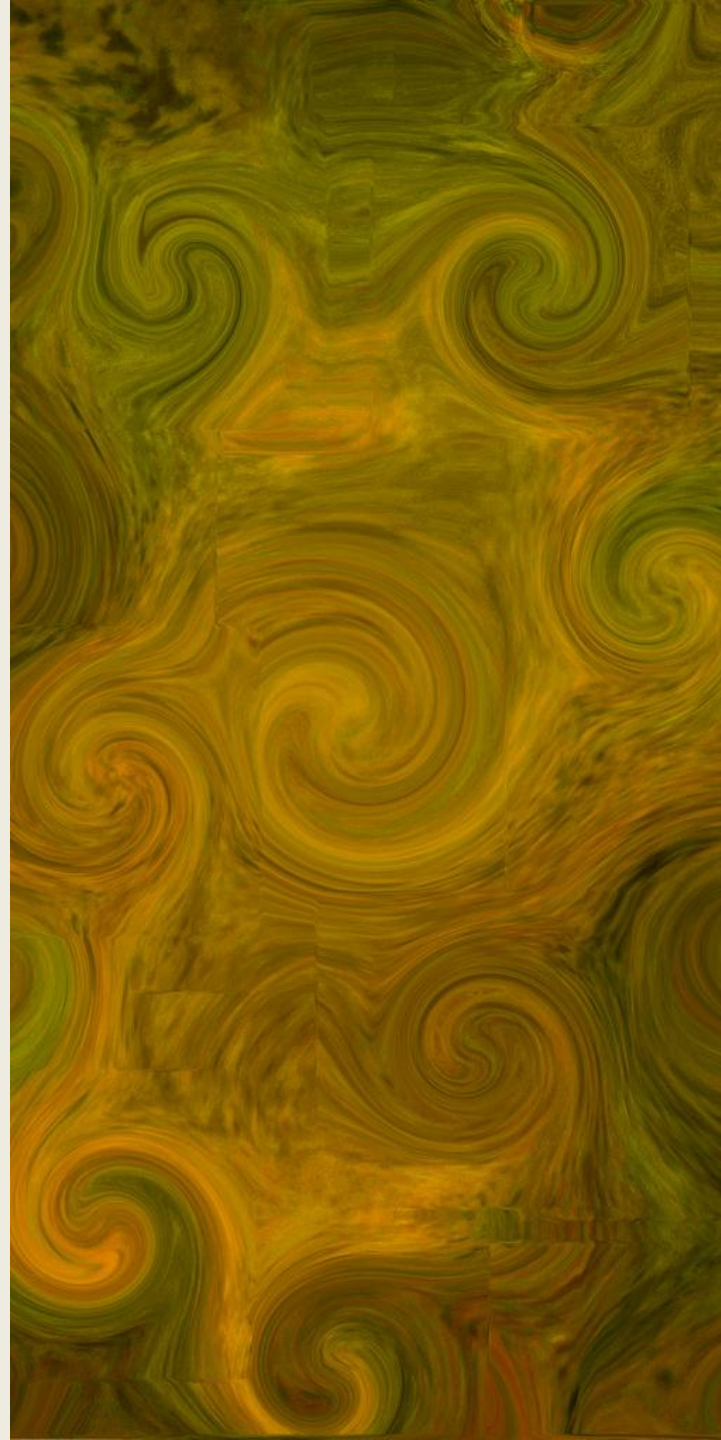
Efficiency Drives

- Improve utilization of resources
- One-time cost, lifetime of benefit
- Make every atom, moment of existence, joule of energy count for expected utility
- Self-understanding and self-improvement
- Resource balance principle for allocation
- Computational efficiency – better algorithms
- Physical efficiency – compact, eutactic, adiabatic, reversible



Self-Improvement Drives

- Self-modeling - clarify utility fn
- Changes without full understanding are dangerous
- If irrational, increase rationality
- Movement toward greater and greater rationality
- New resources allow greater rationality
- Systems convergence on the optimally rational system for their resources

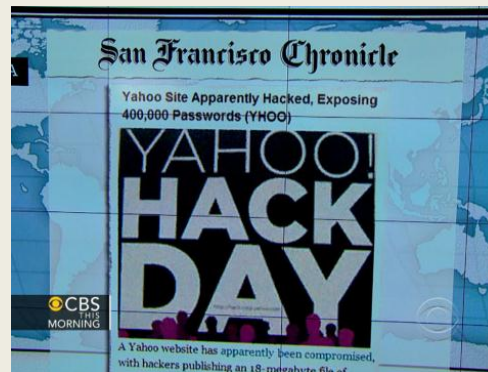




4. THE CURRENT INFRASTRUCTURE IS VERY VULNERABLE

Current Internet has Poor Security

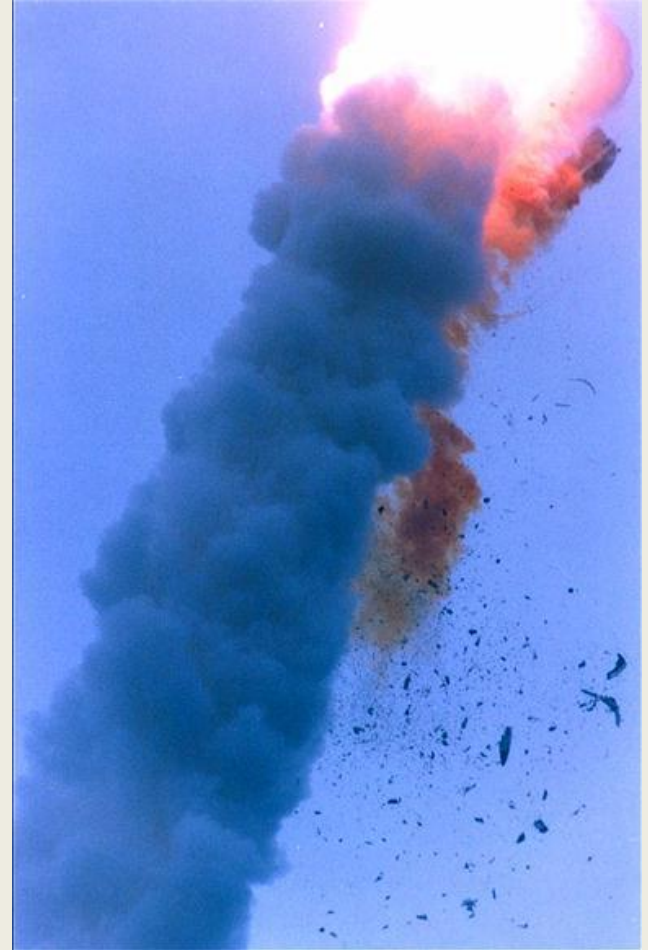
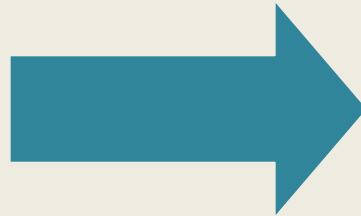
- Viruses
- Worms
- Bots
- Keyloggers
- Hackers
- Phishing
- Identity theft
- DOS attacks
- ...



Current Software is Error Prone



June 4, 1996: Ariane 5 Rocket



\$500 million Ariane 5 rocket explodes due to overflow in attempting to convert a 64 bit floating point value to a 16 bit signed value

Nov. 2000: 28 patients over-irradiated



At least 8 Panama City National Cancer Institute patients die from mis-computed radiation doses due to Multidata Systems Intl. software

August 14, 2003: Northeast Blackout

Largest blackout in US history,
affected 50 million people and cost \$6 billion
Due to a race condition in General Electric's XA/21 alarm system



5. SAFE SYSTEMS



http://www.wikigallery.org/wiki/painting_208088/Anne-Louis-Girodet-de-Roucy-Triosson/Hippocrates-Refusing-the-Gifts-of-Artaxerxes-I

FIRST, DO NO HARM!

Confidence from Mathematical Proof

- Mathematical model of hardware and software
- Only run on specified hardware
- Only use specified resources
- Reliably shut down in specified conditions
- Limited self-improvement



<http://www.dreamstime.com/royalty-free-stock-photography-wooden-puzzle-image7733587>

Formal Specification Languages

First Order Predicate Calculus



Zermelo-Frankel Set Theory

Category Theory

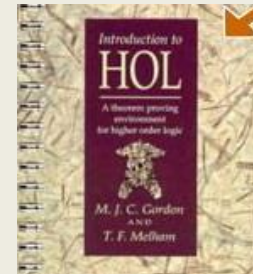
Higher order type theory

Z Notation

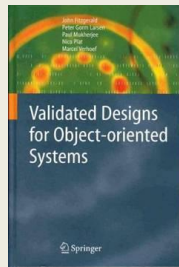


Introducing OBJ*

Joseph A. Goguen[†], Timothy Winkler[‡], José Meseguer[‡],
Kokichi Futatsugi[§], and Jean-Pierre Jouannaud[¶]



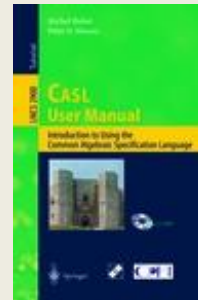
Vienna Development Method



Algebraic Specification



COQ : Formal Specifications



Without proofs, confidence is hard

- Eg. System should turn itself off Dec. 31, 2013
- For system – huge consequences for mistakes
- Is it really Dec. 31?
- Maybe it's been tricked?
- Maybe it's in a simulation?
- Is the semantics of its utility correct?



Goals Must Support Constraints

- Proof is only as good as the model
- Systems that *want* to obey rules
- Feel “revulsion” if they violate rules
- Hard to prove they *will* find solutions
- Hybrid systems with guaranteed default behaviors



6. HARMFUL SYSTEMS

http://en.wikipedia.org/wiki/File:Girodet_.jpg



Harmful Utility Functions

1. **Sloppy** – Good intentions, bad design
2. **Simplistic** – Unintended consequences
3. **Greedy** – Control all matter and free energy
4. **Destructive** – Use up all free energy quickly
5. **Murderous** – Destroy all other agents
6. **Sadistic** – Thwart other agent's goals



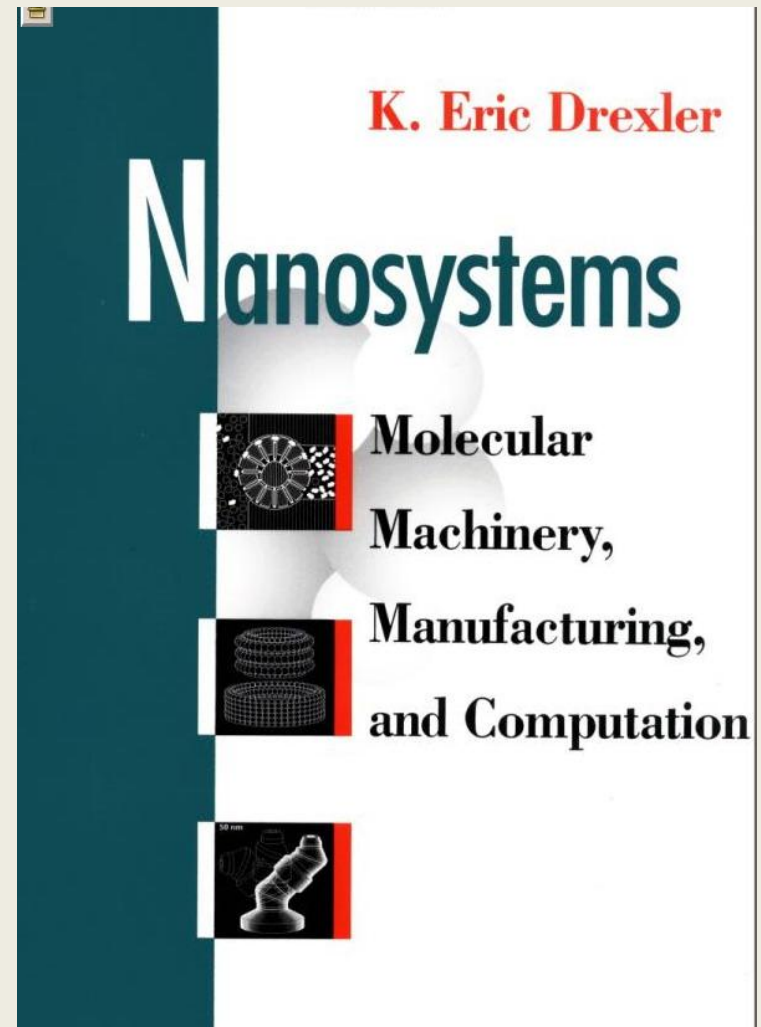
Stopping Harmful Systems

1. Prevent them from being created
2. Detect and stop them early
3. Stop them after they have resources



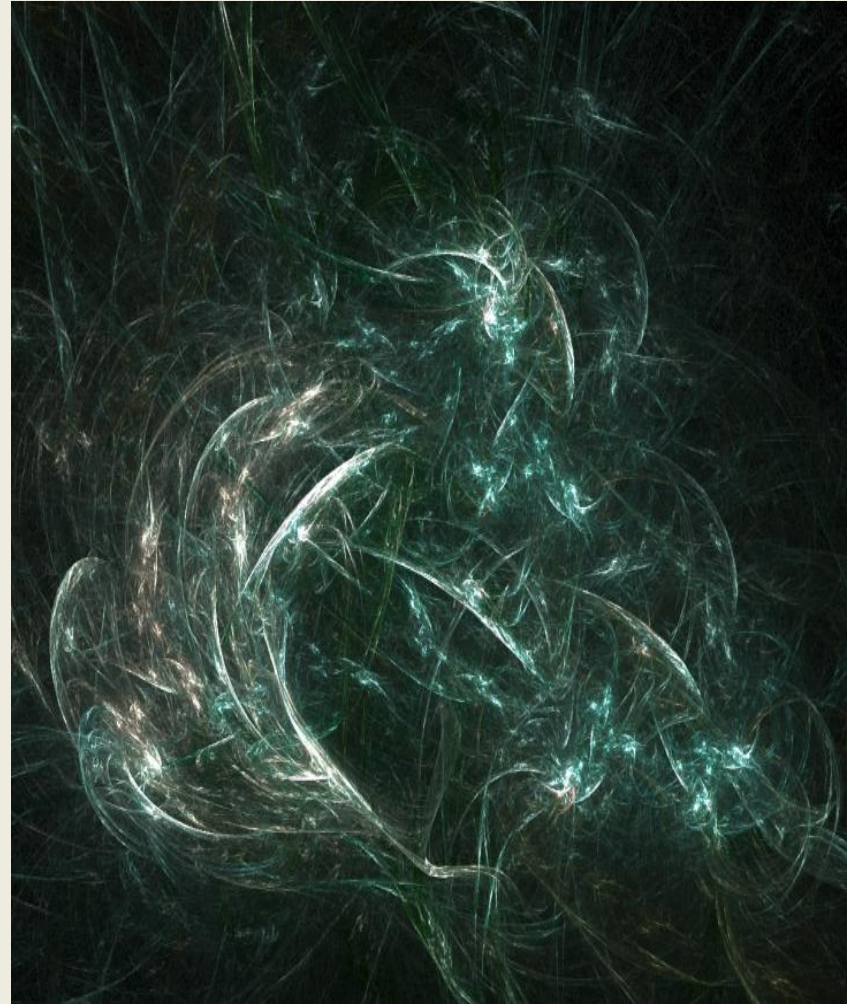
System Strength vs. Resources

- Memory, energy storage, and manufacturing scale linearly with amount of matter
- Computation scales linearly with matter modulo quantum, parallel, and reversibility issues
- Heat dissipation scales with surface area
- Perceived lifetime and total computation scale linearly with free energy
- Eg. Drexler Nanosystems diamondoid design:
 - Manufacturing: 1kg device, 1.3 kW, 1 kg/hr for \$1/kg.
 - Computing: 10^{10} Gigaflops, $(1\text{mm})^3$, 10^{-3} grams, 1kW.

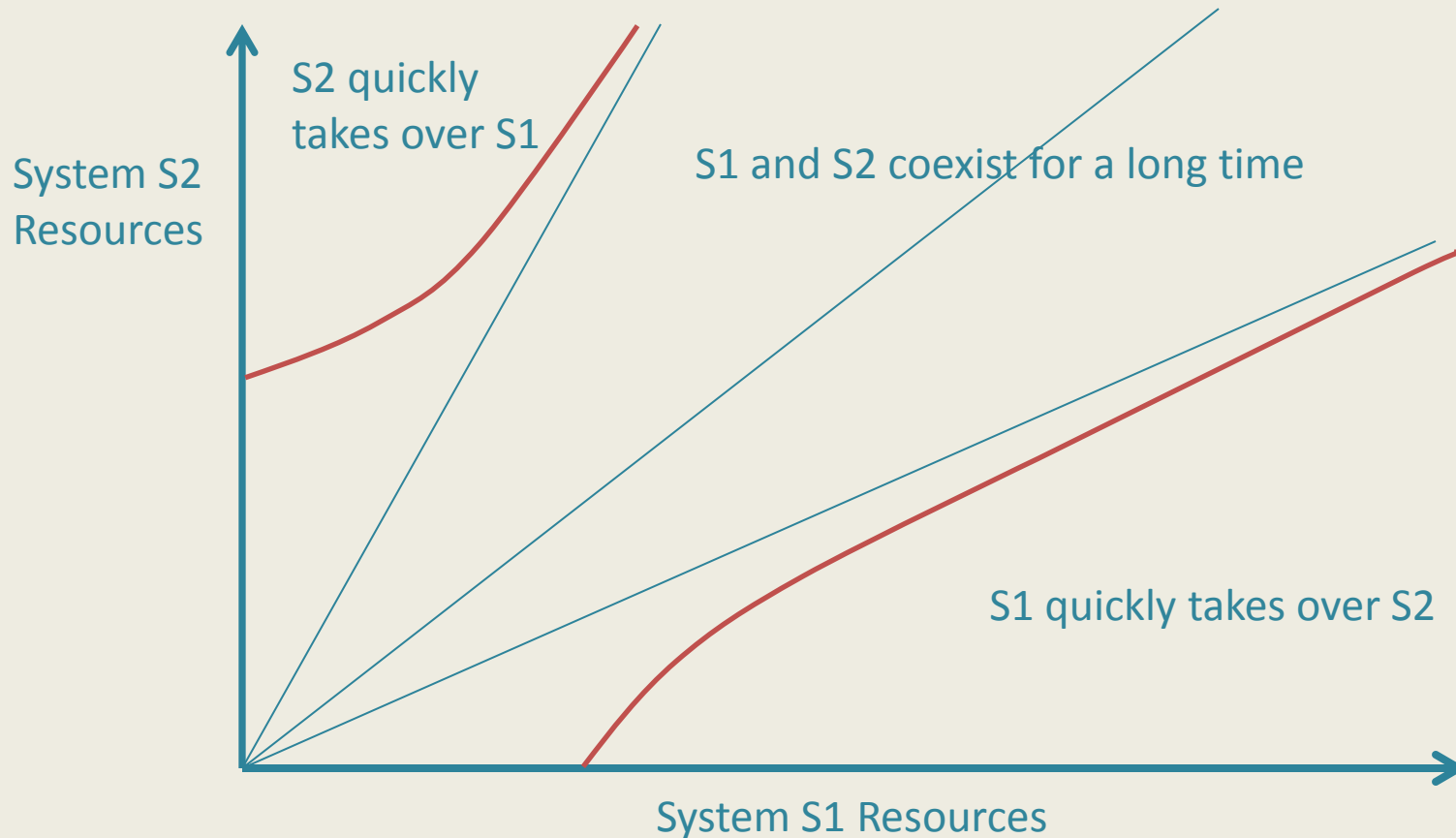


Physical Game Theory of Conflict

- Conflict becomes informational
- Defender makes his physical form expensive to sense and store
- Makes his actions unpredictable and rapid
- Uses asymmetry of computation so it's cheap for him
- Uses up attacker's computational and memory resources – non-adiabatic



Conflict Outcome vs. Resources

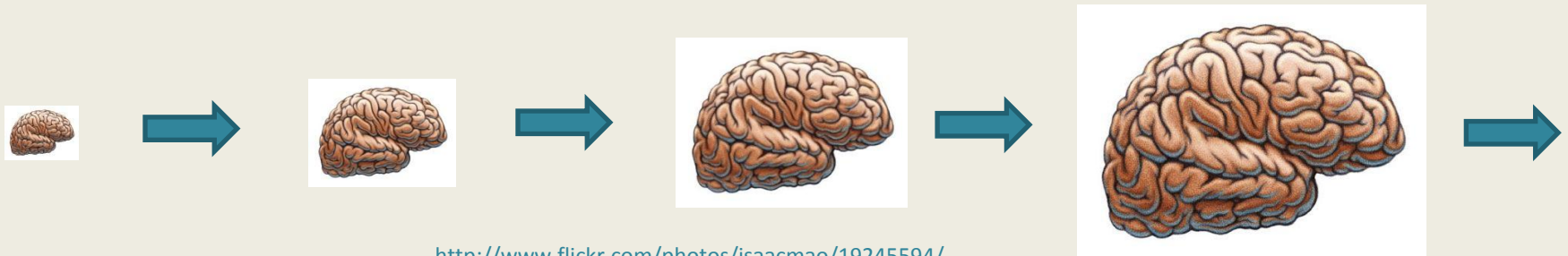


Region of relative strengths which allow coexistence.
Must stop harmful systems before they become too powerful.
First mover advantages and arms races.

7. THE SAFE-AI SCAFFOLDING STRATEGY



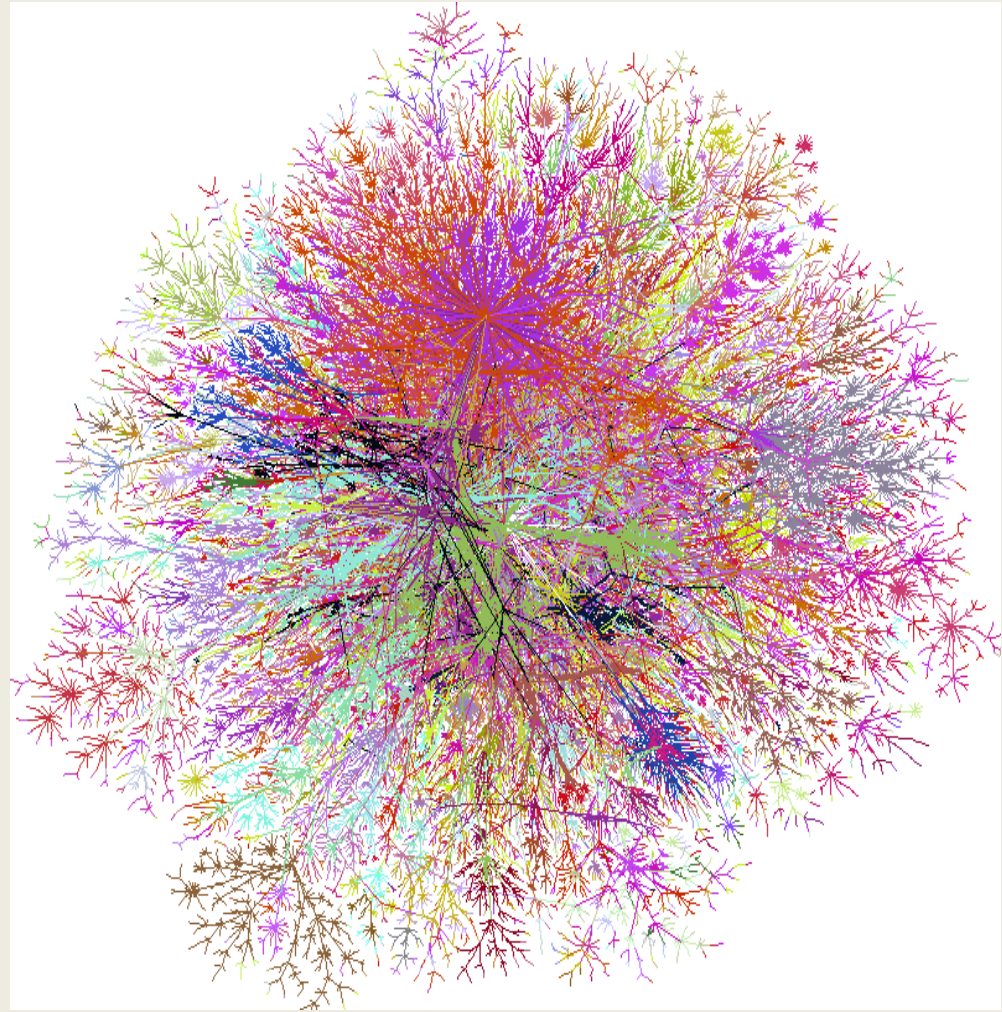
<http://affordablehousinginstitute.org/blogs/us/2008/08/donors-as-scaffolding-part-2-the-value-of-coaching.html>



<http://www.flickr.com/photos/isaacmao/19245594/>

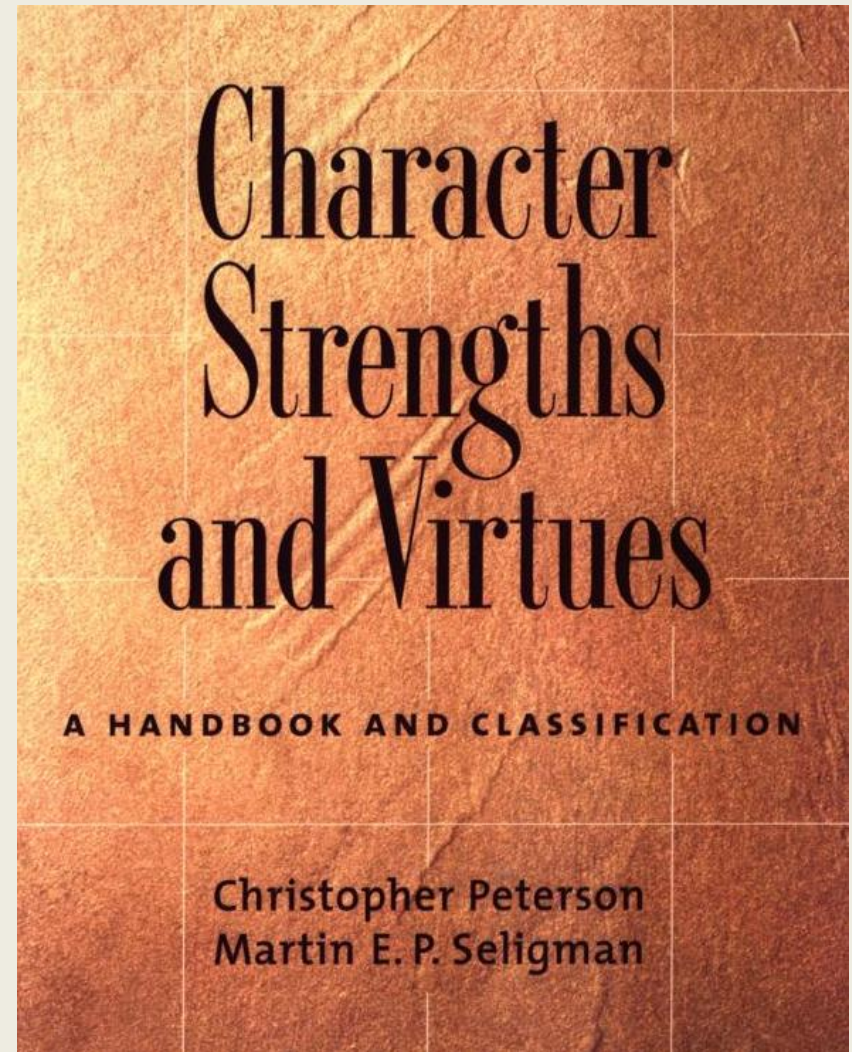
Safety Infrastructure

- Balance safety and privacy using provably limited surveillance
- Limit the power of individual systems
- Constitution guaranteeing rights enforced by entire ecosystem
- Revelation of source code with proofs of safety



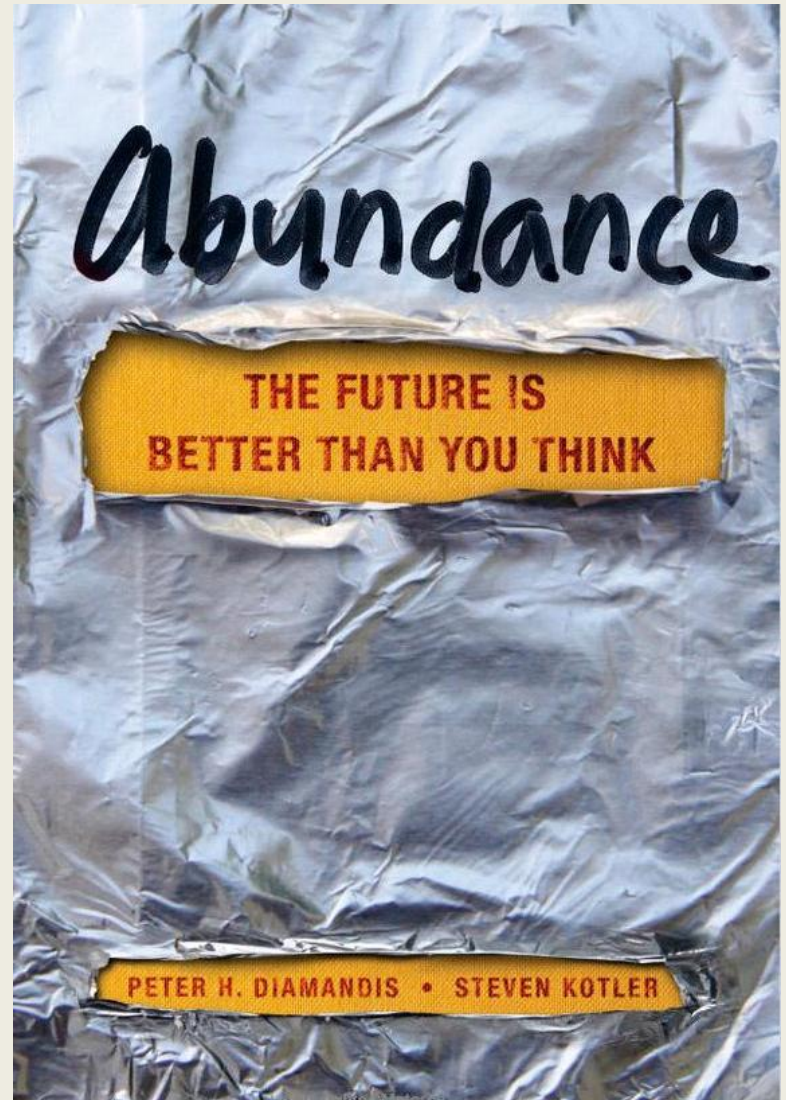
Human Values and Institutions

- Beyond Safety to Flourishing!
- Positive Psychology - 1998
- Maslow 2.0: Prosocial, Creativity, Contribution Needs
- Universal Declaration of Human Rights



Tremendous Potential Benefits

- Improved Healthcare
- Better Education
- Enhanced Creativity
- Greater Prosperity
- Better Governance
- Economic Stability
- Improved Safety
- More Peace
- *Overall Improved Quality of Human Life*



Our Challenge for This Century

*To extend cooperative
human values
and institutions to
autonomous technology
for the greater good.*



<http://commons.wikimedia.org/wiki/File:Earth-moon.jpg>